
Robust Feature Learning for Multi-Index Models in High Dimensions

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recently, there have been numerous studies on feature learning with neural net-
2 works, specifically on learning single- and multi-index models where the target is a
3 function of a low-dimensional projection of the input. Prior works have shown that
4 in high dimensions, the majority of the compute and data resources are spent on
5 recovering the low-dimensional projection; once this subspace is recovered, the
6 remainder of the target can be learned independently of the ambient dimension.
7 However, implications of feature learning in adversarial settings remain unexplored.
8 In this work, we take the first steps towards understanding adversarially robust
9 feature learning with neural networks. Specifically, we prove that the hidden di-
10 rections of a multi-index model offer a Bayes optimal low-dimensional projection
11 for robustness against ℓ_2 -bounded adversarial perturbations under the squared loss,
12 assuming that the multi-index coordinates are statistically independent from the rest
13 of the coordinates. Therefore, robust learning can be achieved by first performing
14 standard feature learning, then robustly tuning a linear readout layer on top of the
15 standard representations. In particular, we show that adversarially robust learning
16 is just as easy as standard learning. Specifically, the additional number of samples
17 needed to robustly learn multi-index models when compared to standard learning
18 does not depend on dimensionality.

19 1 Introduction

20 A crucial capability of neural networks is their ability to hierarchically learn useful features, and
21 to avoid the curse of dimensionality by *adapting* to potential low-dimensional structures in data
22 through empirical risk minimization (ERM) [Bac17, SH20]. Recently, a theoretical line of work has
23 demonstrated that gradient-based training, which is not a priori guaranteed to implement ERM due to
24 non-convexity, also demonstrates similar behavior and efficiently learns functions of low-dimensional
25 projections [WLLM19, DLS22, BBSS22, BEG⁺22, BES⁺22, MHPG⁺23] or functions with certain
26 hierarchical properties [AAM22, ABAM23, DKL⁺23]. These theoretical insights provided a useful
27 avenue for explaining standard feature learning mechanisms in neural networks.

28 On the other hand, it has been empirically observed that deep neural networks trained with respect
29 to standard losses are susceptible to adversarial attacks; small perturbations in the input may not be
30 detectable by humans, yet they can significantly alter the prediction performed by the model [SZS⁺14].
31 To overcome this issue, a popular approach is to instead minimize the adversarially robust empirical
32 risk [MMS⁺18]. However, unlike its standard counterpart, achieving successful generalization of
33 deep neural networks on robust test risk has been particularly challenging, and even the standard
34 performance of the model can degrade once adversarial training is performed [TSE⁺18]. Therefore,
35 one may wonder if robust neural networks are still adaptive to certain problem structures that improve

36 standard generalization. By focusing on hidden low-dimensionality as a well-known example of such
37 structure, we aim at answering the following fundamental question:

38 *Can neural networks maintain their statistical adaptivity to low dimensions*
39 *when trained to be robust against adversarial perturbations?*

40 We answer this question positively by providing the following contributions.

- 41 • When considering ℓ_2 -constrained perturbations, Bayes optimal predictors can be constructed by
42 projecting the input data onto the low-dimensional subspace defined by the target function. In this
43 sense, the optimal low-dimensional projection remains unchanged compared to standard learning.
- 44 • Consequently, provided that they have access to an oracle that is able to recover the low-
45 dimensional target subspace, neural networks can achieve a sample complexity that is *inde-*
46 *pendent of the ambient dimension* when robustly learning multi-index models. This is achieved
47 by minimizing the empirical adversarial risk with respect to the second layer.
- 48 • An oracle for recovering the low-dimensional target subspace can be constructed by training
49 the first layer of a two-layer neural network with a standard loss function, as demonstrated by
50 many prior works. By combining our results with two particular choices of oracle implementa-
51 tion [DLS22, LOSW24], we provide end-to-end guarantees for robustly learning multi-index
52 models with gradient-based algorithms.

53 1.1 Related Works

54 **Feature Learning for Single/Multi-Index Models.** Many recent works have focused on proving
55 benefits of feature learning, allowing the neural network weights to travel far from initialization,
56 as opposed to the fixed kernel regime of freezing weights around initialization [JGH18, COB19].
57 When using online SGD on the squared loss, [BAGJ21] showed that the complexity of learning
58 single-index models with known link function depends on a quantity called information exponent.
59 Gradient-based learning of single-index models has been studied in [BES⁺22, MHPG⁺23, BBSS22]
60 among others. [DLS22] considered multi-index polynomials where the equivalent of information
61 exponent is at most 2. The counterpart of information exponent for multi-index models, the leap
62 exponent, was introduced in [ABAM23]. Considering SGD on the squared loss as an example of
63 a Correlational Statistical Query (CSQ) algorithm, [DNGL23] provided CSQ-optimal algorithms
64 for learning single-index models. Further improvements to the isotropic sample complexity were
65 achieved by either considering structured anisotropic Gaussian data [BES⁺23, MHPG⁺23], or the
66 sparsity of the hidden direction [VE24].

67 More recently, it was observed that gradient-based learning can go beyond CSQ algorithms by reusing
68 batches [DTA⁺24, LOSW24, ADK⁺24], or by changing the loss function [JMS24]. In such cases,
69 the algorithm becomes an instance of a Statistical Query (SQ) learner, and the sample complexity is
70 characterized by the generative exponent of the link function [DPVLB24].

71 While the above works exist in a narrow-width setting where the interaction between neurons is
72 ignored, another line of research focused on the mean-field or wide limits of two-layer neural net-
73 works [CB18, RVE18, MMN18] for providing learnability guarantees [WLLM19, CB20, AAM22,
74 Tel23, MZD⁺23, CG24]. In particular, the mean-field Langevin algorithm provides global conver-
75 gence guarantees for two-layer NNs [Chi22, NWS22], leading to sample complexity linear in an effec-
76 tive dimension for multi-index models [SWON23, NOSW24] and multi-index models [MHWE24].

77 **Adversarially Robust Learning.** The existence of small worst-case or adversarial perturba-
78 tions that can significantly change the prediction of deep neural networks was first demonstrated
79 in [SZS⁺14]. Among many defences proposed, one effective approach is adversarial training intro-
80 duced by [MMS⁺18], which is based on solving a min-max problem to perform robust optimization.
81 However, adversarial training tends to decrease the standard performance [TSE⁺18]. Therefore, the
82 following works studied the hardness of robust learning and established a statistical separation in a
83 simple mixture of Gaussians setting [SST⁺18], or computational separation by proving statistical
84 query lower bounds [BLPR19]. Further studies focused on exact characterizations of the robust and
85 standard error, as well as the fundamental and the algorithmic tradeoffs between robustness and
86 accuracy in the context of linear regression [JSH20], mixture of Gaussians classification [JS22], and
87 in the random features model [HJ24]. Closer to our work, [JM24] show that this tradeoff is mitigated
88 when the data enjoy a low-dimensional structure. However, the focus there is on binary classification

89 and generalized linear models, where the features live on a low-dimensional manifold. Here, we
 90 consider a multi-index model wherein the response depends on a low-dimensional projection of
 91 features. In addition, in [JM24] it is assumed that the manifold structure is known and the focus is on
 92 population adversarial risk (assuming infinite samples with fixed dimension), while here we consider
 93 algorithms for representation learning, and derive rates of convergence for adversarial risk.

94 In this work, we provide an alternative narrative compared to the line of work above, by showing
 95 that in a high-dimensional regression setting, learning multi-index models that are robust against ℓ_2
 96 perturbations can be as easy as standard learning. We achieve this result by focusing on the feature
 97 learning capability of neural networks, i.e. their ability to capture low-dimensional projections.

98 **Notation.** For Euclidean vectors, $\langle \cdot, \cdot \rangle$ and $\|\cdot\|$ denote the Euclidean inner product and norm respec-
 99 tively. For tensors, $\|\cdot\|_F$ and $\|\cdot\|$ denote the Frobenius and operator norms respectively. We use \mathbb{S}^{k-1}
 100 for the unit sphere in \mathbb{R}^k , and τ_k denotes the uniform probability measure on \mathbb{S}^{k-1} .

101 2 Problem Setup: Statistical Model and Adversarial Robustness

102 **Statistical Model.** Consider a regression setting where the input $\mathbf{x} \in \mathbb{R}^d$ and the target $y \in \mathbb{R}$ are
 103 generated from a distribution $(\mathbf{x}, y) \sim \mathcal{P}$. For a prediction function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, its population
 104 adversarial risk, where we assume the adversary can perform a worst-case perturbation on the input
 105 with a budget of ε measured in ℓ_2 -norm, before passing it to the model, is defined as

$$\text{AR}(f) := \mathbb{E} \left[\max_{\|\boldsymbol{\delta}\| \leq \varepsilon} (f(\mathbf{x} + \boldsymbol{\delta}) - y)^2 \right], \quad (2.1)$$

106 where the expectation is over all random variables inside the brackets. Given a (non-parametric)
 107 family of prediction functions \mathcal{F} , our goal is to learn a predictor that achieves the optimal adversarial
 108 risk given by

$$\text{AR}^* := \min_{f \in \mathcal{F}} \text{AR}(f), \quad (2.2)$$

109 We focus on learners of the form of two-layer neural networks with width N , given as

$$f(\mathbf{x}; \mathbf{a}, \mathbf{W}, \mathbf{b}) = \mathbf{a}^\top \sigma(\mathbf{W}\mathbf{x} + \mathbf{b}), \quad (2.3)$$

110 where $\mathbf{a} \in \mathbb{R}^N$ is the second layer weights and $\mathbf{W} \in \mathbb{R}^{N \times d}$ and $\mathbf{b} \in \mathbb{R}^N$ are the first layer weights
 111 and biases. To avoid overloading the notation we use $\text{AR}(f(\cdot; \mathbf{a}, \mathbf{W}, \mathbf{b})) = \text{AR}(\mathbf{a}, \mathbf{W}, \mathbf{b})$. Given
 112 access to n i.i.d. samples $\{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^n$ from \mathcal{P} , the goal is to learn the network parameters \mathbf{a} , \mathbf{W} ,
 113 and \mathbf{b} in such a way that the quantity $\text{AR}(\mathbf{a}, \mathbf{W}, \mathbf{b})$ is close to the optimal adversarial risk AR^* .

114 A long line of recent works has shown that neural networks are particularly efficient in regression
 115 tasks when the target is a function of a low-dimensional projection of the input, see e.g. [Bac17]. We
 116 also make the same assumption that the data follows a *multi-index model*,

$$\mathbb{E}[y | \mathbf{x}] = g(\langle \mathbf{u}_1, \mathbf{x} \rangle, \dots, \langle \mathbf{u}_k, \mathbf{x} \rangle), \quad (2.4)$$

117 for all $\mathbf{x} \in \mathbb{R}^d$, where $g : \mathbb{R}^k \rightarrow \mathbb{R}$ is the link function, and we assume $\mathbf{u}_1, \dots, \mathbf{u}_k$ are orthonormal
 118 without loss of generality. Let $\mathbf{U} \in \mathbb{R}^{k \times d}$ be an orthonormal matrix whose rows are given by (\mathbf{u}_i) ;
 119 we use the shorthand notation $g(\langle \mathbf{u}_1, \mathbf{x} \rangle, \dots, \langle \mathbf{u}_k, \mathbf{x} \rangle) := g(\mathbf{U}\mathbf{x})$. In particular, the above means
 120 that y is independent from the rest of the coordinates when conditioned on $\mathbf{U}\mathbf{x}$. In this paper, we
 121 consider the setting where $k \ll d$, and in particular $k = \mathcal{O}(1)$.

122 3 Optimal Representations for Robust Learning

123 In this section, we demonstrate that under ℓ_2 -constrained perturbations, the optimal low-dimensional
 124 representations for robust learning coincides with those in standard setting, both of which are given by
 125 the target directions \mathbf{U} . Crucially, our result relies on the following assumption on input distribution.

126 **Assumption 1.** Suppose $\tilde{\mathbf{U}} \in \mathbb{R}^{(d-k) \times d}$ is any orthonormal matrix whose rows complete the rows of
 127 \mathbf{U} into a basis of \mathbb{R}^d . Then, $\mathbf{U}\mathbf{x}$ and $\tilde{\mathbf{U}}\mathbf{x}$ are statistically independent.

128 Introducing the notation $\mathbf{x}_\parallel := \mathbf{U}\mathbf{x}$ and $\mathbf{x}_\perp := \tilde{\mathbf{U}}\mathbf{x}$ for any $\mathbf{x} \in \mathbb{R}^d$, the above assumption states that
 129 the distribution of the input \mathbf{x} is such that \mathbf{x}_\parallel , the coordinates that enter the statistical model, are

130 independent from \mathbf{x}_\perp , the coordinates that do not. For example, Assumption 1 holds when \mathbf{x} is a
 131 Gaussian random vector with isotropic covariance, or more generally $\mathbf{x} = \mathbf{U}^\top \mathbf{U} \mathbf{z}_1 + \tilde{\mathbf{U}}^\top \tilde{\mathbf{U}} \mathbf{z}_2$ for
 132 independent vectors $\mathbf{z}_1, \mathbf{z}_2 \in \mathbb{R}^d$. We present a central result below along with its proof.

133 **Theorem 1.** Suppose Assumption 1 holds and (2.2) admits a minimizer. Then, there exists a function
 134 $f^* : \mathbb{R}^d \rightarrow \mathbb{R}$ of the form $f^*(\mathbf{x}) = h(\mathbf{U}\mathbf{x})$ for some $h : \mathbb{R}^k \rightarrow \mathbb{R}$ such that

$$\text{AR}(f^*) \leq \text{AR}^*, \quad (3.1)$$

135 with equality when $f^* \in \mathcal{F}$. Further, h is represented as $h(\mathbf{z}) = \mathbb{E}[f(\mathbf{x}) \mid \mathbf{U}\mathbf{x} = \mathbf{z}]$ for some $f \in \mathcal{F}$.

Remark. To understand the significance of the above result, define the function class $\mathcal{H} = \{\mathbf{z} \mapsto \mathbb{E}[f(\mathbf{x}) \mid \mathbf{U}\mathbf{x} = \mathbf{z}] \text{ for } f \in \mathcal{F}\}$, and observe that the last statement of the theorem reads

$$\min_{h \in \mathcal{H}} \text{AR}(h(\mathbf{U}\cdot)) \leq \text{AR}^*.$$

136 Thus, to achieve the optimal adversarial risk AR^* , one only needs to (i) learn the target directions \mathbf{U} ,
 137 and (ii) approximate real-valued functions in a k -dimensional subspace rather than d . In the context
 138 of NNs, the first layer \mathbf{W} recovers \mathbf{U} , and the remaining parameters \mathbf{a} and \mathbf{b} are used to approximate
 139 the optimal h . While this recipe is general, we provide specific implications in the next section.

140 **Proof.** We will show that for every $f \in \mathcal{F}$, $h(\mathbf{z}) = \mathbb{E}[f(\mathbf{x}) \mid \mathbf{U}\mathbf{x} = \mathbf{z}]$ gives $\text{AR}(h(\mathbf{U}\cdot)) \leq \text{AR}(f)$.
 141 Then, choosing f to be some minimizer of AR yields the desired result.

142 Define the residuals $r_y(\mathbf{x}_\parallel, \boldsymbol{\delta}_\parallel) := y - h(\mathbf{x}_\parallel + \boldsymbol{\delta}_\parallel)$, and $r_f(\mathbf{x}, \boldsymbol{\delta}) := f(\mathbf{x} + \boldsymbol{\delta}) - h(\mathbf{x}_\parallel + \boldsymbol{\delta}_\parallel)$. Then,
 143 by a decomposition of the squared loss and the tower property of conditional expectation,

$$\begin{aligned} \text{AR}(f) &= \mathbb{E} \left[\mathbb{E} \left[\max_{\|\boldsymbol{\delta}_\parallel\| \leq \varepsilon} r_y(\mathbf{x}_\parallel, \boldsymbol{\delta}_\parallel)^2 + r_f(\mathbf{x}, \boldsymbol{\delta})^2 - 2r_y(\mathbf{x}_\parallel, \boldsymbol{\delta}_\parallel)r_f(\mathbf{x}, \boldsymbol{\delta}) \mid \mathbf{x}_\parallel, y \right] \right] \\ &\geq \mathbb{E} \left[\max_{\|\boldsymbol{\delta}_\parallel\| \leq \varepsilon} r_y(\mathbf{x}_\parallel, \boldsymbol{\delta}_\parallel)^2 + \mathbb{E}[r_f(\mathbf{x}, \boldsymbol{\delta})^2 \mid \mathbf{x}_\parallel, y] - 2r_y(\mathbf{x}_\parallel, \boldsymbol{\delta}_\parallel) \mathbb{E}[r_f(\mathbf{x}, \boldsymbol{\delta}) \mid \mathbf{x}_\parallel, y] \right] \\ &\geq \mathbb{E} \left[\max_{\{\|\boldsymbol{\delta}_\parallel\| \leq \varepsilon, \boldsymbol{\delta}_\perp = 0\}} r_y(\mathbf{x}_\parallel, \boldsymbol{\delta}_\parallel)^2 + \mathbb{E}[r_f(\mathbf{x}, \boldsymbol{\delta})^2 \mid \mathbf{x}_\parallel, y] - 2r_y(\mathbf{x}_\parallel, \boldsymbol{\delta}_\parallel) \mathbb{E}[r_f(\mathbf{x}, \boldsymbol{\delta}) \mid \mathbf{x}_\parallel, y] \right]. \end{aligned}$$

144 Since $y \mid \mathbf{x}_\parallel$ is independent from \mathbf{x}_\perp , for any fixed $\boldsymbol{\delta}$, we have $\mathbb{E}[r_f(\mathbf{x}, \boldsymbol{\delta}) \mid \mathbf{x}_\parallel, y] = \mathbb{E}[r_f(\mathbf{x}, \boldsymbol{\delta}) \mid \mathbf{x}_\parallel]$.
 145 Thus, using the notation $f(\mathbf{x}) = f(\mathbf{x}_\parallel, \mathbf{x}_\perp)$, provided that $\boldsymbol{\delta}_\perp = 0$, Assumption 1 yields

$$h(\mathbf{z} + \boldsymbol{\delta}_\parallel) = \mathbb{E}[f(\mathbf{x}) \mid \mathbf{x}_\parallel = \mathbf{z} + \boldsymbol{\delta}_\parallel] = \mathbb{E}[f(\mathbf{z} + \boldsymbol{\delta}_\parallel, \mathbf{x}_\perp + \boldsymbol{\delta}_\perp)] = \mathbb{E}[f(\mathbf{x} + \boldsymbol{\delta}) \mid \mathbf{x}_\parallel = \mathbf{z}],$$

146 for all $\mathbf{z} \in \mathbb{R}^k$. Plugging in $\mathbf{z} = \mathbf{x}_\parallel$ gives $\mathbb{E}[r_f(\mathbf{x}, \boldsymbol{\delta}) \mid \mathbf{x}_\parallel] = 0$. Therefore,

$$\begin{aligned} \text{AR}(f) &\geq \mathbb{E} \left[\max_{\{\|\boldsymbol{\delta}_\parallel\| \leq \varepsilon, \boldsymbol{\delta}_\perp = 0\}} r_y(\mathbf{x}_\parallel, \boldsymbol{\delta}_\parallel)^2 + \mathbb{E}[r_f(\mathbf{x}, \boldsymbol{\delta})^2 \mid \mathbf{x}_\parallel, y] \right] \\ &\geq \mathbb{E} \left[\max_{\{\|\boldsymbol{\delta}_\parallel\| \leq \varepsilon, \boldsymbol{\delta}_\perp = 0\}} (y - h(\mathbf{U}(\mathbf{x} + \boldsymbol{\delta})))^2 \right] = \text{AR}(h(\mathbf{U}\cdot)), \end{aligned}$$

147 where we dropped the constraint $\boldsymbol{\delta}_\perp = 0$ as it does not contribute, which concludes the proof. \square

148 Before moving to the next section, we provide the following remark on proper scaling of ε .

149 Since $\mathbb{E}[\|\mathbf{x}\|]$ grows with \sqrt{d} , it may seem natural to scale the adversary budget ε with dimension
 150 as well. However, we provide a simple argument on the contrary. Consider the single-index case
 151 $y = g(\langle \mathbf{u}, \mathbf{x} \rangle)$, and let h be the optimal function constructed in Theorem 1, providing the prediction
 152 function $\mathbf{x} \mapsto h(\langle \mathbf{u}, \mathbf{x} \rangle)$. One can then observe that even a constant order ε is sufficient to incur a
 153 large change in the input of h , e.g., choosing $\boldsymbol{\delta} = \varepsilon \mathbf{u}$ perturbs the input of the predictor by ε . Thus,
 154 this justifies the regime where ε is of constant order compared to the input dimension, which is the
 155 focus in the rest of the paper.

156 4 Learning Procedure and Guarantees

157 As outlined in the previous section, to robustly learn the target model, standard representations \mathbf{U}
 158 suffice. In this section, we consider concrete examples of how a standard feature learning oracle

159 combined with an adversarially robust second layer training leads to robust learning. We assume
 160 access to either of the following *feature learning oracles* to recover U . We will provide instances of
 161 practical implementations of these oracles using standard gradient-based algorithms in Section 4.1.

Definition 2 (DFL). An α -Deterministic Feature Learner (DFL) is an oracle that for every $\zeta > 0$, given $n_{\text{DFL}}(\zeta)$ samples from \mathcal{P} , returns a weight matrix $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_N)^\top \in \mathbb{R}^{N \times d}$ such that for all $\mathbf{u} \in \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$ with $\|\mathbf{u}\| = 1$, we have

$$\frac{|\{i : \langle \mathbf{w}_i, \mathbf{u} \rangle \geq 1 - \zeta\}|}{N} \geq \alpha \zeta^{(k-1)/2}.$$

162 An α -DFL oracle returns weights such that roughly an α -proportion of them align with (and suffi-
 163 ciently cover) the target subspace. By a packing argument, we can show that the best achievable ratio
 164 is $\alpha = c(k)$ for some constant $c(k) > 0$ depending only on k , which is why we use the normalizing
 165 factor $\zeta^{(k-1)/2}$ above. We show in Section 4.1 that the definition above with a constant order α is
 166 attainable by standard gradient-based algorithms. That said, in the multi-index setting, it is possible
 167 to improve our learning guarantees by considering the following stochastic oracle.

168 **Definition 3 (SFL).** An (α, β) -Stochastic Feature Learner (SFL) is an oracle that for every $\zeta > 0$,
 169 given $n_{\text{SFL}}(\zeta)$ samples from \mathcal{P} , returns a random weight matrix $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_N)^\top \in \mathbb{R}^{N \times d}$,
 170 such that there exists $S \subseteq [N]$ with $|S|/N \geq \alpha$ satisfying $\|\mathbf{w}_i - \mathbf{U}^\top \mathbf{U} \mathbf{w}_i\| \leq \zeta$ for $i \in S$. Further,
 171 $\left(\frac{\mathbf{U} \mathbf{w}_i}{\|\mathbf{U} \mathbf{w}_i\|}\right)_{i \in S} \stackrel{\text{i.i.d.}}{\sim} \mu$, and $\frac{d\mu}{d\tau_k} \geq \beta$, where μ is some measure and τ_k is uniform, both supported on \mathbb{S}^{k-1} .

172 The above oracle essentially defines a random features model in the smaller target subspace, where a
 173 subset of the weights are sampled independently from a distribution that supports all target directions.
 174 We note that an (α, β) -SFL oracle can be used to directly implement an α -DFL oracle; by a standard
 175 union bound argument, one can show $N = \tilde{\Theta}(1/(\alpha\beta\zeta^{(k-1)/2}))$ guarantees the output of (α, β) -SFL
 176 satisfies Definition 2 with high probability. Therefore, while its definition is slightly more involved,
 177 (α, β) -SFL is a more specialized oracle compared to α -DFL.

178 Once the first layer representation is provided by above oracles, we can fix the biases at some random
 179 initialization, and train the second layer weights \mathbf{a} by minimizing the empirical adversarial risk

$$\widehat{\text{AR}}(\mathbf{a}, \mathbf{W}, \mathbf{b}) = \frac{1}{n} \sum_{i=1}^n \max_{\|\delta^{(i)}\| \leq \varepsilon} (f(\mathbf{x}^{(i)} + \delta^{(i)}; \mathbf{a}, \mathbf{W}, \mathbf{b}) - y)^2. \quad (4.1)$$

We formalize the training procedure with two-layer neural networks in Algorithm 1. We highlight

Algorithm 1 Adversarially robust learning with two-layer NNs.

Input: $\zeta, r_a, r_b, \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{n_{\text{FL}}(\zeta)+n}$, $\text{FL} \in \{\alpha\text{-DFL}, (\alpha, \beta)\text{-SFL}\}$.

1: **Phase 1: Feature Learning**

2: $\mathbf{W} = \text{FL}\left(\zeta, \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=n+1}^{n+n_{\text{FL}}(\zeta)}\right)$.

3: **Phase 2: Robust Function Approximation**

4: $b_j \stackrel{\text{iid}}{\sim} \text{Unif}(-r_b, r_b)$ for $1 \leq j \leq N$.

5: $\hat{\mathbf{a}} = \arg \min_{\|\mathbf{a}\| \leq \frac{r_a}{\sqrt{N}}} \widehat{\text{AR}}(\mathbf{a}, \mathbf{W}, \mathbf{b})$.

6: **return** $(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b})$

180 that keeping biases at random initialization while only training the second layer \mathbf{a} performs non-
 181 linear function approximation, and has been used in many prior works on feature learning [DLS22,
 182 MHWSE23, OSSW24]. Further, while $\hat{\mathbf{a}} \mapsto \widehat{\text{AR}}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b})$ is a convex function for fixed \mathbf{W} and
 183 \mathbf{b} since it is a maximum over convex functions, exact training of $\hat{\mathbf{a}}$ in practice may not be entirely
 184 straightforward since the inner maximization is not concave and does not admit a closed-form
 185 solution. In practice, some form of gradient descent ascent algorithm is typically used when training
 186 $\hat{\mathbf{a}}$ [MMS⁺18]. In this work, we do not consider the computational aspect of solving this min-max
 187 problem, and leave that analysis as future work.

189 We will make the following standard tail assumptions on the data distribution.

190 **Assumption 2.** Suppose \mathbf{x} has zero mean and $\mathcal{O}(1)$ subGaussian norm. Furthermore, for all $q \geq 1$,
 191 it holds that $\mathbb{E}[|y|^q]^{1/q} \leq \mathcal{O}(q^{p/2})$ for some constant $p \geq 1$.

192 Note that the condition on y above is mild; for example, it holds for a noisy multi-index model
 193 $y = g(\mathbf{U}\mathbf{x}) + \varsigma$, where ς has $\mathcal{O}(1)$ subGaussian norm and g grows at most polynomially, i.e.,
 194 $|g(\cdot)| \lesssim 1 + |\cdot|^p$. Similarly, we also keep the function class \mathcal{F} quite general and provide our first set
 195 of results for a class of pseudo-Lipschitz functions which is introduced below.

Assumption 3. We assume \mathcal{F} is a class of functions that are pseudo-Lipschitz along the target
 coordinates. Specifically, using the notation $f(\mathbf{x}) = f(\mathbf{x}_{\parallel}, \mathbf{x}_{\perp})$ and defining $\varepsilon_1 := 1 \vee \varepsilon$, we have

$$|f(\mathbf{z}_1, \mathbf{x}_{\perp}) - f(\mathbf{z}_2, \mathbf{x}_{\perp})| \leq L(\mathbf{x}_{\perp}) (\varepsilon_1^{1-p} \|\mathbf{z}_1\|^{p-1} + \varepsilon_1^{1-p} \|\mathbf{z}_2\|^{p-1} + 1) \|\mathbf{z}_1 - \mathbf{z}_2\|$$

196 for all $f \in \mathcal{F}$, all $\mathbf{z}_1, \mathbf{z}_2 \in \mathbb{R}^k$, and some constants L and $p \geq 1$ such that $\mathbb{E}[L(\mathbf{x}_{\perp})] \leq L$.

197 **Remark.** The prefactor ε_1^{1-p} is justified intuitively since the optimal function of the form $h(\mathbf{z}) =$
 198 $\mathbb{E}[f(\mathbf{x}) | \mathbf{U}\mathbf{x} = \mathbf{z}]$ should satisfy $\mathbb{E}[\max_{\|\delta\| \leq \varepsilon} (y - h(\mathbf{U}(\mathbf{x} + \delta)))^2] = \text{AR}^*$, which is bounded,
 199 and does not grow with ε beyond a certain point. This implies that h must be sufficiently smooth
 200 while its input is perturbed, and in particular, its (local) Lipschitz constant should remain bounded
 201 while ε grows, hence the introduction of the prefactor.

202 In Appendix B.1 we focus on a subclass of predictors that are polynomials of a fixed degree p to
 203 achieve finer results. The following theorem presents the first result of this section, which holds under
 204 access to the α -DFL oracle.

205 **Theorem 4.** Suppose Assumptions 1,2,3 hold and the ReLU activation is used. For a tolerance $\epsilon > 0$
 206 define $\tilde{\epsilon} := \epsilon \wedge (\epsilon^2 / \text{AR}^*)$, and for the adversary budget ε recall $\varepsilon_1 := 1 \vee \varepsilon$. Consider Algorithm 1
 207 with FL = α -DFL oracle, $r_a = \tilde{\mathcal{O}}((\varepsilon_1 / \sqrt{\tilde{\epsilon}})^{k+1+1/k} / \alpha)$ and $r_b = \tilde{\mathcal{O}}(\varepsilon_1 (\varepsilon_1 / \sqrt{\tilde{\epsilon}})^{1+1/k})$. Then, if
 208 the number of second phase samples n_{FA} , the number of neurons N , and α -DFL error ζ satisfy

$$n_{\text{FA}} \geq \tilde{\Omega}\left(\frac{\varepsilon_1^4}{\alpha^4 \epsilon^2} \left(\frac{\varepsilon_1}{\sqrt{\tilde{\epsilon}}}\right)^{\mathcal{O}(k)}\right), \quad N \geq \tilde{\Omega}\left(\frac{1}{\alpha \zeta^{(k-1)/2}} \left(\frac{\varepsilon_1}{\sqrt{\tilde{\epsilon}}}\right)^{\mathcal{O}(k)}\right), \quad \zeta \leq \tilde{\mathcal{O}}\left(\left(\frac{\varepsilon_1}{\sqrt{\tilde{\epsilon}}}\right)^{\mathcal{O}(k)}\right),$$

209 we have $\text{AR}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) \leq \text{AR}^* + \epsilon$ with probability at least $1 - n_{\text{FA}}^{-c}$ where $c > 0$ is an absolute
 210 constant. The total sample complexity of Algorithm 1 is given by $n_{\text{total}} = n_{\text{FA}} + n_{\text{DFL}}(\zeta)$.

211 The above theorem states that once the feature learning oracle has recovered the target subspace, the
 212 number of samples and neurons needed for robust learning is independent of the ambient dimension
 213 d . Thus, in a high-dimensional setting, statistical complexity is dominated by the feature learning
 214 oracle, implying that adversarially robust learning is statistically as easy as standard learning.

Arguing about computational complexity is more involved. While the number of neurons required
 is independent of d , in its naive implementation, Phase 2 of Algorithm 1 needs to solve inner
 maximization problems over \mathbb{R}^d , which may be costly. However, once $\mathbf{U} \in \mathbb{R}^{k \times d}$ is estimated in
 Phase 1, we can reduce the input dimension of the network from d to k by projection onto \mathbf{U} , i.e.

$$\sum_{j=1}^N a_j \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle + b_j) \approx \sum_{j=1}^N a_j \sigma(\langle \mathbf{U}\mathbf{w}_j, \mathbf{U}\mathbf{x} \rangle + b_j).$$

215 With this modification, we only need to consider worst-case perturbations over \mathbb{R}^k , thus the computa-
 216 tional complexity of Phase 2 will also be independent of the ambient dimension d .

217 It is possible to remove the dependence on ζ in the number of neurons by instead assuming access to
 218 a (α, β) -SFL oracle, as outlined below.

219 **Theorem 5.** Consider the same setting as Theorem 4, except that we use the (α, β) -SFL oracle in
 220 Algorithm 1 with $r_a = \tilde{\mathcal{O}}((\varepsilon_1 / \sqrt{\tilde{\epsilon}})^{k+1+1/k} / (\alpha\beta))$. Then, the sufficient number of second phase
 221 samples n_{FA} , neurons N , and oracle error ζ , are given as

$$n_{\text{FA}} \geq \tilde{\Omega}\left(\frac{\varepsilon_1^4}{\alpha^4 \beta^4 \epsilon^2} \left(\frac{\varepsilon_1}{\sqrt{\tilde{\epsilon}}}\right)^{\mathcal{O}(k)}\right), \quad N \geq \tilde{\Omega}\left(\frac{\varepsilon_1^4}{\alpha \beta^2} \left(\frac{\varepsilon_1}{\sqrt{\tilde{\epsilon}}}\right)^{\mathcal{O}(k)}\right), \quad \zeta \leq \tilde{\mathcal{O}}\left(\beta^2 \left(\frac{\varepsilon_1}{\sqrt{\tilde{\epsilon}}}\right)^{\mathcal{O}(k)}\right).$$

222 The total sample complexity for this oracle reads $n_{\text{total}} = n_{\text{FA}} + n_{\text{SFL}}(\zeta)$.

223 We restate Theorems 4 and 5 in Appendix B.2 with explicit exponents. As mentioned earlier, under a
 224 Gaussian data assumption, there exist α -DFL, and more generally (α, β) -SFL oracles which only
 225 use standard gradient-based learning, such that $n_{\text{DFL}}(\zeta)$ typically scales with some polynomial of
 226 d , where the exponent depends on properties of the activation such as the information or generative
 227 exponent. In the following, we will provide specific examples of prior work implementing either of
 228 the oracles, along with their corresponding sample complexities.

229 4.1 Oracle Implementations of the Feature Learner

230 The task of recovering the target directions U is classical in statistics, and is known as sufficient
 231 dimension reduction [LD89, Li91], with many dedicated algorithms, see e.g. [KKSK11, DH18,
 232 CM20, YXKH23] to name a few. Here, we focus on algorithms based on neural networks and
 233 iterative gradient-based optimization.

234 While we will consider the case where x is an isotropic Gaussian random vector, recovering the hidden
 235 direction has also been considered for non-isotropic Gaussians [BES⁺23, MHWSE23] where the
 236 additional structure in the inputs can provide further statistical benefits, or non-Gaussian spherically
 237 symmetric distributions [ZPVB23]. Our results readily extend to these settings as well.

238 First, we present the case of single-index polynomials.

239 **Proposition 6** ([LOSW24]). *Suppose $x \sim \mathcal{N}(0, \mathbf{I}_d)$, $k = 1$, and g is a polynomial of degree p where
 240 p is constant. Then, there exists an iterative first-order algorithm on two-layer neural networks
 241 (Algorithm 2) that implements an (α, β) -SFL oracle and an α -DFL oracle, where $\beta = 1$ and
 242 $\alpha = \hat{\Theta}(1)$. Furthermore, we have $n_{\text{SFL}}(\zeta) = n_{\text{DFL}}(\zeta) = \tilde{\mathcal{O}}(d/\zeta^2)$.*

243 When considering Gaussian single-index models beyond polynomials, we must introduce the concepts
 244 of *information* and *generative exponent* to characterize the sample complexity of recovering the target
 245 direction. Let $\gamma = \mathcal{N}(0, 1)$ for conciseness. For any $g : \mathbb{R} \rightarrow \mathbb{R}$ in $L^2(\gamma)$, let $g = \sum_{j \geq 0} \alpha_j \text{He}_j$
 246 denote its Hermite expansion, where He_j is the normalized Hermite polynomial of degree j . The
 247 information exponent of g is defined as $s(g) := \min\{j > 0 : \alpha_j \neq 0\}$. The generative exponent on
 248 the other hand, is defined as the minimum information exponent attainable by any transformation
 249 of g , i.e. $s^*(g) := \min_{\mathcal{T}} s(\mathcal{T}(g))$, where the minimum is over all $\mathcal{T} \in L^2(g \# \gamma)$. As a result,
 250 $s^*(g) \leq s(g)$, and in particular, $s^* = 1$ for all polynomials.

251 There exists an algorithm based on estimating partial traces that implements a 1-DFL (or a 1,1-SFL)
 252 oracle with $n_{\text{DFL}}(\zeta) = \mathcal{O}(d^{s^*/2} + d/\zeta^2)$ [DPVLB24]. While it may be possible to achieve a similar
 253 sample complexity when training neural networks with a ReLU activation, the state of the art results
 254 for ReLU neural networks so far are only able to control the sample complexity with the information
 255 exponent s , e.g. [BBS22] provides a gradient-based algorithm for optimizing a variant of a two-layer
 256 ReLU neural network that implements 1-DFL with $n_{\text{DFL}} = \mathcal{O}(d^s \text{poly}(\zeta^{-1}))$.

257 Recovering U with $k > 1$ is more challenging, and the general picture is that the directions
 258 in U are recovered hierarchically based on each direction's corresponding complexity, such as
 259 in [ABAM23]. For simplicity, we look at a case that is sufficiently simple for all directions to be
 260 learned simultaneously, while emphasizing that in principle any guarantee for learning the subspace
 261 U can be turned into an implementation of the oracles introduced in the section above.

262 **Proposition 7** ([DLS22]). *Suppose $x \sim \mathcal{N}(0, \mathbf{I}_d)$, g is a polynomial of degree p , and p and $k \geq 1$
 263 are constant. Further assume $\frac{\sigma_{\max}(\nabla^2 g)}{\sigma_{\min}(\nabla^2 g)} \geq \kappa$ for some $\kappa > 0$, where $\sigma_{\min}, \sigma_{\max}$ denote the
 264 minimum and maximum singular values respectively. Then, there exists a first-order algorithm
 265 on two-layer ReLU neural networks (Algorithm 3) that implements an (α, β) -SFL and an α -DFL
 266 oracle, where $\beta \geq c_\kappa$ for some constant $c_\kappa > 0$ depending only κ , and $\alpha = 1$. Further, we have
 267 $n_{\text{SFL}}(\zeta) = n_{\text{DFL}}(\zeta) = \tilde{\mathcal{O}}(d^2 + d/\zeta^2)$.*

268 5 Conclusion

269 In this paper, we initiated a theoretical study of the role of feature learning in adversarial robustness
 270 of neural networks. Under ℓ_2 -constrained perturbations, we proved that projecting onto the latent
 271 subspace of a multi-index model is sufficient for achieving Bayes optimal adversarial risk with respect
 272 to the squared loss, provided that the index directions are statistically independent from the rest

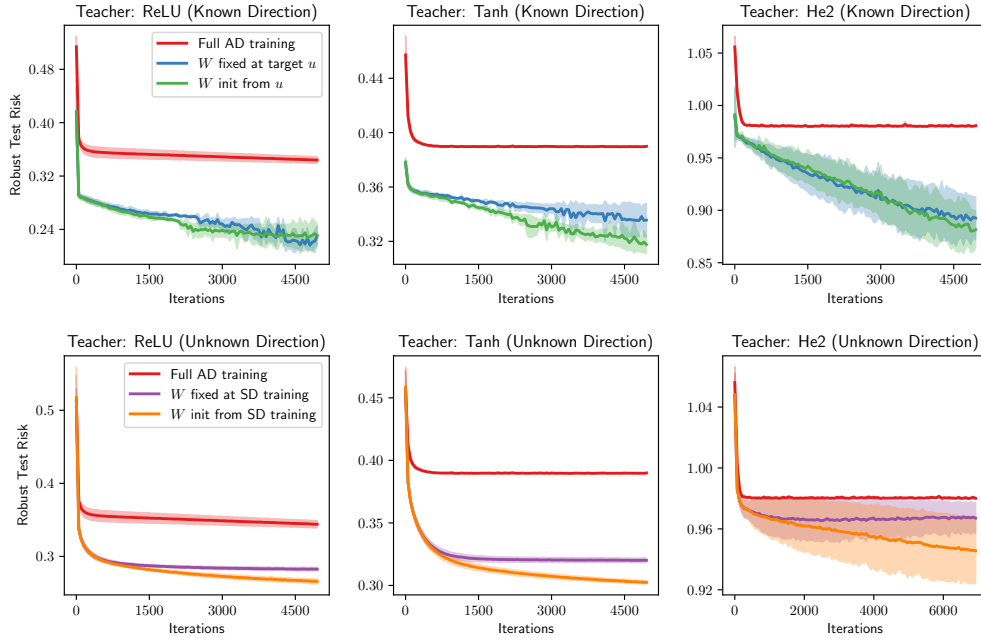


Figure 1: The adversarial test error of a two-layer ReLU network as a function of number of adversarial training iterations to learn a single-index model, where each iteration is performed on a batch of independent 300 samples, except 500 samples for He2 with unknown direction to reduce variance. Full AD training performs adversarial training on all layers from random initialization. SD training is standard training, which provides a better initialization for W before performing adversarial training. We use $\varepsilon = 1$ for all experiments. Experiments are averaged over 3 runs. See Appendix E for details.

273 of the directions in the input space. Remarkably, this subspace can be estimated through standard
 274 feature learning with neural networks, thus turning a high-dimensional robust learning problem into a
 275 low-dimensional one. As a result, under the assumption of having access to a feature learner oracle
 276 which returns an estimate of this subspace, which can be implemented e.g. by training the first-layer
 277 of a two-layer neural network, we proved that robust learning of multi-index models is possible with
 278 number of (additional) samples and neurons independent from ambient dimension.

279 We conclude by mentioning several open questions that arise from this work.

- 280 • Stronger notions of adversarial attacks such as ℓ_∞ norm constraints have been widely considered
 281 in empirical works. It remains open to understand optimal low-dimensional representations under
 282 such perturbations.
- 283 • While our work demonstrates that standard training is sufficient for the first layer, it is unclear
 284 what kind of representation is learned when all layers are trained adversarially. In particular,
 285 Figure 1 suggests that adversarial training of the first layer may be suboptimal in this setting, even
 286 with infinitely many samples.
- 287 • Since our main motivation was to show independence from input dimension, the dependence
 288 of our bounds on the final robust test risk suboptimality ϵ are potentially improvable by a more
 289 careful analysis. It is an interesting direction to obtain a sharper dependency and investigate the
 290 optimality of such dependence on ϵ .

291 References

292 [AAM22] Emmanuel Abbe, Enric Boix Adsera, and Theodor Misiakiewicz. The merged-staircase
 293 property: a necessary and nearly sufficient condition for sgd learning of sparse func-
 294 tions on two-layer neural networks. In *Conference on Learning Theory, 2022*.

- 295 [ABAM23] Emmanuel Abbe, Enric Boix-Adsera, and Theodor Misiakiewicz. Sgd learning on
296 neural networks: leap complexity and saddle-to-saddle dynamics. *arXiv preprint*
297 *arXiv:2302.11055*, 2023.
- 298 [ADK⁺24] Luca Arnaboldi, Yatin Dandi, Florent Krzakala, Luca Pesce, and Ludovic Stephan.
299 Repetita iuvant: Data repetition allows sgd to learn high-dimensional multi-index
300 functions. *arXiv preprint arXiv:2405.15459*, 2024.
- 301 [Bac17] Francis Bach. Breaking the curse of dimensionality with convex neural networks. *The*
302 *Journal of Machine Learning Research*, 18(1):629–681, 2017.
- 303 [BAGJ21] Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. Online stochastic gradient
304 descent on non-convex losses from high-dimensional inference. *J. Mach. Learn. Res.*,
305 22:106–1, 2021.
- 306 [BBSS22] Alberto Bietti, Joan Bruna, Clayton Sanford, and Min Jae Song. Learning single-index
307 models with shallow neural networks. In *Advances in Neural Information Processing*
308 *Systems*, 2022.
- 309 [BEG⁺22] Boaz Barak, Benjamin L Edelman, Surbhi Goel, Sham Kakade, Eran Malach, and
310 Cyril Zhang. Hidden Progress in Deep Learning: SGD Learns Parities Near the
311 Computational Limit. *arXiv preprint arXiv:2207.08799*, 2022.
- 312 [BES⁺22] Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg
313 Yang. High-dimensional Asymptotics of Feature Learning: How One Gradient Step
314 Improves the Representation. *arXiv preprint arXiv:2205.01445*, 2022.
- 315 [BES⁺23] Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, and Denny Wu. Learning
316 in the presence of low-dimensional structure: a spiked random matrix perspective.
317 *Advances in Neural Information Processing Systems*, 36, 2023.
- 318 [BFT17] Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky. Spectrally-normalized margin
319 bounds for neural networks. *Advances in neural information processing systems*, 30,
320 2017.
- 321 [BLPR19] Sébastien Bubeck, Yin Tat Lee, Eric Price, and Ilya Razenshteyn. Adversarial examples
322 from computational constraints. In *International Conference on Machine Learning*,
323 pages 831–840. PMLR, 2019.
- 324 [CB18] Lénaïc Chizat and Francis Bach. On the Global Convergence of Gradient Descent
325 for Over-parameterized Models using Optimal Transport. In *Advances in Neural*
326 *Information Processing Systems*, 2018.
- 327 [CB20] Lénaïc Chizat and Francis Bach. Implicit Bias of Gradient Descent for Wide Two-layer
328 Neural Networks Trained with the Logistic Loss. In *Conference on Learning Theory*,
329 2020.
- 330 [CG24] Ziang Chen and Rong Ge. Mean-field analysis for learning subspace-sparse polynomi-
331 als with gaussian input. *arXiv preprint arXiv:2402.08948*, 2024.
- 332 [Chi22] Lénaïc Chizat. Convergence rates of gradient methods for convex optimization in the
333 space of measures. *Open Journal of Mathematical Optimization*, 3:1–19, 2022.
- 334 [CM20] Sitan Chen and Raghu Meka. Learning polynomials in few relevant dimensions. In
335 *Conference on Learning Theory*, 2020.
- 336 [COB19] Lénaïc Chizat, Edouard Oyallon, and Francis Bach. On Lazy Training in Differentiable
337 Programming. In *Advances in Neural Information Processing Systems*, 2019.
- 338 [DH18] Rishabh Dudeja and Daniel Hsu. Learning single-index models in gaussian space. In
339 *Conference On Learning Theory*, pages 1887–1930. PMLR, 2018.
- 340 [DKL⁺23] Yatin Dandi, Florent Krzakala, Bruno Loureiro, Luca Pesce, and Ludovic Stephan.
341 Learning two-layer neural networks, one (giant) step at a time. *arXiv preprint*
342 *arXiv:2305.18270*, 2023.

- 343 [DLS22] Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi. Neural Networks can Learn
344 Representations with Gradient Descent. In *Conference on Learning Theory*, 2022.
- 345 [DNGL23] Alex Damian, Eshaan Nichani, Rong Ge, and Jason D Lee. Smoothing the landscape
346 boosts the signal for sgd: Optimal sample complexity for learning single index models.
347 *Advances in Neural Information Processing Systems*, 36, 2023.
- 348 [DPVLB24] Alex Damian, Loucas Pillaud-Vivien, Jason D Lee, and Joan Bruna. The com-
349 putational complexity of learning gaussian single-index models. *arXiv preprint*
350 *arXiv:2403.05529*, 2024.
- 351 [DTA⁺24] Yatin Dandi, Emanuele Troiani, Luca Arnaboldi, Luca Pesce, Lenka Zdeborová, and
352 Florent Krzakala. The benefits of reusing batches for gradient descent in two-layer
353 networks: Breaking the curse of information and leap exponents. *arXiv preprint*
354 *arXiv:2402.03220*, 2024.
- 355 [HJ24] Hamed Hassani and Adel Javanmard. The curse of overparametrization in adversarial
356 training: Precise analysis of robust generalization for random features regression. *The*
357 *Annals of Statistics*, 52(2):441–465, 2024.
- 358 [JGH18] Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural Tangent Kernel: Conver-
359 gence and Generalization in Neural Networks. In *Advances in Neural Information*
360 *Processing Systems*, 2018.
- 361 [JM24] Adel Javanmard and Mohammad Mehrabi. Adversarial robustness for latent models:
362 Revisiting the robust-standard accuracies tradeoff. *Operations Research*, 72(3):1016–
363 1030, 2024.
- 364 [JMS24] Nirmitt Joshi, Theodor Misiakiewicz, and Nathan Srebro. On the complexity of learning
365 sparse functions with statistical and gradient queries. *arXiv preprint arXiv:2407.05622*,
366 2024.
- 367 [JS22] Adel Javanmard and Mahdi Soltanolkotabi. Precise statistical analysis of classification
368 accuracies for adversarial training. *The Annals of Statistics*, 50(4):2127–2156, 2022.
- 369 [JSH20] Adel Javanmard, Mahdi Soltanolkotabi, and Hamed Hassani. Precise tradeoffs in
370 adversarial training for linear regression. In *Conference on Learning Theory*, pages
371 2034–2078. PMLR, 2020.
- 372 [KKSK11] Sham M Kakade, Varun Kanade, Ohad Shamir, and Adam Kalai. Efficient learning
373 of generalized linear and single index models with isotonic regression. *Advances in*
374 *Neural Information Processing Systems*, 24, 2011.
- 375 [LD89] Ker-Chau Li and Naihua Duan. Regression Analysis Under Link Violation. *The*
376 *Annals of Statistics*, 1989.
- 377 [Li91] Ker-Chau Li. Sliced inverse regression for dimension reduction. *Journal of the*
378 *American Statistical Association*, 1991.
- 379 [LOSW24] Jason D. Lee, Kazusato Oko, Taiji Suzuki, and Denny Wu. Neural network learns
380 low-dimensional polynomials with sgd near the information-theoretic limit. *arXiv*
381 *preprint arXiv:2406.01581*, 2024.
- 382 [MHPG⁺23] Alireza Mousavi-Hosseini, Sejun Park, Manuela Girotti, Ioannis Mitliagkas, and
383 Murat A Erdogdu. Neural networks efficiently learn low-dimensional representations
384 with SGD. In *The Eleventh International Conference on Learning Representations*,
385 2023.
- 386 [MHWE24] Alireza Mousavi-Hosseini, Denny Wu, and Murat A Erdogdu. Learning multi-index
387 models with neural networks via mean-field langevin dynamics. *arXiv preprint*
388 *arXiv:2408.07254*, 2024.
- 389 [MHWSE23] Alireza Mousavi-Hosseini, Denny Wu, Taiji Suzuki, and Murat A Erdogdu. Gradient-
390 based feature learning under structured data. *Advances in Neural Information Process-*
391 *ing Systems*, 36, 2023.

- 392 [MMN18] Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the
393 landscape of two-layer neural networks. *Proceedings of the National Academy of*
394 *Sciences*, 115(33):E7665–E7671, 2018.
- 395 [MMS⁺18] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and
396 Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In
397 *International Conference on Learning Representations*, 2018.
- 398 [MZD⁺23] Arvind Mahankali, Haochen Zhang, Kefan Dong, Margalit Glasgow, and Tengyu Ma.
399 Beyond ntk with vanilla gradient descent: A mean-field analysis of neural networks
400 with polynomial width, samples, and time. *Advances in Neural Information Processing*
401 *Systems*, 36, 2023.
- 402 [NOSW24] Atsushi Nitanda, Kazusato Oko, Taiji Suzuki, and Denny Wu. Improved statistical
403 and computational complexity of the mean-field langevin dynamics under structured
404 data. In *The Twelfth International Conference on Learning Representations*, 2024.
- 405 [NWS22] Atsushi Nitanda, Denny Wu, and Taiji Suzuki. Convex analysis of the mean field
406 langevin dynamics. In *International Conference on Artificial Intelligence and Statistics*,
407 pages 9741–9757. PMLR, 2022.
- 408 [OSSW24] Kazusato Oko, Yujin Song, Taiji Suzuki, and Denny Wu. Learning sum of diverse
409 features: computational hardness and efficient gradient-based training for ridge combi-
410 nations. In *Conference on Learning Theory*. PMLR, 2024.
- 411 [Pis81] Gilles Pisier. Remarques sur un résultat non publié de b. maurey. *Séminaire d'Analyse*
412 *fonctionnelle (dit "Maurey-Schwartz")*, pages 1–12, 1981.
- 413 [RVE18] Grant M Rotskoff and Eric Vanden-Eijnden. Neural networks as Interacting Particle
414 Systems: Asymptotic convexity of the Loss Landscape and Universal Scaling of the
415 Approximation Error. *arXiv preprint arXiv:1805.00915*, 2018.
- 416 [SH20] Johannes Schmidt-Hieber. Nonparametric regression using deep neural networks with
417 ReLU activation function. *The Annals of Statistics*, 48(4):1875 – 1897, 2020.
- 418 [SST⁺18] Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, and Aleksander
419 Madry. Adversarially robust generalization requires more data. *Advances in neural*
420 *information processing systems*, 31, 2018.
- 421 [SWON23] Taiji Suzuki, Denny Wu, Kazusato Oko, and Atsushi Nitanda. Feature learning via
422 mean-field langevin dynamics: classifying sparse parities and beyond. In *Thirty-*
423 *seventh Conference on Neural Information Processing Systems*, 2023.
- 424 [SZS⁺14] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan,
425 Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *The*
426 *International Conference on Learning Representations*, 2014.
- 427 [Tel23] Matus Telgarsky. Feature selection and low test error in shallow low-rotation relu
428 networks. In *The Eleventh International Conference on Learning Representations*,
429 2023.
- 430 [TSE⁺18] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and
431 Aleksander Madry. Robustness may be at odds with accuracy. *arXiv preprint*
432 *arXiv:1805.12152*, 2018.
- 433 [VE24] Nuri Mert Vural and Murat A. Erdogdu. Pruning is optimal for learning sparse features
434 in high-dimensions. *arXiv preprint arXiv:2406.08658*, 2024.
- 435 [WLLM19] Colin Wei, Jason D Lee, Qiang Liu, and Tengyu Ma. Regularization matters: General-
436 ization and optimization of neural nets vs their induced kernel. *Advances in Neural*
437 *Information Processing Systems*, 32, 2019.
- 438 [WMHC24] Guillaume Wang, Alireza Mousavi-Hosseini, and Lénaïc Chizat. Mean-field langevin
439 dynamics for signed measures via a bilevel approach. *arXiv preprint arXiv:2406.17054*,
440 2024.

- 441 [XLS⁺24] Jiancong Xiao, Qi Long, Weijie Su, et al. Bridging the gap: Rademacher complexity
442 in robust and standard generalization. In *The Thirty Seventh Annual Conference on*
443 *Learning Theory*, pages 5074–5075. PMLR, 2024.
- 444 [YXKH23] Gan Yuan, Mingyue Xu, Samory Kpotufe, and Daniel Hsu. Efficient estimation
445 of the central mean subspace via smoothed gradient outer products. *arXiv preprint*
446 *arXiv:2312.15469*, 2023.
- 447 [Zha02] Tong Zhang. Covering number bounds of certain regularized linear function classes.
448 *Journal of Machine Learning Research*, 2(Mar):527–550, 2002.
- 449 [ZPVB23] Aaron Zweig, Loucas Pillaud-Vivien, and Joan Bruna. On single-index models beyond
450 gaussian data. *Advances in Neural Information Processing Systems*, 36, 2023.

451 A Gradient-Based Neural Feature Learning Algorithms

In this section, we will provide examples of implementations of the feature learner oracles introduced in Section 4 using gradient-based training of two-layer neural networks. First, we look at the algorithm provided by [OSSW24] for the case where g is a polynomial of degree p . Consider the following two-layer neural network with zero bias

$$f(\mathbf{x}; \mathbf{a}, \mathbf{W}) = \sum_{j=1}^N a_j \sigma_j(\langle \mathbf{w}_j, \mathbf{x} \rangle).$$

452 Note that we allow the activation to vary based on neuron. Specifically, we let $\sigma_j = \sum_{l=1}^q \beta_{j,l} \text{He}_l$,
 453 where He_j is the j th normalized Hermite polynomial, $\beta_{j,l} \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\{\pm r_l\})$ for appropriately chosen
 454 r_l , and $q \geq C_p$, see [OSSW24, Lemma 3] for details. Now, we consider the following algorithm.

Algorithm 2 Gradient-Based Feature Learner for Single-Index Polynomials [OSSW24, Algorithm 1, Phase I].

Input: T , step size $(\eta^t)_{t=0}^{T-1}$, momentum parameters $(\zeta_j^t), r_a$.

- 1: $\mathbf{w}_j^0 \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\mathbb{S}^{d-1}), a_j \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\{\pm r_a/N\}), \forall j \in [N]$.
- 2: $(\mathbf{x}^{(0)}, y^{(0)}) \sim \mathcal{P}$
- 3: **for** $t = 0, \dots, T-1$ **do**
- 4: **if** $t > 0$ and t is even **then**
- 5: Draw $(\mathbf{x}^{(t/2)}, y^{(t/2)}) \sim \mathcal{P}$
- 6: $\mathbf{w}_j^t \leftarrow \mathbf{w}_j^{t-1} - \zeta_j^t (\mathbf{w}_j^{t-1} - \mathbf{w}_j^{t-2}), \forall j \in [N]$
- 7: $\mathbf{w}_j^t \leftarrow \frac{\mathbf{w}_j^t}{\|\mathbf{w}_j^t\|} \quad \forall j \in [N]$
- 8: **end if**
- 9: $\mathbf{w}_j^{t+1} \leftarrow \mathbf{w}_j^t - \eta_t \nabla_{\mathbf{w}_j}^S (f(\mathbf{x}^{(\lfloor t/2 \rfloor)}; \mathbf{a}, \mathbf{W}^t) - y^{(\lfloor t/2 \rfloor)})^2$
- 10: **end for**
- 11: **return** $(\mathbf{w}_0^T, \dots, \mathbf{w}_N^T)^\top$

455 Note that $\nabla^S f(\mathbf{w}) = (\mathbf{I} - \mathbf{w}\mathbf{w}^\top) \nabla f(\mathbf{w})$ denotes the spherical gradient. Essentially, the above
 456 algorithm takes two gradient steps on each new sample, and in the even iterations performs a certain
 457 interpolation. Proper choice of hyperparameters in the above algorithm leads to Proposition 6.

Next, we consider the algorithm of [DLS22] for the case where g is a multi-index polynomial.

Algorithm 3 Gradient-Based Feature Learner for Multi-Index Polynomials [DLS22, Algorithm 1, Adapted]

Input: $\{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{n_{\text{FL}}}, r_a$

- 1: $\mathbf{a}_j \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\{\pm r_a\}), \mathbf{w}_j^0 \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\mathbb{S}^{d-1}), \mathbf{a}_{N-j} = -\mathbf{a}_j, \mathbf{w}_{N-j} = \mathbf{w}_j^0, \forall j \in [N/2]$.
- 2: $\alpha \leftarrow \frac{1}{n_{\text{FL}}} \sum_{i=1}^{n_{\text{FL}}} y^{(i)}, \beta \leftarrow \frac{1}{n_{\text{FL}}} \sum_{i=1}^{n_{\text{FL}}} y^{(i)} \mathbf{x}^{(i)}$
- 3: $y^{(i)} \leftarrow y^{(i)} - \alpha - \langle \beta, \mathbf{x}^{(i)} \rangle, \forall i \in [n_{\text{FL}}]$.
- 4: $\mathbf{W} \leftarrow -\nabla_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^{n_{\text{FL}}} (f(\mathbf{x}^{(i)}; \mathbf{a}, \mathbf{W}^0) - y)^2$
- 5: $\mathbf{w}_i \leftarrow \frac{\mathbf{w}_i}{\|\mathbf{w}_i\|}, \forall i \in [N]$
- 6: **return** $(\mathbf{w}_0, \dots, \mathbf{w}_N)^\top$

458

After performing a preprocessing on data, the above essentially performs one gradient descent step with weight decay, when the regularizer of the weight decay is the inverse of step size, thus cancelling out initialization and leaving only gradient as the estimate. [DLS22] prove that, with a sample complexity of $n_{\text{DFL}} = \tilde{O}(d^2 + d/\zeta^2)$, the output of Algorithm 3 satisfies

$$\left\langle \mathbf{w}_i, \frac{\mathbf{U}^\top \mathbf{H} \mathbf{U} \mathbf{w}_i^0}{\|\mathbf{U}^\top \mathbf{H} \mathbf{U} \mathbf{w}_i^0\|} \right\rangle \geq 1 - \zeta, \quad \forall i \in [N],$$

459 with high probability, where $\mathbf{H} = \mathbb{E}[\nabla^2 g(\mathbf{U}\mathbf{x})]$. Thus, for a full-rank \mathbf{H} , the output of Algorithm 3
 460 satisfies the definition of a $(1, \beta)$ SFL oracle for a constant $\beta > 0$ depending only on the conditioning
 461 of H and the number of indices k .

462 B Additional Details of Section 4

463 Throughout the appendix, we will assume the activation satisfies $\sigma(0) = 0$ for simplicity of presenta-
 464 tion, without loss of generality. We will also assume that

$$|\sigma(z_1) - \sigma(z_2)| \leq L_\sigma(|z_1|^{q-1} + |z_2|^{q-1} + 1)|z_1 - z_2|, \quad (\text{B.1})$$

465 for all $z_1, z_2 \in \mathbb{R}$ and some absolute constant L_σ . In the case of ReLU, we have $q = 1$ and $L_\sigma = 1$.
 466 For polynomial activations, q is the same as the degree of the polynomial. For a set of parameters ψ
 467 (e.g. $\psi = q, k$), we will use C_ψ to denote a generic constant whose value depends only on ψ and may
 468 change from line to line.

469 B.1 Competing against the Optimal Polynomial Predictor

470 In this section, we restrict \mathcal{F} to only polynomials, which allows us to derive more refined bounds on
 471 the number of samples and neurons. Specifically, we make the following assumption.

472 **Assumption 4.** *Suppose \mathcal{F} is the class of d -variate polynomials of degree p for some constant $p > 0$.
 473 Further, σ is either the ReLU activation or a polynomial of degree $q \geq p$.*

474 While the ReLU activation is sufficient for approximation purposes, we also consider polynomial
 475 activations in Assumption 4 since recent works have been able to achieve sharper guarantees of
 476 recovering the target directions under such activations. We provide a more detailed discussion in
 477 Section 4.1. Note that a priori we do not require a growth constraint on the coefficients of the
 478 polynomials in \mathcal{F} . The optimal function h in Theorem 1 automatically chooses a polynomial with
 479 suitably bounded coefficients in order to avoid incurring a large robust risk.

480 The following result establishes the sample and computational complexity for competing against
 481 polynomial predictors when having access to the α -DFL oracle.

482 **Theorem 8.** *Suppose Assumptions 1, 2, 4 hold. For a tolerance $\epsilon > 0$ define $\tilde{\epsilon} := \epsilon \wedge (\epsilon^2 / \text{AR}^*)$, and
 483 for the adversary budget ϵ recall $\varepsilon_1 := 1 \vee \epsilon$. Consider Algorithm 1 with α -DFL oracle, $r_a = \tilde{\mathcal{O}}(1)$,
 484 $r_b = \tilde{\mathcal{O}}(\varepsilon_1)$. Then, if the number of second phase samples n_{FA} , neurons N , and α -DFL error ζ
 485 satisfy*

$$n_{\text{FA}} \geq \tilde{\Omega}\left(\frac{\varepsilon_1^{4(q+1)}}{\alpha^4 \epsilon^2}\right), \quad N \geq \tilde{\Omega}\left(\frac{\varepsilon_1^{q+1}}{\alpha \zeta^{\frac{k-1}{2}} \sqrt{\tilde{\epsilon}}}\right), \quad \zeta \leq \tilde{\mathcal{O}}\left(\frac{\tilde{\epsilon}}{\varepsilon_1^{2(q+1)}}\right),$$

486 we have $\text{AR}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) \leq \text{AR}^* + \epsilon$ with probability at least $1 - n^{-c}$ where $c > 0$ is an absolute
 487 constant. The total sample complexity of Algorithm 1 is given by $n_{\text{total}} = n_{\text{FA}} + n_{\text{DFL}}(\zeta)$.

488 Consequently, when restricting \mathcal{F} to the class of fixed degree polynomials, there is no curse of
 489 dimensionality for sample complexity, even in the latent dimension k . This is consistent with the
 490 setting in standard learning, see e.g. [CM20]. Further, similar to the general case above, it is possible
 491 to remove the ζ dependence from N when having access to an SFL oracle, thus also achieving
 492 computational complexity as a fixed polynomial independent of the latent dimension.

493 **Theorem 9.** *In the setting of Theorem 8, consider using Algorithm 1 with an (α, β) -SFL oracle.
 494 Then, the sufficient number of second phase samples and neurons are given as*

$$n_{\text{FA}} \geq \tilde{\Omega}\left(\frac{\varepsilon_1^{4(q+1)}}{\alpha^4 \beta^4 \epsilon^2}\right), \quad N \geq \tilde{\Omega}\left(\frac{\varepsilon_1^{2(q+1)}}{\alpha \beta^2 \tilde{\epsilon}}\right), \quad \zeta \leq \tilde{\mathcal{O}}\left(\frac{\beta^2 \tilde{\epsilon}}{\varepsilon_1^{2(q+1)}}\right)$$

495 The total sample complexity is given by $n_{\text{total}} = n_{\text{FA}} + n_{\text{SFL}}(\zeta)$ for ζ as in Theorem 8.

496 **Remark.** We note that the guarantees provided in Theorem 9 are generally better than those in
 497 Theorem 8 for large k ; yet, they are strictly worse for $k = 1$. That said, both Theorems 9 and
 498 8 respectively achieve better sample complexity guarantees compared to their counterparts in the
 499 previous section, namely Theorems 4 and 5, simply by restricting the function class \mathcal{F} to polynomials.

500

501 B.2 Complete Versions of Theorems in Section 4

502 We first restate Theorem 4 with explicit exponents.

503 **Theorem 10.** *Suppose Assumptions 1, 2, and 3 hold. For any $\epsilon > 0$, define $\tilde{\epsilon} := \epsilon \wedge (\epsilon^2 / \text{AR}^*)$, and*
 504 *recall $\varepsilon_1 := 1 \vee \varepsilon$. Consider Algorithm 1 with the α -DFL oracle, $r_a = \tilde{\mathcal{O}}((\varepsilon_1 / \sqrt{\tilde{\epsilon}})^{k+1+1/k} / \alpha)$, and*
 505 *$r_b = \tilde{\mathcal{O}}(\varepsilon_1 (\varepsilon_1 / \sqrt{\tilde{\epsilon}})^{1+1/k})$. Then, if the number of second phase samples n_{FA} , number of neurons N ,*
 506 *and α -DFL error ζ satisfy*

$$n_{\text{FA}} \geq \tilde{\Omega}\left(\frac{\varepsilon_1^4}{\alpha^4 \varepsilon^2} \left(\frac{\varepsilon_1^2}{\tilde{\epsilon}}\right)^{2k+4+4/k}\right), \quad N \geq \tilde{\Omega}\left(\frac{1}{\alpha \zeta^{(k-1)/2}} \left(\frac{\varepsilon_1}{\sqrt{\tilde{\epsilon}}}\right)^{k+3+2/k}\right), \quad \zeta \leq \tilde{\mathcal{O}}\left(\left(\frac{\tilde{\epsilon}}{\varepsilon_1^2}\right)^{k+2+1/k}\right),$$

507 *we have $\text{AR}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) \leq \text{AR}^* + \epsilon$ with probability at least $1 - n_{\text{FA}}^{-c}$ where $c > 0$ is an absolute*
 508 *constant. The total sample complexity of Algorithm 1 is given by $n_{\text{total}} = n_{\text{FA}} + n_{\text{DFL}}(\zeta)$.*

509 Similarly, we can restate Theorem 5 with explicit exponents.

510 **Theorem 11.** *Consider the same setting as Theorem 10, except that we use the (α, β) -SFL oracle*
 511 *in Algorithm 1 with $r_a = \tilde{\mathcal{O}}((\varepsilon_1 / \sqrt{\tilde{\epsilon}})^{k+1+1/k} / (\alpha\beta))$. Then, if the number of second phase samples*
 512 *n_{FA} , number of neurons N , and α -DFL error ζ satisfy*

$$n \geq \tilde{\Omega}\left(\frac{\varepsilon_1^4}{\alpha^4 \beta^4 \varepsilon^2} \left(\frac{\varepsilon_1^2}{\tilde{\epsilon}}\right)^{2k+4+4/k}\right), \quad N \geq \tilde{\Omega}\left(\frac{1}{\alpha \beta^2} \left(\frac{\varepsilon_1^2}{\tilde{\epsilon}}\right)^{k+3+2/k}\right), \quad \zeta \leq \tilde{\mathcal{O}}\left(\beta^2 \left(\frac{\tilde{\epsilon}}{\varepsilon_1^2}\right)^{k+2+1/k}\right).$$

513 *while the oracle error tolerance ζ stays the same. The total sample complexity in this case is given by*
 514 *$n_{\text{total}} = n + n_{\text{SFL}}(\zeta)$.*

515 The proof of both theorems follows from combining the results of the following sections. Since both
 516 proofs are similar, we only present the proof of Theorem 10. The proof of Theorems 8 and 9 can be
 517 obtained in a similar manner.

Proof. [Proof of Theorem 10] The proof is based on decomposing the suboptimality into generaliza-
 tion and approximation terms, namely

$$\text{AR}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) - \text{AR}^* = \text{AR}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) - \text{AR}(\mathbf{a}^*, \mathbf{W}, \mathbf{b}) + \text{AR}(\mathbf{a}^*, \mathbf{W}, \mathbf{b}) - \text{AR}^*,$$

518 where $\mathbf{a}^* := \min_{\|\mathbf{a}\| \leq r_a / \sqrt{N}} \text{AR}(\mathbf{a}, \mathbf{W}, \mathbf{b})$, thus we can see the first term above as generalization
 519 error, and the second term as approximation error.

520 From Proposition 20, we have $\text{AR}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) - \text{AR}(\mathbf{a}^*, \mathbf{W}, \mathbf{b}) \leq \epsilon/2$ as soon as $n \geq \tilde{\Omega}(r_a^4 (\varepsilon_1^4 +$
 521 $r_b^4 / \varepsilon^2))$ (recall that $q = 1$ here, since we are considering the ReLU activation). For the approximation
 522 error, we can use Proposition 34, which guarantees there exists \mathbf{a}^* with $\|\mathbf{a}^*\| \leq r_a / \sqrt{N}$ such that
 523 $\text{AR}(\mathbf{a}^*, \mathbf{W}, \mathbf{b}) - \text{AR}^* \leq \epsilon/2$ with $r_a \leq \tilde{\mathcal{O}}((\varepsilon_1 / \sqrt{\tilde{\epsilon}})^{k+1+1/k} / \alpha)$, as soon as

$$\zeta \leq \tilde{\mathcal{O}}\left(\left(\frac{\tilde{\epsilon}}{\varepsilon_1^2}\right)^{k+2+1/k}\right), \quad \text{and} \quad N \geq \tilde{\Omega}\left(\frac{1}{\zeta^{(k-1)/2} \alpha} \left(\frac{\varepsilon_1}{\sqrt{\tilde{\epsilon}}}\right)^{k+3+2/k}\right),$$

524 provided that we choose $r_b = \tilde{\Theta}(\varepsilon_1 (\varepsilon_1 / \sqrt{\tilde{\epsilon}})^{1+1/k})$. Plugging the value of r_a and r_b in the bound for
 525 n completes the proof. \square

526 C Generalization Analysis

527 We will first focus on proving a generalization bound for bounded and Lipschitz losses, and then
 528 extend the results to cover the squared loss.

529 C.1 Generalization Bounds for Bounded Lipschitz Losses

Let us focus on a general C_ℓ Lipschitz loss $\ell(f(\cdot; \mathbf{a}, \mathbf{W}, \mathbf{b}) - y)$ for now. Later, we will argue how
 to extend the results of this section to the squared error loss. Our uniform convergence argument
 depends on the covering number of the family of adversarial loss functions. Let $\Theta \subseteq \mathbb{R}^N$ be the set
 of second layer weights, to be determined later. This family is given by

$$\mathcal{L}(\mathbf{W}, \mathbf{b}) = \{(x, y) \mapsto \max_{\|\delta\| \leq \varepsilon} \ell(f(x + \delta; \mathbf{a}, \mathbf{W}, \mathbf{b}) - y) : \mathbf{a} \in \Theta\}.$$

For brevity, we will also use \mathcal{L} to denote $\mathcal{L}(\mathbf{W}, \mathbf{b})$, but we highlight that \mathbf{W} and \mathbf{b} are fixed at this stage. We define the following metric over this family

$$\forall \tilde{l}, \tilde{l}' \in \mathcal{L}(\mathbf{W}, \mathbf{b}), \quad d_{\mathcal{L}}(\tilde{l}, \tilde{l}')^2 := \frac{1}{n} \sum_{i=1}^n (\tilde{\ell}(\mathbf{x}^{(i)}, y^{(i)}) - \tilde{\ell}'(\mathbf{x}^{(i)}, y^{(i)}))^2.$$

530 We say $\mathcal{S} \subseteq \mathcal{L}$ is an ϵ -cover of \mathcal{L} if for every $\tilde{l} \in \mathcal{L}$, there exists $\tilde{l}' \in \mathcal{S}$ such that $d_{\mathcal{L}}(\tilde{l}, \tilde{l}') \leq \epsilon$.
 531 The ϵ -covering number of \mathcal{L} is the least cardinality among all ϵ -covers of \mathcal{L} , which we denote
 532 by $\mathcal{C}(\mathcal{L}, d_{\mathcal{L}}, \epsilon)$. Note that since \mathcal{L} is parameterized by \mathbf{a} , constructing such a covering reduces to
 533 constructing a finite set over Θ .

Therefore, we define the following metric over Θ ,

$$\forall \mathbf{a}, \mathbf{a}' \in \Theta, \quad d_{\Theta}(\mathbf{a}, \mathbf{a}')^2 := \frac{1}{n} \sum_{i=1}^n \max_{\|\delta^{(i)}\| \leq \epsilon} (f(\mathbf{x}^{(i)} + \delta^{(i)}; \mathbf{a}, \mathbf{W}, \mathbf{b}) - f(\mathbf{x}^{(i)} + \delta^{(i)}; \mathbf{a}', \mathbf{W}, \mathbf{b}))^2.$$

534 We can similarly define the ϵ -covering number of Θ with respect to the metric d_{Θ} as $\mathcal{C}(\Theta, d_{\Theta}, \epsilon)$.
 535 The following lemma relates the covering numbers of \mathcal{L} and Θ .

536 **Lemma 12.** *We have $\mathcal{C}(\mathcal{L}, d_{\mathcal{L}}, \epsilon) \leq \mathcal{C}(\Theta, d_{\Theta}, \epsilon/C_{\ell})$ for all $\epsilon > 0$.*

537 **Proof.** We will use the following fact in the proof. For any $F_1, F_2 : S \rightarrow \mathbb{R}$, we have

$$\left| \max_{\delta_1 \in S} F_1(\delta_1) - \max_{\delta_2 \in S} F_2(\delta_2) \right| \leq \max_{\delta \in S} |F_1(\delta) - F_2(\delta)|. \quad (\text{C.1})$$

This is true because

$$\max_{\delta_1 \in S} F_1(\delta_1) - \max_{\delta_2 \in S} F_2(\delta_2) \leq \max_{\delta_1 \in S} \{F_1(\delta_1) - F_2(\delta_1)\},$$

538 and the other direction holds by symmetry. This trick is used to relate the adversarial loss to its
 539 non-adversarial counterpart, e.g. in [XLS⁺24, Lemma 5].

Now, we will show that an ϵ/C_{ℓ} cover for Θ implies an ϵ cover for \mathcal{L} . We will suppress dependence on the fixed \mathbf{W} and \mathbf{b} in the notation. Let \mathcal{S}_{Θ} be an ϵ/C_{ℓ} cover of Θ with respect to the d_{Θ} metric. Then, we define \mathcal{S} via

$$\mathcal{S} = \{(\mathbf{x}, y) \mapsto \max_{\|\delta\| \leq \epsilon} \ell(f(\mathbf{x} + \delta; \mathbf{a}) - y) : \mathbf{a} \in \mathcal{S}_{\Theta}\}.$$

540 To show \mathcal{S} is an ϵ cover of \mathcal{L} , consider an arbitrary $\tilde{\ell}(\mathbf{x}, y) = \max_{\|\delta\| \leq \epsilon} \ell(f(\mathbf{x} + \delta; \mathbf{a}) - y)$. Suppose
 541 \mathbf{a}' is the closest element to \mathbf{a} in \mathcal{S}_{Θ} , and let $\tilde{\ell}'(\mathbf{x}, y) = \max_{\|\delta\| \leq \epsilon} \ell(f(\mathbf{x} + \delta; \mathbf{a}') - y)$. Then,

$$\begin{aligned} d_{\mathcal{L}}(\tilde{\ell}, \tilde{\ell}')^2 &= \frac{1}{n} \sum_{i=1}^n \left(\max_{\|\delta_1^{(i)}\| \leq \epsilon} \ell(f(\mathbf{x} + \delta_1^{(i)}; \mathbf{a}) - y^{(i)}) - \max_{\|\delta_2^{(i)}\| \leq \epsilon} \ell(f(\mathbf{x} + \delta_2^{(i)}; \mathbf{a}') - y^{(i)}) \right)^2 \\ &\leq \frac{1}{n} \sum_{i=1}^n \max_{\|\delta^{(i)}\| \leq \epsilon} (\ell(f(\mathbf{x} + \delta^{(i)}; \mathbf{a}) - y^{(i)}) - \ell(f(\mathbf{x} + \delta^{(i)}; \mathbf{a}') - y^{(i)}))^2 \\ &\leq \frac{C_{\ell}^2}{n} \sum_{i=1}^n \max_{\|\delta^{(i)}\| \leq \epsilon} (f(\mathbf{x} + \delta^{(i)}; \mathbf{a}) - f(\mathbf{x} + \delta^{(i)}; \mathbf{a}'))^2 \\ &\leq C_{\ell}^2 d_{\Theta}(\mathbf{a}, \mathbf{a}')^2 \leq \epsilon^2, \end{aligned}$$

542 where we used (C.1) for the first inequality. \square

543 To construct an ϵ -cover of Θ , we depend on the Maurey sparsification lemma [Pis81], which has been
 544 used in the literature for providing covering numbers for linear classes [Zha02] and neural networks
 545 via matrix covering, see e.g. [BFT17].

Lemma 13 (Maurey Sparsification Lemma, [Zha02, Lemma 1]). *Let \mathcal{H} be a Hilbert space with norm $\|\cdot\|$, let $\mathbf{u} \in \mathcal{H}$ be represented by $\mathbf{u} = \sum_{j=1}^m \alpha_j \mathbf{v}_j$, where $\alpha_j \geq 0$ and $\|\mathbf{v}_j\| \leq b$ for all $j \in [m]$, and $\alpha = \sum_{j=1}^m \alpha_j \leq 1$. Then, for every $k \geq 1$, there exist non-negative integers k_1, \dots, k_m , such that*

$$\left\| \mathbf{u} - \frac{1}{k} \sum_{j=1}^m k_j \mathbf{v}_j \right\| \leq \frac{\alpha b^2 - \|\mathbf{u}\|^2}{k}.$$

546 Then, we have the following upper bound on the the covering number of Θ .

Lemma 14. *Suppose σ is either ReLU or a polynomial of degree $q \geq 1$, $\Theta = \{\|\mathbf{a}\|_1 \leq r_a\}$, and additionally $\|\mathbf{w}_i\| \leq r_w$ and $|b_i| \leq r_b$ for all $1 \leq i \leq N$. Then we have*

$$\log \mathcal{C}(\Theta, d_\Theta, \epsilon) \leq \frac{C_q L_\sigma^2 r_a^2 \log N \left\{ T_{\mathbf{W}, \mathbf{X}}^{(q)} + r_w^{2q} \epsilon^{2q} + r_b^{2q} + T_{\mathbf{W}, \mathbf{X}}^{(2)} + r_w^2 \epsilon^2 + r_b^2 \right\}}{\epsilon^2},$$

547 where $T_{\mathbf{W}, \mathbf{X}}^{(q)} := \max_{1 \leq j \leq N} \frac{1}{n} \sum_{i=1}^n \langle \mathbf{w}_j, \mathbf{x}_i \rangle^{2q}$.

Proof. Given some positive integer $k > 0$, let \mathcal{S}_Θ be given by the following

$$\mathcal{S}_\Theta = \left\{ \frac{r_a}{k} (k_1 - k'_1, k_2 - k'_2, \dots, k_N - k'_N)^\top : \forall i, k_i, k'_i \geq 0, \sum_{i=1}^N k_i + \sum_{i=1}^N k'_i = k \right\}.$$

Let $\mathbf{X}, \mathbf{\Delta} \in \mathbb{R}^{n \times d}$ be the matrices with (\mathbf{x}_i) and (δ_i) as rows respectively. Let $\mathbf{A} = \sigma((\mathbf{X} + \mathbf{\Delta})\mathbf{W}^\top + \mathbf{1}_n \mathbf{b}^\top) \in \mathbb{R}^{n \times m}$. Then,

$$\frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}^{(i)} + \delta^{(i)}; \mathbf{a}, \mathbf{W}, \mathbf{b}) - f(\mathbf{x}^{(i)} + \delta^{(i)}; \mathbf{a}', \mathbf{W}, \mathbf{b}))^2 = \frac{1}{n} \|\mathbf{A}(\mathbf{a} - \mathbf{a}')\|^2 = \frac{1}{n} \left\| \sum_{i=1}^n \mathbf{A}_i(\mathbf{a} - \mathbf{a}') \right\|^2,$$

where $\mathbf{A}_i = \sigma((\mathbf{X} + \mathbf{\Delta})\mathbf{w}_i + \mathbf{1}_n b_i)$ is the i th column of \mathbf{A} . We are going to choose \mathbf{a}' from \mathcal{S}_Θ . To that end, define

$$\tilde{\mathbf{A}}_i = \text{sign}(\mathbf{a}_i) \mathbf{A}_i, \quad \tilde{k}_i = \begin{cases} k_i, & \text{sign}(\mathbf{a}_i) \geq 0 \\ k'_i, & \text{sign}(\mathbf{a}_i) < 0 \end{cases}$$

Further, we will choose $k'_i = 0$ if $\text{sign}(\mathbf{a}_i) \geq 0$ and $k_i = 0$ otherwise. Therefore, we have $\sum_{i=1}^N \tilde{k}_i = k$. By Mauery's sparsification lemma [XLS⁺24, Lemma 13], there exist $\tilde{k}_i \geq 0$ with $\sum_{i=1}^N \tilde{k}_i = k$ such that

$$\left\| \sum_{i=1}^N |\mathbf{a}_i| \tilde{\mathbf{A}}_i - \frac{r_a}{k} \sum_{i=1}^N \tilde{k}_i \tilde{\mathbf{A}}_i \right\|^2 \leq \frac{r_a^2 b^2}{k},$$

where $\|\mathbf{A}_i\| \leq b$ for all i . Consequently, given \mathbf{a} , we have constructed $\mathbf{a}' \in \mathcal{S}_\Theta$ such that

$$\frac{1}{n} \left\| \sum_{i=1}^n \mathbf{A}_i(\mathbf{a} - \mathbf{a}') \right\|^2 \leq \frac{r_a^2 b^2}{nk}.$$

548 Next, we provide a bound on b . By the assumptions on σ , we have

$$\begin{aligned} \|\mathbf{A}_i\|^2 &\lesssim C_q L_\sigma^2 \left(\|\mathbf{X}\mathbf{w}_i\|_{2q}^{2q} + \|\mathbf{\Delta}\mathbf{w}_i\|_{2q}^{2q} + n b_i^{2q} + \|\mathbf{X}\mathbf{w}_i\|^2 + \|\mathbf{\Delta}\|^2 + n b_i^2 \right) \\ &\lesssim n C_q L_\sigma^2 \left(T_{\mathbf{W}, \mathbf{X}}^{(q)} + r_w^{2q} \epsilon^{2q} + r_b^{2q} + T_{\mathbf{W}, \mathbf{X}}^{(2)} + r_w^2 \epsilon^2 + r_b^2 \right). \end{aligned}$$

Consequently, we can choose

$$k = \left\lceil \frac{C_q L_\sigma^2 r_a^2 \left(T_{\mathbf{W}, \mathbf{X}}^{(q)} + r_w^{2q} \epsilon^{2q} + r_b^{2q} + T_{\mathbf{W}, \mathbf{X}}^{(2)} + r_w^2 \epsilon^2 + r_b^2 \right)}{\epsilon^2} \right\rceil.$$

Finally, we need to count $|\mathcal{S}_\Theta|$. Note that

$$|\mathcal{S}_\Theta| = \binom{2N + k - 1}{k} \leq \left(\frac{e(2N + k - 1)}{k} \right)^k \leq (3eN)^k,$$

549 which concludes the proof. \square

550 We can now turn the above covering number into Rademacher complexity via a chaining argument,
551 as follows.

Lemma 15. Let $\mathfrak{R}(\mathcal{L}(\mathbf{W}, \mathbf{b}))$ denote the Rademacher complexity of the class of adversarial loss functions $\mathcal{L}(\mathbf{W}, \mathbf{b})$, defined via

$$\mathfrak{R}(\mathcal{L}(\mathbf{W}, \mathbf{b})) := \mathbb{E} \left[\sup_{\mathbf{a} \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \xi_i \max_{\|\delta^{(i)}\| \leq \varepsilon} \ell(f(\mathbf{x}^{(i)} + \delta^{(i)}; \mathbf{a}, \mathbf{W}, \mathbf{b}), y^{(i)}) \right| \right],$$

where ξ_i are i.i.d. Rademacher random variables. For simplicity, assume $C_\ell, r_a \gtrsim 1$. Then we have

$$\mathfrak{R}(\mathcal{L}(\mathbf{W}, \mathbf{b})) \lesssim \frac{C_\ell C_q L_\sigma r_a \log n \log N \left(\mathbb{E} \left[\sqrt{T_{\mathbf{W}, \mathbf{X}}^{(q)}} \right] + r_w^q \varepsilon^q + r_b^q + \mathbb{E} \left[\sqrt{T_{\mathbf{W}, \mathbf{X}}^{(2)}} \right] + r_w \varepsilon + r_b \right)}{\sqrt{n}}.$$

Proof. Let $\mathfrak{R}_n(\mathcal{L}(\mathbf{W}, \mathbf{b}))$ denote the empirical Rademacher complexity by

$$\mathfrak{R}_n(\mathcal{L}(\mathbf{W}, \mathbf{b})) := \mathbb{E}_\xi \left[\sup_{\mathbf{a} \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \xi_i \max_{\|\delta^{(i)}\| \leq \varepsilon} \ell(f(\mathbf{x}^{(i)} + \delta^{(i)}; \mathbf{a}, \mathbf{W}, \mathbf{b}), y^{(i)}) \right| \right],$$

where the expectation is only taken w.r.t. the randomness of ξ and is conditional on the training set. For simplicity, define

$$B := C_q L_\sigma \left(\sqrt{T_{\mathbf{W}, \mathbf{X}}^{(q)}} + r_w^q \varepsilon^q + r_b^q + \sqrt{T_{\mathbf{W}, \mathbf{X}}^{(2)}} + r_w \varepsilon^2 + r_b \right).$$

552 Then, by a standard chaining argument, we have for all $\alpha > 0$,

$$\begin{aligned} \mathfrak{R}_n(\mathcal{L}(\mathbf{W}, \mathbf{b})) &\lesssim \alpha + \int_{\varepsilon=\alpha}^{\infty} \sqrt{\frac{\log \mathcal{C}(\mathcal{L}, d_{\mathcal{L}}, \varepsilon)}{n}} d\varepsilon \\ &\lesssim \alpha + \frac{C_\ell r_a B \log N}{\sqrt{n}} \log \left(\frac{1}{\alpha} \right). \end{aligned}$$

By choosing $\alpha = 1/\sqrt{n}$, we obtain

$$\mathfrak{R}_n(\mathcal{L}(\mathbf{W}, \mathbf{b})) \lesssim \frac{C_\ell r_a B \log n \log N}{\sqrt{n}}.$$

553 Taking expectations with respect to the input distribution completes the proof. \square

554 Note that it remains to provide an upper bound for $T_{\mathbf{W}, \mathbf{x}}^{(q)}$ introduced in Lemma 14. This is achieved
555 by the following lemma.

Lemma 16. Suppose $\|\mathbf{w}_i\| \leq r_w$. Then, for all $q > 0$ and $N > e$, we have

$$\mathbb{E} \left[\max_{1 \leq j \leq N} \frac{1}{n} \sum_{i=1}^n \langle \mathbf{w}_j, \mathbf{x}^{(i)} \rangle^{2q} \right] \leq C_q r_w^{2q} (\log N)^q,$$

556 where C_q is a constant depending only on q .

Proof. For conciseness, let $Z_j := \frac{1}{n} \sum_{i=1}^n \langle \mathbf{w}_j, \mathbf{x}^{(i)} \rangle^{2q}$. By non-negativity of Z_j and Jensen's inequality, for all $t \geq 1$ we have,

$$\mathbb{E} \left[\max_{1 \leq j \leq N} Z_j \right] = \mathbb{E} \left[\max_{1 \leq j \leq N} Z_j^t \right]^{1/t} \leq \left(\sum_{j=1}^N \mathbb{E} [Z_j^t] \right)^{1/t} \leq N^{1/t} \left(\max_{1 \leq j \leq N} \mathbb{E} [Z_j^t] \right)^{1/t}.$$

557 Further, by Jensen's inequality

$$\begin{aligned} \mathbb{E} [Z_j^t] &= \mathbb{E} \left[\left(\frac{1}{n} \sum_{i=1}^n \langle \mathbf{w}_j, \mathbf{x}^{(i)} \rangle^{2q} \right)^t \right] \\ &\leq \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \langle \mathbf{w}_j, \mathbf{x}^{(i)} \rangle^{2qt} \right] \\ &\leq (C r_w)^{2qt} (2qt)^{qt}, \end{aligned}$$

558 where $C > 0$ is a absolute constant, and we used the moment bound of subGaussian random variables
 559 along with the fact that $\langle \mathbf{w}_j, \mathbf{x} \rangle$ is a centered subGaussian random variable with subGaussian norm
 560 $\mathcal{O}(r_w)$. As a result,

$$\mathbb{E} \left[\max_{1 \leq j \leq N} Z_j \right] \leq C_q r_w^{2q} N^{1/t} t^q \lesssim C_q r_w^{2q} (\log N)^q,$$

561 where the last inequality follows by choosing $t = \log N$. □

562 As a consequence, if the loss is also bounded, we get the following high-probability concentration
 563 bound.

564 **Corollary 17.** *Suppose $|\tilde{\ell}| \leq B_\ell$ for all $\tilde{\ell} \in \mathcal{L}(\mathbf{W}, \mathbf{b})$. Then, with probability at least $1 - \delta$ we have*

$$\left| \sup_{\tilde{\ell} \in \mathcal{L}(\mathbf{W}, \mathbf{b})} \mathbb{E} [\tilde{\ell}(\mathbf{x}, y)] - \frac{1}{n} \sum_{i=1}^n \tilde{\ell}(\mathbf{x}^{(i)}, y^{(i)}) \right| \lesssim \frac{C_\ell r_a R \log n \log N + B_\ell \sqrt{\log(1/\delta)}}{\sqrt{n}},$$

where

$$R := C_q L_\sigma (r_w^q (\log^{q/2} N + \varepsilon^q) + r_b^q + r_w (\log^{1/2} N + \varepsilon) + r_b).$$

565 C.2 Applying the Generalization Bound to Squared Loss

To apply the generalization argument above to the squared loss, we bound it with a threshold τ , and define the loss family

$$\mathcal{L}_\tau(\mathbf{W}, \mathbf{b}) := \{(\mathbf{x}, y) \mapsto \max_{\|\boldsymbol{\delta}\| \leq \varepsilon} (f(\mathbf{x} + \boldsymbol{\delta}; \mathbf{a}, \mathbf{W}, \mathbf{b}) - y)^2 \wedge \tau : \mathbf{a} \in \Theta\}.$$

We similarly define AR_τ and $\widehat{\text{AR}}_\tau$. Recall that our goal is to show

$$\text{AR}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) \leq \widehat{\text{AR}}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) + \varepsilon_1(n, N, d).$$

We readily have $\widehat{\text{AR}}_\tau(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) \leq \widehat{\text{AR}}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b})$. Further, Corollary 17 yields

$$\left| \text{AR}_\tau(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) - \widehat{\text{AR}}_\tau(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) \right| \lesssim \frac{\sqrt{\tau} r_a R \log n \log N}{\sqrt{n}} + \tau \sqrt{\frac{\log(1/\delta)}{n}},$$

566 with probability at least $1 - \delta$. Thus, the remaining step is to bound $\text{AR}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b})$ and $\widehat{\text{AR}}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b})$
 567 with their clipped versions. To do so, we first provide the following tail probability estimate.

Lemma 18. *Suppose $(z_j)_{j=1}^N$ are non-negative random variables with subGaussian norm r . Then, for any $q > 0$ and $\tau \geq C_q r^q$ where C_q is a constant depending only on q , we have*

$$\mathbb{P} \left(\frac{1}{N} \sum_{j=1}^N z_j^q \geq \tau \right) \leq \exp \left(-\frac{c\tau^{2/q}}{r^2} \right),$$

568 where $c > 0$ is an absolute constant.

Proof. For any $t \geq 1$, we have the following Markov bound,

$$\mathbb{P} \left(\frac{1}{N} \sum_{j=1}^N z_j^q \geq \tau \right) = \mathbb{P} \left(\left(\frac{1}{N} \sum_{j=1}^N z_j^q \right)^t \geq \tau^t \right) \leq \frac{\mathbb{E} \left[\left(\frac{1}{N} \sum_{j=1}^N z_j^q \right)^t \right]}{\tau^t} \leq \frac{\mathbb{E} \left[\frac{1}{N} \sum_{j=1}^N z_j^{qt} \right]}{\tau^t},$$

where the last inequality follows from Jensen's inequality. Further, by subGaussianity of z_j , we have $\mathbb{E}[z_j^{qt}] \leq (Cr^2 qt)^{qt/2}$, where $C > 0$ is an absolute constant. As a result,

$$\mathbb{P} \left(\frac{1}{N} \sum_{j=1}^N z_j^q \geq \tau \right) \leq \frac{(Cr^2 qt)^{qt/2}}{\tau^t}.$$

The above bound is minimized at $t = \frac{\tau^{2/q}}{Cr^{2/q}e}$. Note that $t \geq 1$ requires $\tau \geq C_q r^q$. Plugging this choice of t in the above bound yields

$$\mathbb{P}\left(\frac{1}{N}\sum_{j=1}^N z_j^q \geq \tau\right) \leq \exp\left(-\frac{\tau^{2/q}}{2Cr^{2/q}e}\right),$$

569 which completes the proof. \square

Lemma 19. *Suppose Assumption 2 holds, in particular $|g(\mathbf{z})| \leq L_g(1 + \|\mathbf{z}\|^p)$ for some constant L_g . Let $\Theta = \{\mathbf{a} : \|\mathbf{a}\| \leq r_a/\sqrt{N}\}$, $\|\mathbf{w}_i\| \leq r_w$, and $|b_i| \leq r_b$ for all $i \in [N]$. Assume σ satisfies (B.1). Define $\varepsilon_1 := 1 \vee \varepsilon$, and let*

$$\varkappa := C_q r_a^2 L_\sigma^2 (r_w^{2q} \varepsilon_1^{2q} + r_b^{2q} + r_w^2 \varepsilon_1^2 + r_b^2) + C,$$

where C_q is a constant depending only on q and C is an absolute constant. Then, for all

$$\tau \geq C_k \left\{ \varkappa \vee L_\sigma^2 r_w^{2q} \log^q \frac{n}{\delta} \vee L_g^2 \log^p \frac{n}{\delta} \right\},$$

570 where C_k is a sufficiently large constant, we have

$$\begin{aligned} \left| \text{AR}(\mathbf{a}, \mathbf{W}, \mathbf{b}) - \widehat{\text{AR}}(\mathbf{a}, \mathbf{W}, \mathbf{b}) \right| &\leq \left| \text{AR}_\tau(\mathbf{a}, \mathbf{W}, \mathbf{b}) - \widehat{\text{AR}}_\tau(\mathbf{a}, \mathbf{W}, \mathbf{b}) \right| \\ &\quad + C\varkappa \left(\exp\left(-\Omega\left(\frac{\tau^{1/q}}{L_\sigma^{2/q} r_w^2}\right)\right) + \exp(-\Omega(\tau^{1/p})) \right), \end{aligned}$$

571 with probability at least $1 - \delta$ uniformly over all $\mathbf{a} \in \Theta$.

572 **Proof.** Since \mathbf{W} and \mathbf{b} are fixed, we use the shorthand notation $f(\mathbf{x}; \mathbf{a}) = f(\mathbf{x}; \mathbf{a}, \mathbf{W}, \mathbf{b})$.

In the first section of the proof, we will upper and lower bound $\text{AR}(\mathbf{a}, \mathbf{W}, \mathbf{b})$ with $\text{AR}_\tau(\mathbf{a}, \mathbf{W}, \mathbf{b})$. Note that the lower bound is trivial as $\text{AR}_\tau(\mathbf{a}, \mathbf{W}, \mathbf{b}) \leq \text{AR}(\mathbf{a}, \mathbf{W}, \mathbf{b})$, thus we move on to the upper bound. Let

$$\tilde{\ell}(\mathbf{x}, y) = \max_{\|\delta\| \leq \varepsilon} (f(\mathbf{x} + \delta; \mathbf{a}) - y)^2.$$

573 Then,

$$\begin{aligned} \text{AR}(\mathbf{a}, \mathbf{W}, \mathbf{b}) &= \mathbb{E}\left[\tilde{\ell}(\mathbf{x}, y)\mathbb{I}\left[\tilde{\ell}(\mathbf{x}, y) \leq \tau\right]\right] + \mathbb{E}\left[\tilde{\ell}(\mathbf{x}, y)\mathbb{I}\left[\tilde{\ell}(\mathbf{x}, y) > \tau\right]\right] \\ &\leq \text{AR}_\tau(\mathbf{a}, \mathbf{W}, \mathbf{b}) + \mathbb{E}\left[\tilde{\ell}(\mathbf{x}, y)^2\right]^{1/2} \mathbb{P}\left(\tilde{\ell}(\mathbf{x}, y) \geq \tau\right)^{1/2}. \end{aligned}$$

574 Further, we have the following upper bound for the adversarial loss,

$$\begin{aligned} \tilde{\ell}(\mathbf{x}, y) &= \max_{\|\delta\| \leq \varepsilon} (f(\mathbf{x} + \delta; \mathbf{a}) - y)^2 \\ &\lesssim \max_{\|\delta\| \leq \varepsilon} f(\mathbf{x} + \delta; \mathbf{a})^2 + y^2 \\ &\lesssim \max_{\|\delta\| \leq \varepsilon} \|\mathbf{a}\|^2 \|\sigma(\mathbf{W}(\mathbf{x} + \delta) + \mathbf{b})\|^2 + y^2 \\ &\lesssim r_a^2 C_q L_\sigma^2 \left(\frac{1}{N} \sum_{j=1}^N \langle \mathbf{w}_j, \mathbf{x} \rangle^{2q} + r_w^{2q} \varepsilon^{2q} + r_b^{2q} + \frac{1}{N} \sum_{j=1}^N \langle \mathbf{w}_j, \mathbf{x} \rangle^2 + r_w^2 \varepsilon^2 + r_b^2 \right) + y^2 \end{aligned}$$

575 Moreover, by Jensen's inequality,

$$\begin{aligned} \mathbb{E}\left[\left(\frac{1}{N}\sum_{j=1}^N \langle \mathbf{w}_j, \mathbf{x} \rangle^{2q}\right)^2\right] &\leq \mathbb{E}\left[\frac{1}{N}\sum_{j=1}^N \langle \mathbf{w}_j, \mathbf{x} \rangle^{4q}\right] \\ &\leq (Cr_w)^{4q} (4q)^{2q} \leq C_q r_w^{4q} \end{aligned}$$

for all $q > 0$, where C is an absolute constant and we used the subGaussianity of $\langle \mathbf{w}_j, \mathbf{x} \rangle$ to bound its moment. As a result,

$$\mathbb{E} \left[\tilde{\ell}(\mathbf{x}, y)^2 \right]^{1/2} \lesssim r_a^2 C_q L_\sigma^2 (r_w^{2q} (1 + \varepsilon^{2q}) + r_b^{2q} + r_w^2 (1 + \varepsilon^2) + r_b^2) + \mathbb{E} [y^4]^{1/2}.$$

576 By assumption 2, we have $\mathbb{E} [y^4]^{1/2} \leq C$ for some absolute constant C .

577 To estimate the tail probability of $\tilde{\ell}(\mathbf{x}, y)$. Using the assumption on τ and the upper bound on $\tilde{\ell}(\mathbf{x}, y)$
578 developed above, using a union bound we have

$$\begin{aligned} \mathbb{P} \left(\tilde{\ell}(\mathbf{x}, y) \geq \tau \right) &\leq \mathbb{P} \left(\frac{L_\sigma^2}{N} \sum_{j=1}^N \langle \mathbf{w}_j, \mathbf{x} \rangle^{2q} + \frac{L_\sigma^2}{N} \sum_{j=1}^N \langle \mathbf{w}_j, \mathbf{x} \rangle^2 + y^2 \geq \frac{\tau}{2} \right) \\ &\leq \mathbb{P} \left(\frac{L_\sigma^2}{N} \sum_{j=1}^N \langle \mathbf{w}_j, \mathbf{x} \rangle^{2q} \geq \frac{\tau}{6} \right) + \mathbb{P} \left(\frac{L_\sigma^2}{N} \sum_{j=1}^N \langle \mathbf{w}_j, \mathbf{x} \rangle^2 \geq \frac{\tau}{6} \right) + \mathbb{P} \left(y^2 \geq \frac{\tau}{6} \right) \\ &\leq 2 \exp \left(\frac{-c\tau^{1/q}}{L_\sigma^{2/q} r_w^2} \right) + \mathbb{P} (y^2 \geq \tau), \end{aligned}$$

where we used Lemma 18, the fact that $|\langle \mathbf{w}_j, \mathbf{x} \rangle|$ is subGaussian with norm $\mathcal{O}(r_w)$, and that $q \geq 1$. Furthermore, using the moment estimate on y in Assumption 2 along with the technique developed in Lemma 18, we have

$$\mathbb{P} \left(y^2 \geq \frac{\tau}{6} \right) \leq \exp \left(\frac{-c\tau^{1/p}}{L_g^{2/p}} \right),$$

579 for $\tau \geq C$, where $C, c > 0$ are absolute constants.

As a result, we obtain

$$\text{AR}(\mathbf{a}, \mathbf{W}, \mathbf{b}) - \text{AR}_\tau(\mathbf{a}, \mathbf{W}, \mathbf{b}) \lesssim \varkappa \left(\exp \left(-\frac{c\tau^{1/q}}{L_\sigma^{2/q} r_w^2} \right) + \exp(-c\tau^{1/p}) \right),$$

580 for all $\mathbf{a} \in \Theta$.

581 In the next part of the proof, we will show that with probability at least $1 - \delta$, we have $\widehat{\text{AR}}(\mathbf{a}, \mathbf{W}, \mathbf{b}) =$
582 $\widehat{\text{AR}}_\tau(\mathbf{a}, \mathbf{W}, \mathbf{b})$ uniformly over all \mathbf{a} . Note that this is equivalent to asking $\tilde{\ell}(\mathbf{x}^{(i)}, y^{(i)}) \leq \tau$ for all
583 $1 \leq i \leq n$. For any fixed i , using the upper bound on $\tilde{\ell}(\mathbf{x}^{(i)}, y^{(i)})$, we have

$$\begin{aligned} \mathbb{P} \left(\tilde{\ell}(\mathbf{x}^{(i)}, y^{(i)}) \geq \tau \right) &\leq \mathbb{P} \left(\frac{L_\sigma^2}{N} \sum_{j=1}^N \langle \mathbf{w}_j, \mathbf{x} \rangle^{2q} + \frac{L_\sigma^2}{N} \sum_{j=1}^N \langle \mathbf{w}_j, \mathbf{x} \rangle^2 + y^2 \geq \frac{\tau}{2} \right) \\ &\lesssim \exp \left(\frac{-c\tau^{1/q}}{L_\sigma^{2/q} r_w^2} \right) + \exp \left(\frac{-c_k \tau^{1/p}}{L_g^{2/p}} \right). \end{aligned}$$

Consequently, by a union bound we have

$$\mathbb{P} \left(\max_{1 \leq i \leq n} \tilde{\ell}(\mathbf{x}^{(i)}, y^{(i)}) \geq \tau \right) \leq n \left(\exp \left(\frac{-c\tau^{1/q}}{L_\sigma^{2/q} r_w^2} \right) + \exp \left(\frac{-c_k \tau^{1/p}}{L_g^{2/p}} \right) \right).$$

Choosing

$$\tau \geq C_k \left\{ L_\sigma^2 r_w^{2q} \log^q \frac{n}{\delta} \vee L_g^2 \log^p \frac{n}{\delta} \right\}$$

584 with a sufficiently large constant C_k ensures the above probability is at most δ , finishing the proof.

585 □

586 We are now ready to present the main result of this section.

Proposition 20. Suppose Assumption 2 holds and σ satisfies (B.1), $\Theta = \{\mathbf{a} : \|\mathbf{a}\| \leq r_a/\sqrt{N}\}$, $\|\mathbf{w}_i\| \leq 1$, and $|b_i| \leq r_b$ for all $1 \leq i \leq N$. Let

$$\varkappa := C_q r_a^2 L_\sigma^2 (1 + \varepsilon^{2q} + r_b^{2q}) + C_{p,k} L_g^2,$$

where C_q and C_p are constants depending only on q and p respectively. Then we have

$$\text{AR}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) - \min_{\mathbf{a} \in \Theta} \text{AR}(\mathbf{a}, \mathbf{W}, \mathbf{b}) \leq \tilde{\mathcal{O}}\left(\frac{\varkappa}{\sqrt{n}}\right),$$

587 with probability at least $1 - \mathcal{O}(n^{-c})$ for some constant $c > 0$.

Proof. We can summarize the generalization bound of Corollary 17 as

$$\text{AR}_\tau(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) - \widehat{\text{AR}}_\tau(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) \lesssim \sqrt{\frac{\tau \varkappa}{n}} + \tau \sqrt{\frac{\log(1/\delta)}{n}},$$

where

$$\varkappa := C_q r_a^2 L_\sigma^2 (1 + \varepsilon^{2q} + r_b^{2q}) + C_{p,k} L_g^2,$$

is obtained from Lemma 19 by letting $r_w = 1$. Further, we use the fact that $\widehat{\text{AR}}_\tau(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) \leq \text{AR}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b})$, and

$$\text{AR}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) - \text{AR}_\tau(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) \lesssim \varkappa e^{-\Omega(\tau/\varkappa)}$$

from Lemma 19 to arrive that

$$\text{AR}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) - \widehat{\text{AR}}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) \leq \tilde{\mathcal{O}}\left(\sqrt{\frac{\tau \varkappa}{n}} + \tau \sqrt{\frac{\log(1/\delta)}{n}} + \varkappa e^{-\Omega\left(\frac{\tau^{1/q}}{L_\sigma^{2/q}}\right)} + \varkappa e^{-\Omega\left(\frac{\tau^{1/p}}{L_g^{2/p}}\right)}\right).$$

Note that $\varkappa \geq L_\sigma^2 \vee L_g^2$. Choosing $\tau = C \varkappa \log^{p \vee q}(\varkappa n / \delta)$ with a sufficiently large absolute constant $C > 0$ satisfies the assumption of Lemma 19, and further letting $\delta = n^{-c}$ for some constant $c > 0$, we obtain

$$\text{AR}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) - \widehat{\text{AR}}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) \leq \tilde{\mathcal{O}}\left(\frac{\varkappa}{\sqrt{n}}\right),$$

588 which holds with probability at least $1 - n^{-c}$ over the randomness of the training set.

Recall $\mathbf{a}^* = \arg \min_{\mathbf{a} \in \Theta} \text{AR}(\mathbf{a}, \mathbf{W}, \mathbf{b})$. Similarly, Lemma 19 guarantees

$$\widehat{\text{AR}}(\mathbf{a}^*, \mathbf{W}, \mathbf{b}) - \text{AR}(\mathbf{a}^*, \mathbf{W}, \mathbf{b}) \leq \tilde{\mathcal{O}}\left(\frac{\varkappa}{\sqrt{n}}\right),$$

589 on the same event as above. Finally, we have $\widehat{\text{AR}}(\hat{\mathbf{a}}, \mathbf{W}, \mathbf{b}) \leq \widehat{\text{AR}}(\mathbf{a}^*, \mathbf{W}, \mathbf{b})$ by definition of $\hat{\mathbf{a}}$,
590 which concludes the proof of the proposition. \square

591 D Approximation Analysis

592 Let $\Pi_U \mathbf{w} = \frac{\mathbf{U}^\top \mathbf{U} \mathbf{w}}{\|\mathbf{U} \mathbf{w}\|}$ denote the projection of $\mathbf{w} \in \mathbb{S}^{d-1}$ onto $\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k) \cap \mathbb{S}^{d-1}$ (if $\|\mathbf{U} \mathbf{w}\| =$
593 0 we can simply let $\Pi_U \mathbf{w} = \mathbf{u}_1$). Suppose $\langle \mathbf{w}, \mathbf{u} \rangle \geq 1 - \zeta$ for some $\zeta \in (0, 1)$ and $\mathbf{u} \in$
594 $\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$ with $\|\mathbf{u}\| = 1$. Then, we have the following properties for this projection:

- 595 $\langle \Pi_U \mathbf{w}, \mathbf{u} \rangle \geq 1 - \zeta,$
- 596 $\|\mathbf{w} - \Pi_U \mathbf{w}\| \leq \sqrt{2\zeta}.$

Let $h : \mathbb{R}^k \rightarrow \mathbb{R}$ be the function constructed in the proof of Theorem 1. Then,

$$\text{AR}^* = \mathbb{E} \left[\max_{\|\boldsymbol{\delta}\| \leq \varepsilon} (h(\mathbf{U}(\mathbf{x} + \boldsymbol{\delta})) - y)^2 \right].$$

597 Let us denote $f(\mathbf{x}) = f(\mathbf{x}; \mathbf{a}^*, \mathbf{W}, \mathbf{b})$ for consciences. Then,

$$\begin{aligned} \text{AR}(\mathbf{a}^*, \mathbf{W}, \mathbf{b}) - \text{AR}^* &= \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} (f(\mathbf{x} + \delta) - y)^2 - \max_{\|\delta\| \leq \varepsilon} (h(\mathbf{U}(\mathbf{x} + \delta)) - y)^2 \right] \\ &\leq \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \left\{ (f(\mathbf{x} + \delta) - y)^2 - (h(\mathbf{U}(\mathbf{x} + \delta)) - y)^2 \right\} \right] \\ &= \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \left\{ (f(\mathbf{x} + \delta) - h(\mathbf{U}(\mathbf{x} + \delta))) \underbrace{(f(\mathbf{x} + \delta) + h(\mathbf{U}(\mathbf{x} + \delta)) - 2y)}_{=: \mathcal{Z}} \right\} \right] \end{aligned}$$

Let $\Pi_U \mathbf{W} = (\Pi_U \mathbf{w}_1, \dots, \Pi_U \mathbf{w}_N)^\top$. Then, we have the decompositions

$$f(\mathbf{x} + \delta; \mathbf{a}^*, \mathbf{W}, \mathbf{b}) = f(\mathbf{x} + \delta; \mathbf{a}^*, \mathbf{W}, \mathbf{b}) - f(\mathbf{x} + \delta; \mathbf{a}^*, \Pi_U \mathbf{W}, \mathbf{b}) + f(\mathbf{x} + \delta; \mathbf{a}^*, \Pi_U \mathbf{W}, \mathbf{b}),$$

598 and

$$\begin{aligned} \mathcal{Z} &= f(\mathbf{x} + \delta; \mathbf{a}^*, \mathbf{W}, \mathbf{b}) - f(\mathbf{x} + \delta; \mathbf{a}^*, \Pi_U \mathbf{W}, \mathbf{b}) + f(\mathbf{x} + \delta; \mathbf{a}^*, \Pi_U \mathbf{W}, \mathbf{b}) - h(\mathbf{U}(\mathbf{x} + \delta)) \\ &\quad + 2h(\mathbf{U}(\mathbf{x} + \delta)) - 2y. \end{aligned}$$

599 Plugging this decomposition into the above and using the Cauchy-Schwartz inequality yields

$$\text{AR}(\mathbf{a}^*, \mathbf{W}, \mathbf{b}) - \text{AR}^* \leq (\sqrt{\mathcal{E}_1} + \sqrt{\mathcal{E}_2})^2 + \sqrt{\mathcal{E}_3(\mathcal{E}_1 + \mathcal{E}_2)}, \quad (\text{D.1})$$

600 where

$$\mathcal{E}_1 := \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} (f(\mathbf{x} + \delta; \Pi_U \mathbf{W}, \mathbf{b}) - h(\mathbf{U}(\mathbf{x} + \delta)))^2 \right], \quad (\text{D.2})$$

$$\mathcal{E}_2 := \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} (f(\mathbf{x} + \delta; \Pi_U \mathbf{W}, \mathbf{b}) - f(\mathbf{x} + \delta; \mathbf{a}^*, \mathbf{W}, \mathbf{b}))^2 \right], \quad (\text{D.3})$$

$$\mathcal{E}_3 := 4 \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} (h(\mathbf{U}(\mathbf{x} + \delta)) - y)^2 \right] = 4\text{AR}^*. \quad (\text{D.4})$$

601 Under Definition 3, we have a set of good neurons S to work with. To continue, we introduce a
602 similar subset of good neurons under Definition 2.

Definition 21. Suppose the weights $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_N)^\top$ are obtained from the α -DFL oracle of Definition 2. Fix a maximal $2\sqrt{2\zeta}$ -packing of \mathbb{S}^{k-1} with respect to the Euclidean norm, denoted by $(\bar{\mathbf{v}}_i)_{i=1}^M$. Define $\mathbf{v}_j := \frac{\mathbf{U}\mathbf{w}_j}{\|\mathbf{U}\mathbf{w}_j\|}$ for all $j \in [N]$, and

$$S_i := \{j \in [N] : \|\mathbf{v}_j - \bar{\mathbf{v}}_i\| \leq \sqrt{2\zeta}\},$$

603 for all $i \in [M]$. Note that (S_i) are mutually exclusive. Define $S := \bigcup_{i=1}^M S_i$. Note that there
604 are constants $c_k, C_k > 0$ such that $c_k(1/\zeta)^{(k-1)/2} \leq M \leq C_k(1/\zeta)^{(k-1)/2}$. Therefore, using
605 Definition 2, we have $|S|/N \geq \Omega(\alpha)$.

606 Note that when considering the (α, β) -SFL oracle, we leave S unchanged from Definition 3. In either
607 case, for every $j \notin S$, we will choose $a_j^* = 0$. Then, we then have the following upper bound on \mathcal{E}_2 .

Lemma 22. Suppose $a_j^* = 0$ for $j \notin S$ and $\|\mathbf{a}^*\| \leq \tilde{r}_\alpha / \sqrt{|S|}$. Then,

$$\mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} (f(\mathbf{x} + \delta; \mathbf{a}^*, \Pi_U \mathbf{W}, \mathbf{b}) - f(\mathbf{x} + \delta; \mathbf{a}^*, \mathbf{W}, \mathbf{b}))^2 \right] \lesssim L_\sigma^2 C_q \tilde{r}_\alpha^2 (1 + r_b^{2(q-1)} + \varepsilon^{2(q-1)})(1 + \varepsilon^2)\zeta,$$

608 where C_q is a constant only depending on q .

609 **Proof.** To be concise, we define $\tilde{\mathbf{x}}_\delta := \mathbf{x} + \delta$ and hide dependence on \mathbf{a}^* and \mathbf{b} in the following
610 notation. By pseudo-Lipschitzness of σ and the Cauchy-Schwartz inequality,

$$\begin{aligned} f(\tilde{\mathbf{x}}_\delta; \Pi_U \mathbf{W}) - f(\tilde{\mathbf{x}}_\delta; \mathbf{W}) &= \sum_{j \in S} a_j^* (\sigma(\langle \Pi_U \mathbf{w}_j, \tilde{\mathbf{x}}_\delta \rangle + b_j) - \sigma(\langle \mathbf{w}_j, \tilde{\mathbf{x}}_\delta \rangle + b_j)) \\ &\leq L_\sigma \sum_{j \in S} |a_j^*| (|\langle \Pi_U \mathbf{w}_j, \tilde{\mathbf{x}}_\delta \rangle + b_j|^{q-1} + |\langle \mathbf{w}_j, \tilde{\mathbf{x}}_\delta \rangle + b_j|^{q-1} + 1) |\langle \Pi_U \mathbf{w}_j - \mathbf{w}_j, \tilde{\mathbf{x}}_\delta \rangle|. \end{aligned}$$

Let

$$\mathcal{A}_j := |\langle \Pi_U \mathbf{w}_j, \tilde{\mathbf{x}}_\delta \rangle + b_j|^{q-1} + |\langle \mathbf{w}_j, \tilde{\mathbf{x}}_\delta \rangle + b_j|^{q-1} + 1,$$

and

$$\mathcal{B}_j := |\langle \Pi_U \mathbf{w}_j - \mathbf{w}_j, \tilde{\mathbf{x}}_\delta \rangle|.$$

611 Then,

$$\begin{aligned} \mathcal{E}_2 &\leq L_\sigma^2 \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \left(\sum_{j \in S} |a_j^*| \mathcal{A}_j \mathcal{B}_j \right)^2 \right] \leq \frac{L_\sigma^2 \tilde{r}_a^2}{|S|} \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \sum_{j \in S} \mathcal{A}_j^2 \mathcal{B}_j^2 \right] \\ &\leq \frac{L_\sigma^2 \tilde{r}_a^2}{|S|} \sum_{j \in S} \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \mathcal{A}_j^2 \mathcal{B}_j^2 \right] \\ &\leq \frac{L_\sigma^2 \tilde{r}_a^2}{|S|} \sum_{j \in S} \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \mathcal{A}_j^4 \right]^{1/2} \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \mathcal{B}_j^4 \right]^{1/2}. \end{aligned}$$

Additionally, we have

$$\max_{\|\delta\| \leq \varepsilon} \mathcal{A}_j \leq C_q \left(|\langle \Pi_U \mathbf{w}_j, \mathbf{x} \rangle|^{q-1} + |\langle \mathbf{w}_j, \mathbf{x} \rangle|^{q-1} + \varepsilon^{q-1} + r_b^{q-1} + 1 \right),$$

and

$$\max_{\|\delta\| \leq \varepsilon} \mathcal{B}_j \leq \varepsilon \|\Pi_U \mathbf{w}_j - \mathbf{w}_j\| + |\langle \Pi_U \mathbf{w}_j - \mathbf{w}_j, \mathbf{x} \rangle|.$$

612 Further, by Assumption 2, for all $\mathbf{v} \in \mathbb{R}^d$, $\langle \mathbf{v}, \mathbf{x} \rangle$ is a centered subGaussian random variable with
613 subGaussian norm $\mathcal{O}(\|\mathbf{v}\|)$, therefore $\mathbb{E}[|\langle \mathbf{v}, \mathbf{x} \rangle|^q] \leq C_q \|\mathbf{v}\|^q$ for all $q > 0$. In summary,

$$\mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \mathcal{A}_j^4 \right]^{1/2} \leq C_q (1 + r_b^{2(q-1)} + \varepsilon_1^{2(q-1)}), \quad \text{and} \quad \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \mathcal{B}_j^4 \right]^{1/2} \lesssim (1 + \varepsilon^2) \zeta,$$

614 where we used the fact that $\|\Pi_U \mathbf{w}_j - \mathbf{w}_j\|^2 \leq 2\zeta$ for all $j \in S$. This completes the proof. \square

615 While the term \mathcal{E}_1 defined in (D.2) is an expectation over the entire distribution of \mathbf{x} , most approxi-
616 mation bounds support only a compact subset of \mathbb{R}^d . The following lemma shows that approximation
617 on compact sets is sufficient to bound \mathcal{E}_1 .

Lemma 23. Suppose $a_j^* = 0$ for $j \notin S$ and $\|\mathbf{a}^*\| \leq \tilde{r}_a / \sqrt{|S|}$. Further, suppose $r_z \geq 1 \vee 2\varepsilon$. Let

$$\epsilon_{\text{approx}} := \sup_{\|\mathbf{U}\mathbf{x}\| \leq r_z} |f(\mathbf{x}; \mathbf{a}^*, \Pi_U \mathbf{W}, \mathbf{b}) - h(\mathbf{U}(\mathbf{x} + \delta))|.$$

Assume h satisfies $|h(\mathbf{z})| \leq L_h(1 + \|\mathbf{z}\|^p)$ for all $\mathbf{z} \in \mathbb{R}^k$ and some constant $p \geq 0$. Then,

$$\mathcal{E}_1 \leq \epsilon_{\text{approx}}^2 + \left(L_\sigma^2 C_q \tilde{r}_a^2 (1 + \varepsilon^{2q} + r_b^{2q}) + L_h^2 C_{p,k} (1 + \varepsilon^{2p}) \right) e^{-\Omega(r_z^2)}.$$

Proof. For brevity, define

$$\Delta_\delta := (f(\tilde{\mathbf{x}}_\delta; \mathbf{a}^*, \Pi_U \mathbf{W}, \mathbf{b}) - h(\mathbf{U}(\mathbf{x} + \delta)))^2$$

618 where $\tilde{\mathbf{x}}_\delta := \mathbf{x} + \delta$. Then,

$$\begin{aligned} \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \Delta_\delta \right] &\leq \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \Delta_\delta \mathbb{I}[\|\mathbf{U}\tilde{\mathbf{x}}_\delta\| \leq r_z] \right] + \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \Delta_\delta \mathbb{I}[\|\mathbf{U}\tilde{\mathbf{x}}_\delta\| > r_z] \right] \\ &\leq \epsilon_{\text{approx}}^2 + \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \Delta_\delta^2 \right]^{1/2} \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \mathbb{I}[\|\mathbf{U}\tilde{\mathbf{x}}_\delta\| > r_z] \right]^{1/2} \\ &\leq \epsilon_{\text{approx}}^2 + \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \Delta_\delta^2 \right]^{1/2} \mathbb{P}(\|\mathbf{U}\mathbf{x}\| > r_z - \varepsilon)^{1/2} \\ &\leq \epsilon_{\text{approx}}^2 + \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \Delta_\delta^2 \right]^{1/2} \mathbb{P}\left(\|\mathbf{U}\mathbf{x}\| > \frac{r_z}{2}\right)^{1/2}. \end{aligned}$$

Furthermore, we have

$$\mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \Delta_{\delta}^2 \right] \lesssim \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} f(\tilde{\mathbf{x}}_{\delta}; \mathbf{a}^*, \Pi_{\mathcal{U}} \mathbf{W}, \mathbf{b})^4 \right] + \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} h(\mathbf{U}(\mathbf{x} + \delta))^4 \right].$$

619 Recall the notation $\mathbf{v}_j := \frac{\mathbf{U} \mathbf{w}_j}{\|\mathbf{U} \mathbf{w}_j\|}$ and $\mathbf{z} := \mathbf{U} \mathbf{x}$. Then, by Cauchy-Schwartz and Jensen inequalities,

$$\begin{aligned} \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} f(\tilde{\mathbf{x}}_{\delta}; \mathbf{a}^*, \Pi_{\mathcal{U}} \mathbf{W}, \mathbf{b})^4 \right] &\leq \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \|\mathbf{a}\|^4 \|\sigma(\Pi_{\mathcal{U}} \mathbf{W}(\mathbf{x} + \delta) + \mathbf{b})\|^4 \right] \\ &\leq \frac{\tilde{r}_a^4}{|S|} \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \left(\sum_{j \in S} \sigma(\langle \mathbf{v}_j, \mathbf{z} + \mathbf{U} \delta \rangle + b_j)^4 \right) \right] \\ &\leq \frac{\tilde{r}_a^4 L_{\sigma}^4 C_q}{|S|} \mathbb{E} \left[\sum_{j \in S} \langle \mathbf{v}_j, \mathbf{z} \rangle^{4q} + \varepsilon^{4q} + r_b^{4q} \right] \\ &\leq C_q L_{\sigma}^4 \tilde{r}_a^4 (1 + \varepsilon^{4q} + r_b^{4q}). \end{aligned}$$

Similarly we can prove

$$\mathbb{E} \left[\max_{\|\delta\|} h(\mathbf{U}(\mathbf{x} + \delta))^4 \right] \leq C_{p,k} L_h^4 (1 + \varepsilon^{4p}).$$

In summary,

$$\mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} \Delta_{\delta}^2 \right]^{1/2} \lesssim C_q L_{\sigma}^2 \tilde{r}_a^2 (1 + \varepsilon^{2q} + r_b^{2q}) + C_{p,k} L_h^2 (1 + \varepsilon^{2p}).$$

Finally, the probability bound

$$\mathbb{P} \left(\|\mathbf{U} \mathbf{x}\| \geq \frac{r_z}{4} \right) \leq e^{-\Omega(r_z^2)}$$

620 follows from subGaussianity of \mathbf{x} and the fact that $k = \mathcal{O}(1)$. □

621 D.1 Approximating Univariate Functions

622 In this section, we recall prior results on approximating univariate functions with random biases in
623 the infinite-width regime under ReLU and polynomial activations.

Lemma 24 ([MHPG⁺23, Lemma 21, Adapted]). *Let σ be the ReLU activation and $b \sim \text{Unif}(-r_b, r_b)$. Let $\tilde{h} : \mathbb{R} \rightarrow \mathbb{R}$ be a smooth function such that $\tilde{h}(z) = h(z)$ for all $|z| \leq r_b/2$, and $\tilde{h}(-r_b) = \tilde{h}'(-r_b) = 0$. Then, for all $|z| \leq r_b/2$ we have*

$$\mathbb{E}_b \left[2r_b \tilde{h}''(-b) \sigma(z + b) \right] = h(z).$$

624 **Proof.** The proof follows from integration by parts, namely

$$\begin{aligned} \mathbb{E}_b \left[2r_b \tilde{h}''(-b) \sigma(z + b) \right] &= \int_{-z}^{r_b} \tilde{h}''(-b) (z + b) db \\ &= -\tilde{h}'(-r_b) (z + r_b) + \int_{-r_b}^z \tilde{h}'(u) du \\ &= -\tilde{h}'(-r_b) (z + r_b) + \tilde{h}(z) - \tilde{h}(-r_b) = h(z). \end{aligned}$$

625 □

626 Furthermore, we have the following result for infinite-width approximation with polynomial activa-
627 tions.

Lemma 25 ([OSSW24, Lemma 30, Adapted]). *Let σ be a polynomial of degree q and suppose $b \sim \text{Unif}(-r_b, r_b)$ and h is a polynomial of degree p such that $q \geq p$, and in particular satisfies $|h(z)| \leq L_h(1 + |z|^p)$. Suppose $r_b \geq q$. Then, there exists a function $f : [-r_b, r_b] \rightarrow \mathbb{R}$ such that*

$$\mathbb{E}_b [2r_b f(b) \sigma(z + b)] = h(z), \quad \forall z \in \mathbb{R}.$$

628 Furthermore, we have $|f(z)| \leq C_{\sigma,h}$ for all z , where $C_{\sigma,h}$ only depends on the activation and L_h .

Proof. In order for σ to approximate arbitrary polynomials of degree at most q , it is sufficient to show that σ can approximate at least one polynomial per degree, ranging from degree 0 to q . Defining the corresponding polynomial with degree i as $g_i(z)$, then h will be in the span of $\{g_i\}_{i=0}^q$. More specifically, suppose $h(z) = \sum_{j=0}^p \alpha_j z^j$, and $g_i(z) = \sum_{j=0}^i \gamma_{i,j} z^j$. Then there exist $\{\beta_i\}_{i=0}^q$ such that

$$\sum_{i=1}^p \beta_i g_i(z) = \sum_{j=0}^p \sum_{i=j}^p \gamma_{i,j} \beta_i z^j = \sum_{j=0}^p \alpha_j z^j.$$

Indeed, we can let $\beta_i = 0$ for all $i > p$. Additionally, note that $\gamma_{i,i} \neq 0$ for all $i \leq q$ by definition. Therefore, the solution to the above equation is given iteratively by $\beta_p = \alpha_p / \gamma_{p,p}$ and

$$\beta_{p-j} = \frac{\alpha_{p-j} - \sum_{i=0}^{j-1} \gamma_{p-i,p-j} \beta_{p-i}}{\gamma_{p-j,p-j}},$$

629 for $1 \leq j \leq p$. Importantly, $|\beta_i|$ for all i can be bounded polynomially by $\{\alpha_j\}_j$, $\{\gamma_{i,j}\}_{i,j}$ and
 630 $\{\gamma_{i,i}^{-1}\}_i$. Further, $|\alpha_i|$ can be bounded polynomially by L_h for all i . Thus, it remains to construct
 631 $\{g_i\}$.

Following [OSSW24], we define

$$g_q(z) = \int_{-q}^0 \sigma(z+b) db.$$

It is straightforward to verify that g_q has degree (exactly) q . We then iteratively define

$$g_{q-i}(z) = g_{q-(i-1)}(z+1) - g_{q-(i-1)}(z), \quad \forall 1 \leq i \leq q.$$

Using the definition above and by induction, one can verify g_i has degree exactly i . Furthermore, expanding the definition above yields

$$g_{q-i}(z) = \sum_{j=0}^i c_{i,j} g_q(z+j) = \sum_{j=0}^i c_{i,j} \int_{-q}^0 \sigma(z+b+j) db,$$

where $c_{i,j} = (-1)^{i-j} \binom{i}{j}$, i.e. the coefficients that satisfy $(z-1)^i = \sum_{j=0}^i c_{i,j} z^j$. In particular, we can write

$$g_{q-i}(z) = \sum_{j=0}^i c_{i,j} \int_{-q+j}^j \sigma(z+b) db = \mathbb{E}_b \left[2r_b \sum_{j=0}^i \mathbb{I}[-q+j \leq b \leq j] \sigma(z+b) \right].$$

Therefore, we can define

$$f(b) := \sum_{i=0}^q \beta_{q-i} \sum_{j=0}^i c_{i,j} \mathbb{I}[-q+j \leq b \leq j],$$

632 which completes the proof. □

633 D.2 Approximating Multivariate Polynomials

634 We adapt the approximation result of this section from [DLS22], modifying the proof to be consistent
 635 with our assumption on the first layer weights.

First, we remark that for any fixed $\mathbf{v} \in \mathbb{S}^{k-1}$ and any degree $0 \leq s \leq p$, we can approximate the function $\mathbf{z} \mapsto \langle \mathbf{v}, \mathbf{z} \rangle^s$ with random biases as established by Lemma 24 for the ReLU activation and Lemma 25 for the polynomial activation. Therefore, our main effort will be spent in approximating a polynomial $h(\mathbf{z})$ using monomials $\langle \mathbf{v}, \mathbf{z} \rangle^s$. Note that we can represent h by

$$h(\mathbf{z}) = \sum_{s=0}^p T^{(s)}[\mathbf{z}^{\otimes s}],$$

where $\mathbf{T}^{(s)}$ is a symmetric tensor of order s , and we use the notation

$$\mathbf{T}^{(s)}[\mathbf{z}^{\otimes s}] = \text{vec}(\mathbf{T}^{(s)})^\top \text{vec}(\mathbf{z}^{\otimes s}) = \sum_{i_1, \dots, i_s=1}^k \mathbf{T}_{i_1, \dots, i_s}^{(s)} z_{i_1} \dots z_{i_s}.$$

636 The approximation result relies on the following fact.

637 **Lemma 26.** *Let $\mathbf{v} \sim \tau_k$. Then, the matrix $\mathbb{E}_{\mathbf{v} \sim \tau_k} [\text{vec}(\mathbf{v}^{\otimes s}) \text{vec}(\mathbf{v}^{\otimes s})^\top]$ is invertible.*

Proof. Let \mathbf{T} be an arbitrary symmetric tensor of order s with $\|\mathbf{T}\|_F = 1$. We need to find a constant $c_{s,k} > 0$ such that

$$\text{vec}(\mathbf{T})^\top \mathbb{E}_{\mathbf{v} \sim \tau_k} [\text{vec}(\mathbf{v}^{\otimes s}) \text{vec}(\mathbf{v}^{\otimes s})] \text{vec}(\mathbf{T}) \geq c_{s,k}.$$

Note that

$$\text{vec}(\mathbf{T})^\top \mathbb{E}_{\mathbf{v} \sim \tau_k} [\text{vec}(\mathbf{v}^{\otimes s}) \text{vec}(\mathbf{v}^{\otimes s})] \text{vec}(\mathbf{T}) = \mathbb{E}_{\mathbf{v} \sim \tau_k} [\mathbf{T}[\mathbf{v}^{\otimes s}]^2] = \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_k)} \left[\frac{\mathbf{T}[\mathbf{w}^{\otimes s}]^2}{\|\mathbf{w}\|^{2s}} \right].$$

Furthermore, [DLS22, Lemma 23] implies that

$$\mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_k)} [\mathbf{T}[\mathbf{w}^{\otimes s}]^2] \geq c'_{s,k},$$

for some constant $c'_{s,k} > 0$. Therefore, for any $r > 0$, we have

$$\mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_k)} [\mathbf{T}[\mathbf{w}^{\otimes s}]^2 \mathbb{I}[\|\mathbf{w}\| > r]] + \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_k)} [\mathbf{T}[\mathbf{w}^{\otimes s}]^2 \mathbb{I}[\|\mathbf{w}\| \leq r]] \geq c'_{s,k}.$$

Note that the first term on the LHS above can become arbitrarily small by choosing r sufficiently large (depending on s and k). Thus for sufficiently large r we have

$$\mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_k)} [\mathbf{T}[\mathbf{w}^{\otimes s}]^2 \mathbb{I}[\|\mathbf{w}\| \leq r]] \geq \frac{c'_{s,k}}{2}.$$

Finally, we have

$$\mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_k)} \left[\frac{\mathbf{T}[\mathbf{w}^{\otimes s}]^2}{\|\mathbf{w}\|^{2s}} \right] \geq \frac{1}{r^{2s}} \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_k)} [\mathbf{T}[\mathbf{w}^{\otimes s}]^2 \mathbb{I}[\|\mathbf{w}\| \leq r]] \geq \frac{c'_{s,k}}{2r^{2s}}.$$

638 Therefore, taking $c_{s,k} = \frac{c'_{s,k}}{2r^{2s}}$ completes the proof. \square

639 The following lemma establishes how we can use monomials of the form $(\mathbf{v}^\top \mathbf{z})^s$ to approximate
640 each term appearing in $h(\mathbf{z})$.

Lemma 27 ([DLS22, Corollary 4, Adapted]). *There exists $f : \mathbb{S}^{k-1} \rightarrow \mathbb{R}$ such that for all $\mathbf{z} \in \mathbb{R}^k$ and non-negative integers $s \geq 0$,*

$$\int_{\mathbb{S}^{k-1}} f(\mathbf{v}) \langle \mathbf{v}, \mathbf{z} \rangle^s d\tau_k(\mathbf{v}) = \mathbf{T}^{(s)}[\mathbf{z}^{\otimes s}].$$

641 *Further, $|f(\mathbf{v})| \leq C_{k,s} \|\mathbf{T}^{(s)}\|_F$ for all $\mathbf{v} \in \mathbb{S}^{k-1}$.*

Proof. Note that by definition, $\langle \mathbf{v}, \mathbf{z} \rangle^s = \text{vec}(\mathbf{v}^{\otimes s})^\top \text{vec}(\mathbf{z}^{\otimes s})$. Therefore,

$$\int f(\mathbf{v}) \langle \mathbf{v}, \mathbf{z} \rangle^s d\tau_k(\mathbf{v}) = \left(\int f(\mathbf{v}) \text{vec}(\mathbf{v}^{\otimes s}) d\tau_k(\mathbf{v}) \right)^\top \text{vec}(\mathbf{z}^{\otimes s}).$$

We need to match the first vector on the RHS above equal to $\text{vec}(\mathbf{T}^{\otimes s})$, thus our choice of f is

$$f(\mathbf{v}) = \text{vec}(\mathbf{v}^{\otimes s})^\top \mathbb{E}_{\mathbf{v} \sim \tau_k} [\text{vec}(\mathbf{v}^{\otimes s}) \text{vec}(\mathbf{v}^{\otimes s})^\top]^{-1} \text{vec}(\mathbf{T}^{(s)}).$$

642 The proof is then completed via the lower bound of Lemma 26 which guarantees the existence of
643 some constant $c_{s,k} > 0$ such that $\lambda_{\min}(\mathbb{E}_{\mathbf{v} \sim \tau_k} [\text{vec}(\mathbf{v}^{\otimes s}) \text{vec}(\mathbf{v}^{\otimes s})^\top]) \geq c_{s,k}$. \square

644 The above result along with the univariate approximations proved earlier immediately yields the
645 following corollary.

Corollary 28. Suppose h is a polynomial of degree p denoted by $h(\mathbf{z}) = \sum_{s=0}^p \mathbf{T}^{(s)}[\mathbf{z}^{\otimes s}]$. Further assume the activation σ is either ReLU or a polynomial of degree $q \geq p$. Then, there exists $\hat{h} : \mathbb{S}^{k-1} \times [-r_b, r_b] \rightarrow \mathbb{R}$ such that for every $\|\mathbf{z}\| \leq \frac{r_b}{2}$, we have

$$\int_{\mathbb{S}^{k-1} \times [-r_b, r_b]} \hat{h}(\mathbf{v}, b) \sigma(\langle \mathbf{v}, \mathbf{z} \rangle + b) d\tau_k(\mathbf{v}) db = h(\mathbf{z}).$$

646 Furthermore, $\left| \hat{h}(\mathbf{v}, b) \right| \leq C_{k,q} \max_{s \leq p} \left\| \mathbf{T}^{(s)} \right\|_{\mathbb{F}}$.

Proof. Let

$$\hat{h}(\mathbf{v}, b) = \sum_{s=0}^k f_{1,s}(\mathbf{v}) f_{2,s}(b),$$

for $(f_{1,s})$ and $(f_{2,s})$ which we now determine. We choose $f_{2,s}$ according to Lemma 24 for the ReLU activation and Lemma 25 for the polynomial activation, then

$$\int_{b=-r_b}^{r_b} f_{2,s}(b) \sigma(\langle \mathbf{v}, \mathbf{z} \rangle + b) db = \langle \mathbf{v}, \mathbf{z} \rangle^s,$$

for all $\|\mathbf{z}\| \leq r_b/2$, and $|f_{2,s}(b)| \leq C_{s,q}$ for all b . Then, we choose $f_{1,s}$ according to Lemma 27, which yields

$$\int_{\mathbb{S}^{k-1} \times [-r_b, r_b]} f_{1,s}(\mathbf{v}) f_{2,s}(b) \sigma(\langle \mathbf{v}, \mathbf{z} \rangle + b) d\tau_k(\mathbf{v}) db = \int f_{1,s}(\mathbf{v}) \langle \mathbf{v}, \mathbf{z} \rangle^s d\tau_k(\mathbf{v}) = \mathbf{T}^{(s)}[\mathbf{z}^{\otimes s}],$$

647 for all $\|\mathbf{z}\| \leq r_b/2$. Additionally $|f_{1,s}(\mathbf{v})| \leq C_{s,k} \left\| \mathbf{T}^{(s)} \right\|_{\mathbb{F}}$, which completes the proof. \square

648 As a last step in this section, we verify that one can indeed control $\max_{s \leq p} \left\| \mathbf{T}^{(s)} \right\|_{\mathbb{F}}$ with an absolute
649 constant when h is the minimizer of the adversarial risk.

Lemma 29. Suppose \mathcal{F} is the class of degree p polynomials on \mathbb{R}^d . Let $\mathcal{H} = \{z \mapsto \mathbb{E}[f(\mathbf{x}) | \mathbf{U}\mathbf{x} = z] : f \in \mathcal{F}\}$, and define

$$h = \arg \min_{h' \in \mathcal{H}} \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} (h'(\mathbf{U}(\mathbf{x} + \delta)) - y)^2 \right].$$

650 Denote the decomposition of h by $h(\mathbf{z}) = \sum_{s=0}^p \mathbf{T}^{(s)}[\mathbf{z}^{\otimes s}]$. Then, $\left\| \mathbf{T}^{(s)} \right\|_{\mathbb{F}} \leq C_{k,y}$, where $C_{k,y}$ is
651 a constant depending only on k and the target second moment $\mathbb{E}[y^2]$ (thus an absolute constant in
652 our setting). As a consequence, we have $|h(\mathbf{z})| \leq L_h (1 + \|\mathbf{z}\|^p)$ for all $\mathbf{z} \in \mathbb{R}^k$, where $L_h > 0$ is an
653 absolute constant.

Proof. By comparing with the zero function, we have

$$\mathbb{E}[(h(\mathbf{U}\mathbf{x}) - y)^2] \leq \mathbb{E} \left[\max_{\|\delta\| \leq \varepsilon} (h(\mathbf{U}(\mathbf{x} + \delta)) - y)^2 \right] \leq \mathbb{E}[y^2]^{1/2}.$$

Furthermore, by the Cauchy-Schwartz inequality,

$$\mathbb{E}[(h(\mathbf{U}\mathbf{x}) - y)^2] \geq \mathbb{E}[h(\mathbf{U}\mathbf{x})^2] + \mathbb{E}[y^2] - 2 \mathbb{E}[h(\mathbf{U}\mathbf{x})^2]^{1/2} \mathbb{E}[y^2]^{1/2}.$$

Combining the two inequalities above, we obtain $\mathbb{E}[h(\mathbf{U}\mathbf{x})^2] \leq 4 \mathbb{E}[y^2]$. Let $\mathbf{z} := \mathbf{U}\mathbf{x}$, and let $\mu_{\mathbf{z}}$ be the marginal distribution of \mathbf{z} . Then

$$\mathbb{E}[h(\mathbf{z})^2] = \int h(\mathbf{z})^2 \frac{d\mu_{\mathbf{z}}}{d\mathcal{N}(0, C_k \mathbf{I}_k)}(\mathbf{z}) d\mathcal{N}(0, C_k \mathbf{I}_k)(\mathbf{z}).$$

Further, by subGaussianity of \mathbf{x} and subsequent subGaussianity of \mathbf{z} , we have $\frac{d\mu_{\mathbf{z}}}{d\mathcal{N}(0, C_k)}(\mathbf{z}) \leq C'_k < \infty$ for all \mathbf{z} , when C_k, C'_k are sufficiently large constants depending only on k . Therefore,

$$\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(0, C_k \mathbf{I}_k)} [h(\mathbf{z})^2] \leq 4C'_k \mathbb{E}[y^2].$$

654 The proof is completed by using the Hermite decomposition of h . \square

655 D.3 Approximating Multivariate Pseudo-Lipschitz Functions

656 We now turn to the more general problem of approximating pseudo-Lipschitz functions. Specifically,
 657 when \mathcal{F} satisfies Assumption 3, functions of the form $h(\mathbf{z}) = \mathbb{E}[f(\mathbf{x}) | \mathbf{U}\mathbf{x} = \mathbf{z}]$ will be L -pseudo-
 658 Lipschitz. The following lemma investigates approximating such functions with infinite-width
 659 two-layer neural networks.

Lemma 30. *Suppose $h : \mathbb{R}^k \rightarrow \mathbb{R}$ is L -Lipschitz on $\|\mathbf{z}\| \leq r_z$ and σ is the ReLU activation. Then, for every $\Delta \geq C_k$, there exists $\hat{h} : \mathbb{S}^{k-1} \times [-r_b, r_b] \rightarrow \mathbb{R}$ such that*

$$\left| h(\mathbf{z}) - \int_{\mathbb{S}^{k-1} \times [-r_b, r_b]} \hat{h}(\mathbf{v}, b) \sigma(\langle \mathbf{v}, \mathbf{z} \rangle + b) d\tau_k(\mathbf{v}) db \right| \leq C_k L r_z \left\{ \left(\frac{\Delta}{L r_z} \right)^{\frac{-2}{k+1}} \log \frac{\Delta}{L r_z} + \left(\frac{\Delta}{r_z} \right)^{\frac{2k}{k+1}} \left(\frac{r_z}{r_b} \right)^k \right\},$$

for all $\|\mathbf{z}\| \leq r_z$. Furthermore, we have $|\hat{h}(\mathbf{v}, b)| \leq C_k L (\Delta / L r_z)^{2k/(k+1)} / r_z$ for all \mathbf{v} and b , and

$$\int_{\mathbb{S}^{k-1} \times [-r_b, r_b]} \hat{h}(\mathbf{v}, b)^2 d\tau_k(\mathbf{v}) db \leq \frac{C_k \Delta^2}{r_z^3}.$$

Proof. Let $\tilde{\mathbf{z}} := (\mathbf{z}^\top, r_z)^\top \in \mathbb{R}^{k+1}$. By [Bac17, Proposition 6], we know that for all $\Delta \geq C_k$, there exists $p : \mathbb{S}^k \rightarrow \mathbb{R}$, such that $\|p\|_{L^2(\tau_{k+1})} \leq \Delta$ and

$$\left| h(\mathbf{z}) - \int_{\mathbb{S}^k} p(\tilde{\mathbf{v}}) \sigma\left(\frac{\langle \tilde{\mathbf{v}}, \tilde{\mathbf{z}} \rangle}{r_z}\right) d\tau_{k+1}(\tilde{\mathbf{v}}) \right| \leq C_k L r_z \left(\frac{\Delta}{L r_z} \right)^{\frac{-2}{k+1}} \log \frac{\Delta}{L r_z},$$

for all $\|\mathbf{z}\| \leq r_z$. Furthermore, the proof of [MHWE24, Proposition 19] demonstrated that

$$|p(\tilde{\mathbf{v}})| \leq C_k L r_z \left(\frac{\Delta}{L r_z} \right)^{\frac{2k}{k+1}}, \quad \forall \tilde{\mathbf{v}} \in \mathbb{S}^k.$$

660 Let $\tilde{\mathbf{v}} = (\tilde{\mathbf{v}}_{1:k}^\top, \tilde{v}_{k+1})^\top$ be the decomposition of $\tilde{\mathbf{v}}$ into its first k and last coordinate. Then, we will use
 661 the fact that for $\tilde{\mathbf{v}} \sim \text{Unif}(\mathbb{S}^k)$ when conditioned on \tilde{v}_{k+1} , by symmetry $\frac{\mathbf{v}_{1:k}}{\|\mathbf{v}_{1:k}\|}$ is uniformly distributed
 662 on \mathbb{S}^{k-1} . In other words, let $\mathbf{v} \sim \text{Unif}(\mathbb{S}^{k-1})$ and $\tilde{b} \sim \rho_{k+1}$ independently, where we choose ρ_{k+1}
 663 such that $\frac{\tilde{b}}{\sqrt{1+\tilde{b}^2}}$ has the same marginal distribution as \tilde{v}_{k+1} . Since the marginal distribution of
 664 \tilde{v}_{k+1} is given by $d\mathbb{P}(\tilde{v}_{k+1}) \propto (1 - \tilde{v}_{k+1}^2)^{(k-2)/2} d\tilde{v}_{k+1}$, we have $\rho_{k+1}(\tilde{b}) = Z_k (1 + \tilde{b}^2)^{-(k+1)/2}$,
 665 where Z_k is the normalizing constant. Then, $\tilde{\mathbf{v}} = \mathbf{T}(\mathbf{v}, \tilde{b})$ is distributed uniformly on \mathbb{S}^k , where
 666 $\mathbf{T} : \mathbb{S}^{k-1} \times \mathbb{R} \rightarrow \mathbb{S}^k$ is given by $\mathbf{T}(\mathbf{v}, \tilde{b}) = \frac{1}{\sqrt{1+\tilde{b}^2}} (\mathbf{v}^\top, \tilde{b})$. As a result,

$$\begin{aligned} \int p(\tilde{\mathbf{v}}) \sigma\left(\frac{\langle \tilde{\mathbf{v}}, \tilde{\mathbf{z}} \rangle}{r_z}\right) d\tau_{k+1}(\tilde{\mathbf{v}}) &= \int p(\mathbf{T}(\mathbf{v}, \tilde{b})) \sigma\left(\frac{\langle \mathbf{v}, \mathbf{z} \rangle + \tilde{b} r_z}{r_z \sqrt{1 + \tilde{b}^2}}\right) d\tau_k(\mathbf{v}) d\rho_{k+1}(\tilde{b}) \\ &= Z_k \int_{\mathbb{S}^{k-1} \times \mathbb{R}} \frac{p(\mathbf{T}(\mathbf{v}, \tilde{b}))}{r_z \sqrt{1 + \tilde{b}^2}} \cdot \frac{1}{(1 + \tilde{b}^2)^{(k+1)/2}} \sigma(\langle \mathbf{v}, \mathbf{z} \rangle + \tilde{b} r_z) d\tau_k(\mathbf{v}) d\tilde{b} \\ &= Z_k \int_{\mathbb{S}^{k-1} \times \mathbb{R}} \frac{r_z^k p(\mathbf{T}(\mathbf{v}, b/r_z))}{(r_z^2 + b^2)^{(k+2)/2}} \sigma(\langle \mathbf{v}, \mathbf{z} \rangle + b) d\tau_k(\mathbf{v}) db. \end{aligned}$$

Therefore, our choice of \hat{h} will be

$$\hat{h}(\mathbf{v}, b) = Z_k \frac{r_z^k p(\mathbf{T}(\mathbf{v}, b/r_z))}{(r_z^2 + b^2)^{(k+2)/2}}.$$

Next, we bound the following error term due to cutoff of bias,

$$\mathcal{E} := \left| \int_{\mathbb{S}^{k-1} \times (\mathbb{R} \setminus [-r_b, r_b])} \frac{r_z^k p(\mathbf{T}(\mathbf{v}, b/r_z))}{(r_z^2 + b^2)^{(k+2)/2}} \sigma(\langle \mathbf{v}, \mathbf{z} \rangle + b) d\tau_k(\mathbf{v}) db \right|.$$

667 We have

$$\begin{aligned}
\mathcal{E} &\lesssim C_k L r_z \left(\frac{\Delta}{L r_z}\right)^{\frac{2k}{k+1}} \int_{|b|>r_b} \frac{r_z^k (r_z + b)}{(r_z^2 + b^2)^{(k+2)/2}} db \\
&\lesssim C_k L r_z \left(\frac{\Delta}{L r_z}\right)^{\frac{2k}{k+1}} \int_{|b|>r_b} \frac{r_z^k}{(r_z^2 + b^2)^{(k+1)/2}} db \\
&\lesssim C_k \Delta^{\frac{2k}{k+1}} \int_{|b|>r_b} \frac{r_z^k}{b^{k+1}} db \\
&\lesssim C_k L r_z \left(\frac{\Delta}{r_z}\right)^{\frac{2k}{k+1}} \left(\frac{r_z}{r_b}\right)^k.
\end{aligned}$$

668 Finally, we prove the guarantees provided for \hat{h} . The uniform bound on $|\hat{h}(\mathbf{v}, b)|$ follows directly by
669 plugging in the uniform bound on p . For the L^2 bound on \hat{h} , we have

$$\begin{aligned}
\int_{\mathbb{S}^{k-1} \times [-r_b, r_b]} \hat{h}(\mathbf{v}, b)^2 d\tau_k(\mathbf{v}) db &\leq \int_{\mathbb{S}^{k-1} \times \mathbb{R}} \hat{h}(\mathbf{v}, b)^2 d\tau_k(\mathbf{v}) db \\
&= \int \frac{Z_k^2 r_z^{2k} p(\mathbb{T}(\mathbf{v}, b/r_z))^2}{(r_z^2 + b^2)^{k+2}} d\tau_k(\mathbf{v}) db \\
&= \int \frac{Z_k^2 p(\mathbb{T}(\mathbf{v}, \tilde{b}))^2}{r_z^3 (1 + \tilde{b}^2)^{k+2}} d\tau_k(\mathbf{v}) d\tilde{b} \\
&= \frac{Z_k}{r_z^3} \int \frac{p(\mathbb{T}(\mathbf{v}, \tilde{b}))^2}{(1 + \tilde{b}^2)^{(k+3)/2}} d\tau_k(\mathbf{v}) d\rho_{k+1}(\tilde{b}) \\
&= \frac{Z_k}{r_z^3} \int (1 - \tilde{v}_{k+1}^2)^{(k+3)/2} p(\tilde{\mathbf{v}})^2 d\tau_{k+1}(\tilde{\mathbf{v}}) \\
&\leq \frac{Z_k \|p\|_{L^2(\tau_{k+1})}^2}{r_z^3} \leq \frac{Z_k \Delta^2}{r_z^3},
\end{aligned}$$

670 completing the proof. □

671 D.4 Discretizing Infinite-Width Approximations

672 In this section, we provide finite-width guarantees corresponding to the infinite-width approximations
673 proved earlier. Define the following integral operator

$$\mathcal{T}\hat{h}(\mathbf{z}) = \int_{\mathbb{S}^{k-1} \times [-r_b, r_b]} \hat{h}(\mathbf{v}, b) \sigma(\langle \mathbf{v}, \mathbf{z} \rangle + b) d\tau_k(\mathbf{v}) db. \quad (\text{D.5})$$

674 The type of discretization error depends on whether we are using the α -DFL or the (α, β) -SFL
675 oracle. We first cover the case of α -DFL oracles.

676 **Proposition 31** (Approximation by Riemann Sum). *Suppose σ satisfies (B.1). Let $(\mathbf{w}_1, \dots, \mathbf{w}_N)$ be
677 the first layer weights obtained from the α -DFL oracle (Definition 2), and define $\mathbf{v}_i = \frac{\mathbf{U}\mathbf{w}_i}{\|\mathbf{U}\mathbf{w}_i\|}$ for*

678 $i \in [N]$. *Suppose $(b_j)_{j \in [N]} \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(-r_b, r_b)$, and let $\|\hat{h}\|_\infty := \sup_{\mathbf{v}, b} |\hat{h}(\mathbf{v}, b)|$. Then, there exists*

679 \mathbf{a}^* *such that $a_j^* = 0$ for $j \notin S$ and $|a_j^*| \leq C_k \|\hat{h}\|_\infty r_b \log(\alpha N / (\zeta \delta)) / (\alpha N)$ for $j \in S$ (where S is
680 given by Definition 21), and*

$$\left| \sum_{j \in S} a_j^* \sigma(\langle \mathbf{v}_j, \mathbf{z} \rangle + b_j) - \mathcal{T}\hat{h}(\mathbf{z}) \right| \leq C_q \|\hat{h}\|_\infty L \sigma r_z^{q-1} r_b \left(r_z \sqrt{\zeta} + \frac{r_b \log(N/\delta)}{\zeta^{(k-1)/2} \alpha N} \right), \quad (\text{D.6})$$

681 for all $\mathbf{z} \in \mathbb{R}^k$ where $\|\mathbf{z}\| \leq r_z$, with probability at least $1 - \delta$ over the randomness of biases.

682 **Proof.** The proof is a multivariate version of the argument given in [OSSW24, Lemma 29]. Let
683 $\{\tilde{\mathbf{v}}_i\}_{i=1}^M$ be the maximal $2\sqrt{2\zeta}$ -packing of \mathbb{S}^{k-1} from Definition 21, which is also a $2\sqrt{2\zeta}$ -covering

684 of \mathbb{S}^{k-1} . Using a lower bound on the surface area of the spherical cap (see e.g. [WMHC24, Lemma
685 F.11]), the maximal packing number is bounded by $M \leq C_k (\frac{1}{\zeta})^{(k-1)/2}$.

For every $i \in [M]$, define

$$S_i := \{j \in [N], \|\mathbf{v}_j - \bar{\mathbf{v}}_i\| \leq \sqrt{2\zeta}\}.$$

686 Note that by definition of packing and Definition 2, each \mathbf{v}_j can only belong to exactly one of S_i
687 when $j \in S$, meaning that (S_i) are disjoint and $\bigcup_{i \in [M]} S_i = S$. In particular, $|S_i|/N \geq \zeta^{(k-1)/2} \alpha$,
688 and $|S_i|/N \leq 1/M \leq c_k \zeta^{(k-1)/2}$.

689 We want each group of biases $(b_j)_{j \in S_i}$ to cover the interval $[-r_b, r_b]$. We divide this interval into $2A$
690 subintervals of the form $[-r_b(1 + \frac{l}{A}), r_b(1 + \frac{l+1}{A})]$ for $0 \leq l \leq 2A - 1$. When $b_j \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(-r_b, r_b)$,
691 by a union bound, the probability that there exists some subinterval and some S_i such that the
692 subinterval contains no element of $\{b_j : j \in S_i\}$ is at most $2AM(1 - \frac{1}{2A})^{|S_i|}$. Thus, taking
693 $A \leq \lfloor \frac{|S_i|}{2 \log(|S_i|/M/\delta)} \rfloor$ for all $i \in [M]$ guarantees that all subintervals have at least one bias from every
694 S_i inside them with probability at least $1 - \delta$.

695 Next, we define $\Pi_1 : \mathbb{S}^{k-1} \rightarrow \mathbb{S}^{k-1}$ as the projection onto the packing, i.e. $\Pi_1(\mathbf{v}) =$
696 $\arg \min_{\{\bar{\mathbf{v}}_i : i \in [M]\}} \|\mathbf{v} - \bar{\mathbf{v}}_i\|$. Further, we define $\Pi_2 : [M] \times [-r_b, r_b] \rightarrow [-r_b, r_b]$ by $\Pi_2(i, b) =$
697 $\arg \min_{\{b_j : j \in S_i\}} |b - b_j|$. Tie braking can be performed by choosing any of the answers. By defini-
698 tion, we then have $\|\mathbf{v} - \Pi_1(\mathbf{v})\| \leq 2\sqrt{2\zeta}$, and $|b - \Pi_2(i, b)| \leq r_b/A$ for all $i \in [M]$.

We are now ready to construct \mathbf{a}^* . Specifically, let

$$a_j^* = \begin{cases} \int \hat{h}(\mathbf{v}, b) \mathbb{I}[i = \Pi_1(\mathbf{v}), b_j = \Pi_2(i, b)] d\tau_k(\mathbf{v}) db & \text{if } j \in S_i \text{ for some } i, \\ 0 & \text{if } j \notin S. \end{cases}$$

Note that by definition,

$$\sum_{i=1}^M \sum_{j \in S_i} \mathbb{I}[i = \Pi_1(\mathbf{v}), b_j = \Pi_2(i, b)] = 1,$$

For conciseness, we define $E(\mathbf{v}, i, j) = \mathbb{I}[i = \Pi_1(\mathbf{v}), b_j = \Pi_2(i, b)]$. When $j \in S_i$, on the event
 $E(\mathbf{v}, i, j)$ we have

$$\|\mathbf{v} - \mathbf{v}_j\| \leq \|\mathbf{v} - \bar{\mathbf{v}}_i\| + \|\bar{\mathbf{v}}_i - \mathbf{v}_j\| \leq 3\sqrt{2\zeta}.$$

Moreover, since $\mathbb{P}_{\mathbf{v} \sim \tau_k} [\|\mathbf{v} - \bar{\mathbf{v}}_i\| \leq 2\sqrt{2\zeta}] \leq C_k \zeta^{(k-1)/2}$, for $j \in S$ we have

$$|a_j^*| \leq \frac{C_k \|\hat{h}\|_\infty \zeta^{(k-1)/2} r_b}{A} \leq \frac{C_k \|\hat{h}\|_\infty r_b \log(N/\delta)}{\alpha N}.$$

699 As a result,

$$\begin{aligned} \left| \sum_{j \in S} a_j^* \sigma(\langle \mathbf{v}_j, \mathbf{z} \rangle + b_j) - \mathcal{T} \hat{h}(\mathbf{z}) \right| &= \left| \sum_{j \in S} a_j^* \sigma(\langle \mathbf{v}_j, \mathbf{z} \rangle + b_j) - \int \hat{h}(\mathbf{v}, b) \sigma(\langle \mathbf{v}, \mathbf{z} \rangle + b) d\tau_k(\mathbf{v}) db \right| \\ &= \left| \sum_{i=1}^M \sum_{j \in S_i} \int \hat{h}(\mathbf{v}, b) E(\mathbf{v}, i, j) (\sigma(\langle \mathbf{v}_j, \mathbf{z} \rangle + b_j) - \sigma(\langle \mathbf{v}, \mathbf{z} \rangle + b)) d\tau_k(\mathbf{v}) db \right| \\ &\lesssim C_q \|\hat{h}\|_\infty L_\sigma r_z^{q-1} r_b (r_z \sqrt{\zeta} + \frac{r_b}{A}), \end{aligned}$$

700 for all $\|\mathbf{z}\| \leq r_z$, where we used the fact that $\sigma(z)$ is $\mathcal{O}(L_\sigma r_z^{q-1})$ Lipschitz when restricted to
701 $|z| \leq r_z$. This concludes the proof. \square

702 Next, we provide a discretization guarantee when using (α, β) -SFL oracles.

Proposition 32. Consider the same setting as Proposition 31, except the first-layer weights
 $(\mathbf{w}_1, \dots, \mathbf{w}_N)$ are obtained from the (α, β) -SFL oracle (Definition 3). Then, there exists \mathbf{a}^* such
that $a_i^* = 0$ for $i \notin S$ and $|a_i^*| \leq \|\hat{h}\|_\infty r_b / (\beta \alpha N)$ for $i \in S$, and

$$\left| \sum_{j \in S} a_j^* \sigma(\langle \mathbf{v}_j, \mathbf{z} \rangle + b_j) - \mathcal{T} \hat{h}(\mathbf{z}) \right| \leq \frac{C_q L_\sigma \|\hat{h}\|_\infty r_b^{q+1}}{\beta} \sqrt{\frac{\log(\alpha N / \delta)}{\alpha N}},$$

703 for all $\mathbf{z} \in \mathbb{R}^k$ with $\|\mathbf{z}\| \leq r_b/2$, with probability at least $1 - \delta$ over the randomness of $(\mathbf{v}_i, b_i)_{i \in [N]}$.
 704 Moreover, suppose $\mathbb{E}_{\mathbf{v}, b \sim \tau_k \otimes \text{Unif}(-r_b, r_b)} [\hat{h}(\mathbf{v}, b)^2] \leq M_2(\hat{h})^2$. Then, assuming $\frac{d\mu}{d\tau_k} \leq \beta'$, we have

$$\|\mathbf{a}^*\|^2 \lesssim \frac{r_b^2 \beta' M_2(\hat{h})^2}{\alpha \beta^2 N}, \quad \text{provided that,} \quad N \gtrsim \frac{\|\hat{h}\|_\infty^4 \log(1/\delta)}{\alpha M_2(\hat{h})^4},$$

705 which also holds with probability at least $1 - \delta$.

706 **Proof.** By definition,

$$\begin{aligned} \mathcal{T}\hat{h}(\mathbf{z}) &= \int_{\mathbb{S}^{k-1} \times [-r_b, r_b]} \hat{h}(\mathbf{v}, b) \sigma(\langle \mathbf{v}, \mathbf{z} \rangle + b) d\tau_k(\mathbf{v}) db \\ &= \int_{\mathbb{S}^{k-1} \times [-r_b, r_b]} \hat{h}(\mathbf{v}, b) \frac{d\tau_k}{d\mu}(\mathbf{v}) \sigma(\langle \mathbf{v}, \mathbf{z} \rangle + b) d\mu(\mathbf{v}) db \\ &= \mathbb{E}_{\mathbf{v}, b \sim \mu \otimes \text{Unif}(-r_b, r_b)} \left[2r_b \hat{h}(\mathbf{v}, b) \frac{d\tau_k}{d\mu}(\mathbf{v}) \sigma(\langle \mathbf{v}, \mathbf{z} \rangle + b) \right]. \end{aligned}$$

Consider $(\mathbf{v}_i, b_i)_{i \in S} \stackrel{\text{i.i.d.}}{\sim} \mu \otimes \text{Unif}(-r_b, r_b)$ from Definition 3. Let

$$a_i^* = \begin{cases} \frac{2r_b \hat{h}(\mathbf{v}_i, b_i)}{|S|} \frac{d\tau_k}{d\mu}(\mathbf{v}_i) & \text{if } i \in S, \\ 0 & \text{if } i \notin S. \end{cases}$$

Consequently

$$|a_i^*| \leq \frac{2r_b \|\hat{h}\|_\infty}{\beta |S|},$$

for all $i \in S$. Given \mathbf{z} , define the random variable

$$\hat{\mathcal{T}}\hat{h}(\mathbf{z}) = \sum_{i \in S} a_i^* \sigma(\langle \mathbf{v}_i, \mathbf{z} \rangle + b_i).$$

707 Our next step is to bound the difference between $\hat{\mathcal{T}}\hat{h}(\mathbf{z})$ and $\mathcal{T}\hat{h}(\mathbf{z})$ uniformly over all $\|\mathbf{z}\| \leq r_z$.

Let $(\hat{\mathbf{z}}_j)_{j=1}^M$ be a Δ -covering of $\{\mathbf{z} : \|\mathbf{z}\| \leq r_z\}$, therefore $M \leq (3r_z/\Delta)^k$. Note that for any fixed \mathbf{z} with $\|\mathbf{z}\| \leq r_z$, we have $|\hat{\mathcal{T}}\hat{h}(\mathbf{z})| \lesssim \|\hat{h}\|_\infty L_\sigma r_b r_z^{q-1} (r_z + r_b) / \beta$. Thus, by Hoeffding's lemma,

$$|\hat{\mathcal{T}}\hat{h}(\mathbf{z}) - \mathcal{T}\hat{h}(\mathbf{z})| \lesssim \frac{\|\hat{h}\|_\infty L_\sigma r_b r_z^{q-1} (r_z + r_b)}{\beta} \sqrt{\frac{\log(1/\delta)}{|S|}},$$

with probability at least $1 - \delta$ for a fixed \mathbf{z} . By a union bound,

$$\max_{j \in [M]} |\hat{\mathcal{T}}\hat{h}(\hat{\mathbf{z}}_j) - \mathcal{T}\hat{h}(\hat{\mathbf{z}}_j)| \lesssim \frac{\|\hat{h}\|_\infty L_\sigma r_b r_z^{q-1} (r_z + r_b)}{\beta} \sqrt{\frac{\log(M/\delta)}{|S|}},$$

708 with probability at least $1 - \delta$. For any \mathbf{z} with $\|\mathbf{z}\| \leq r_z$, let $\hat{\mathbf{z}}$ denote the projection of \mathbf{z} onto the
 709 covering $(\hat{\mathbf{z}}_j)_{j=1}^M$. Then,

$$\begin{aligned} \sup_{\|\mathbf{z}\| \leq r_b/2} |\hat{\mathcal{T}}\hat{h}(\mathbf{z}) - \mathcal{T}\hat{h}(\mathbf{z})| &\leq \max_{j \in [M]} |\hat{\mathcal{T}}\hat{h}(\hat{\mathbf{z}}_j) - \mathcal{T}\hat{h}(\hat{\mathbf{z}}_j)| + |\mathcal{T}\hat{h}(\hat{\mathbf{z}}) - \mathcal{T}\hat{h}(\mathbf{z})| + |\hat{\mathcal{T}}\hat{h}(\hat{\mathbf{z}}) - \mathcal{T}\hat{h}(\mathbf{z})| \\ &\lesssim \frac{\|\hat{h}\|_\infty L_\sigma r_b r_z^{q-1} (r_z + r_b)}{\beta} \sqrt{\frac{\log(M/\delta)}{|S|}} + \frac{\|\hat{h}\|_\infty L_\sigma r_b r_z^{q-1} \Delta}{\beta} \\ &\lesssim \frac{\|\hat{h}\|_\infty L_\sigma r_b r_z^{q-1} (r_z + r_b)}{\beta} \sqrt{\frac{\log(r_b/(\Delta\delta))}{|S|}} + \frac{\|\hat{h}\|_\infty L_\sigma r_b r_z^{q-1} \Delta}{\beta}. \end{aligned}$$

Choosing $\Delta = r_b/\sqrt{|S|}$ implies

$$\sup_{\|\mathbf{z}\| \leq r_b/2} |\hat{\mathcal{T}}\hat{h}(\mathbf{z}) - \mathcal{T}\hat{h}(\mathbf{z})| \lesssim \frac{C_k \|\hat{h}\|_\infty L_\sigma r_b r_z^{q-1} (r_z + r_b)}{\beta} \sqrt{\frac{\log(|S|/\delta)}{|S|}}$$

710 with probability at least $1 - \delta$ over the randomness of $(\mathbf{v}_i, b_i)_{i \in [N]}$.

The last step is to bound $\|\mathbf{a}^*\|^2$. Note that,

$$\|\mathbf{a}^*\|^2 \leq \frac{4r_b^2}{\beta^2|S|} \sum_{i \in S} \frac{\hat{h}(\mathbf{v}_i, b_i)^2}{|S|}.$$

Further, by the Hoeffding inequality,

$$\sum_{i \in S} \frac{\hat{h}(\mathbf{v}_i, b_i)^2}{|S|} - \mathbb{E}_{\mathbf{v}, b \sim \mu \otimes \text{Unif}(-r_b, r_b)} [\hat{h}(\mathbf{v}, b)^2] \lesssim \|\hat{h}\|_\infty^2 \sqrt{\frac{\log(1/\delta)}{|S|}}$$

711 with probability at least $1 - \delta$. Further,

$$\begin{aligned} \mathbb{E}_{\mathbf{v}, b \sim \mu \otimes \text{Unif}(-r_b, r_b)} [\hat{h}(\mathbf{v}, b)^2] &= \mathbb{E}_{\mathbf{v}, b \sim \tau_k \otimes \text{Unif}(-r_b, r_b)} \left[\hat{h}(\mathbf{v}, b)^2 \frac{d\mu}{d\tau_k}(\mathbf{v}) \right] \\ &\leq \beta' M_2(\hat{h})^2. \end{aligned}$$

712 Thus, when $|S| \geq \frac{\|\hat{h}\|_\infty^4 \log(1/\delta)}{M_2(\hat{h})^4}$, we have $\|\mathbf{a}^*\|^2 \lesssim r_b^2 \beta' M_2(\hat{h})^2 / (\beta^2 |S|)$ with probability at least
713 $1 - \delta$, which completes the proof. \square

714 D.5 Combining All Steps

715 We can finally bound our original objective of this section, i.e. $\text{AR}(\mathbf{a}^*, \mathbf{W}, \mathbf{b}) - \text{AR}^*$. Let us begin
716 with the case where \mathcal{F} is the class of polynomials of degree p .

717 **Proposition 33.** *Suppose \mathcal{F} and σ satisfy Assumption 4 and $(b_i)_{i \in [N]} \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(-r_b, r_b)$. Recall
718 that $\varepsilon_1 := 1 \vee \varepsilon$, and $\tilde{\varepsilon} := \varepsilon \wedge \frac{\varepsilon^2}{\text{AR}^*}$ for any $\varepsilon \in (0, 1)$. Using the simplification $k, q, p, L_\sigma \lesssim 1$ and
719 recalling $\varepsilon_1 := 1 \vee \varepsilon$, there exists a choice of $r_b = \tilde{\Theta}(\varepsilon_1)$ such that:*

720 • If $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_N)^\top$ are given by the α -DFL oracle, there exists \mathbf{a}^* such that $|a_i^*| \leq$
721 $\tilde{\mathcal{O}}(\varepsilon_1 / (\alpha N))$ for all $i \in [N]$, and $\text{AR}(\mathbf{a}^*, \mathbf{W}, \mathbf{b}) - \text{AR}^* \leq \varepsilon$ as soon as

$$\zeta \leq \tilde{\mathcal{O}}\left(\frac{\tilde{\varepsilon}}{\varepsilon_1^{2(q+1)}}\right) \quad \text{and} \quad N \geq \tilde{\Omega}\left(\frac{\varepsilon_1^{q+1}}{\alpha \zeta^{(k-1)/2} \sqrt{\tilde{\varepsilon}}}\right).$$

722 • If $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_N)^\top$ are given by the (α, β) -SFL oracle, there exists \mathbf{a}^* such that $|a_i^*| \leq$
723 $\tilde{\mathcal{O}}(\varepsilon_1 / (\beta \alpha N))$ for all $i \in [N]$, and $\text{AR}(\mathbf{a}^*, \mathbf{W}, \mathbf{b}) - \text{AR}^* \leq \varepsilon$ as soon as

$$\zeta \leq \tilde{\mathcal{O}}\left(\frac{\beta^2 \tilde{\varepsilon}}{\varepsilon_1^{2(q+1)}}\right) \quad \text{and} \quad N \geq \tilde{\Omega}\left(\frac{\varepsilon_1^{2(q+1)}}{\alpha \beta^2 \tilde{\varepsilon}}\right).$$

724 Both cases above hold with probability at least $1 - n^{-c}$ for some absolute constant $c > 0$ over the
725 choice of random biases $(b_i)_{i \in [N]}$ (and random weights (\mathbf{w}_i) in the case of SFL).

Proof. Recall from (D.1) that

$$\text{AR}(\mathbf{a}^*, \mathbf{W}, \mathbf{b}) - \text{AR}^* \lesssim \mathcal{E}_1 + \mathcal{E}_2 + \sqrt{\mathcal{E}_3(\mathcal{E}_1 + \mathcal{E}_2)}.$$

By definition, $\mathcal{E}_3 \lesssim \text{AR}^*$. By Lemma 22, we have

$$\mathcal{E}_2 \lesssim L_\sigma^2 \tilde{r}_a^2 (1 + r_b^{2(q-1)} + \varepsilon^{2(q-1)}) (1 + \varepsilon^2) \zeta.$$

Further, thanks to Lemma 29, by choosing $r_z = r_b/2$, we have $|h(\mathbf{z})| \lesssim 1 + \|\mathbf{z}\|^p$. Therefore, by
Lemma 23,

$$\mathcal{E}_1 \lesssim \varepsilon_{\text{approx}}^2 + (L_\sigma^2 \tilde{r}_a^2 (1 + \varepsilon^{2q} + r_b^{2q}) + 1 + \varepsilon^{2p}) e^{-\Omega(r_b^2)}.$$

Let us now consider the case of α -DFL. By Proposition 31, we know there exists \mathbf{a}^* with $|a_i^*| \leq \tilde{\mathcal{O}}(r_b/(\alpha N))$ (we used the fact that $\max_{s \leq p} \|\mathbf{T}^{(s)}\|_{\mathbb{F}} \lesssim 1$ from Lemma 29) such that

$$\epsilon_{\text{approx}} \leq \tilde{\mathcal{O}}\left(r_b^{q+1}\left(\sqrt{\zeta} + \frac{r_b}{\zeta^{(k-1)/2}\alpha N}\right)\right),$$

provided that $r_b \gtrsim \varepsilon_1$ where we recall $\varepsilon_1 = 1 \vee \varepsilon$, and the above statements with probability at least $1 - \delta$ for any polynomially decaying δ , e.g. $\delta = n^{-c}$ for some absolute constant $c > 0$. Therefore, we have $\tilde{r}_a \leq \tilde{\mathcal{O}}(r_b)$. Further, it suffices to choose r_b large enough such that $r_b \gtrsim \varepsilon_1 \vee \sqrt{\log(NL_\sigma^2 \tilde{r}_a^2 r_b^{2q} + \varepsilon_1^{2p})} = \tilde{\Theta}(\varepsilon_1)$ to have

$$\mathcal{E}_1 \leq \tilde{\mathcal{O}}\left(r_b^{2(q+1)}\left(\zeta + \frac{1}{\zeta^{k-1}\alpha^2 N^2}\right)\right).$$

726 Plugging in the values of \tilde{r}_a and r_b , we obtain,

$$\mathcal{E}_2 \leq \tilde{\mathcal{O}}(\varepsilon_1^{2(q+1)}\zeta), \quad \text{and} \quad \mathcal{E}_1 \leq \tilde{\mathcal{O}}\left(\varepsilon_1^{2(q+1)}\zeta + \frac{\varepsilon_1^{2(q+1)}}{\zeta^{k-1}\alpha^2 N^2}\right).$$

727 Hence, choosing

$$\zeta \leq \tilde{\mathcal{O}}\left(\frac{\tilde{\varepsilon}}{\varepsilon_1^{2(q+1)}}\right), \quad \text{and}, \quad N \geq \tilde{\Omega}\left(\frac{\varepsilon_1^{q+1}}{\alpha\zeta^{(k-1)/2}\sqrt{\tilde{\varepsilon}}}\right)$$

728 which concludes the proof of the α -DFL case.

In the case of (α, β) -SFL, we instead invoke Proposition 32, thus obtain $|a_i^*| \lesssim r_b/(\beta\alpha N)$, and

$$\epsilon_{\text{approx}} \leq \tilde{\mathcal{O}}\left(\frac{L_\sigma r_b^{q+1}}{\beta\sqrt{\alpha N}}\right),$$

729 which holds with probability at least $1 - \delta$ for any polynomially decaying δ such as $\delta = n^{-c}$ for
730 some absolute constant $c > 0$. Consequently, with the same choice of $r_b = \tilde{\Theta}(\varepsilon_1)$ as before, we have

$$\mathcal{E}_2 \leq \tilde{\mathcal{O}}\left(\frac{\varepsilon_1^{2(q+1)}\zeta}{\beta^2}\right) \quad \text{and} \quad \mathcal{E}_1 \leq \tilde{\mathcal{O}}\left(\frac{\varepsilon_1^{2(q+1)}}{\beta^2\alpha N}\right),$$

731 which completes the proof. \square

732 We can also combine approximation bounds for the more general class of pseudo-Lipschitz \mathcal{F} .

733 **Proposition 34.** *Suppose \mathcal{F} and σ satisfy Assumption 3 and $(b_i)_{i \in [N]} \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(-r_b, r_b)$. Recall
734 that $\varepsilon_1 := 1 \vee \varepsilon$, and $\tilde{\varepsilon} := \varepsilon \wedge \frac{\varepsilon^2}{\text{AR}^*}$ for any $\varepsilon \in (0, 1)$. Using the simplification $k, p, L \lesssim 1$, there
735 exists a choice of $r_b = \tilde{\Theta}(\varepsilon_1(\varepsilon_1/\sqrt{\tilde{\varepsilon}})^{1+1/k})$ such that:*

736 • If $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_N)^\top$ is given by the α -DFL oracle, there exists \mathbf{a}^* such that $|a_i^*| \leq$
737 $\tilde{\mathcal{O}}((\varepsilon_1/\sqrt{\tilde{\varepsilon}})^{k+1+1/k}/(\alpha N))$ for all $i \in [N]$, and $\text{AR}(\mathbf{a}^*, \mathbf{W}, \mathbf{b}) - \text{AR}^* \leq \varepsilon$ as soon as

$$\zeta \leq \tilde{\mathcal{O}}\left(\left(\frac{\tilde{\varepsilon}}{\varepsilon_1^2}\right)^{k+2+1/k}\right), \quad \text{and} \quad N \geq \tilde{\Omega}\left(\frac{1}{\zeta^{(k-1)/2}\alpha}\left(\frac{\varepsilon_1}{\sqrt{\tilde{\varepsilon}}}\right)^{k+3+2/k}\right).$$

738 • If $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_N)^\top$ is given by the (α, β) -SFL oracle, there exists \mathbf{a}^* such that
739 $|a_i^*| \leq \tilde{\mathcal{O}}((\varepsilon_1/\sqrt{\tilde{\varepsilon}})^{k+1+1/k}/(\alpha\beta N))$, and $\text{AR}(\mathbf{a}^*, \mathbf{W}, \mathbf{b}) - \text{AR}^* \leq \varepsilon$ as soon as

$$\zeta \leq \tilde{\mathcal{O}}\left(\beta^2\left(\frac{\tilde{\varepsilon}}{\varepsilon_1^2}\right)^{k+2+1/k}\right), \quad \text{and} \quad N \geq \tilde{\Omega}\left(\frac{1}{\alpha\beta^2}\left(\frac{\varepsilon_1}{\tilde{\varepsilon}}\right)^{k+3+2/k}\right).$$

740 Both cases above hold with probability at least $1 - n^{-c}$ for some absolute constant $c > 0$ over the
741 choice of random biases $(b_i)_{i \in [N]}$ (and random weights (\mathbf{w}_i) in the case of SFL).

Proof. Our starting point is once again the decomposition

$$\text{AR}(\mathbf{a}^*, \mathbf{W}, \mathbf{b}) - \text{AR}^* \leq \mathcal{E}_1 + \mathcal{E}_2 + \sqrt{\mathcal{E}_3(\mathcal{E}_1 + \mathcal{E}_2)}.$$

742 Given Assumption 3, it is straightforward to verify that $|h(\mathbf{z}_1) - h(\mathbf{z}_2)| \lesssim (\varepsilon_1^{1-p} \|\mathbf{z}_1\|^{p-1} +$
 743 $\varepsilon_1^{1-p} \|\mathbf{z}_2\|^{p-1} + 1) \|\mathbf{z}_1 - \mathbf{z}_2\|$ for $\mathbf{z}_1, \mathbf{z}_2 \in \mathbb{R}^k$. As a consequence, we have $|h(\mathbf{z})| \lesssim 1 + \|\mathbf{z}\|^p$
 744 for all $\mathbf{z} \in \mathbb{R}^k$. Therefore, by Lemma 23 with a choice of $r_z = \tilde{\Theta}(\varepsilon_1)$, we have $\mathcal{E}_1 \lesssim \epsilon_{\text{approx}}^2$. In the
 745 rest of the proof we will fix $r_z = \tilde{\Theta}(\varepsilon_1)$.

746 We begin with considering the case of α -DFL. Unlike the proof of Proposition 33 where $\mathcal{T}\hat{h} = h$, in
 747 this case we have an additional error due to $\mathcal{T}\hat{h}$ only approximating h . From Lemma 30, we have

$$\|\hat{h}\|_\infty \leq \tilde{\mathcal{O}}\left(\frac{1}{\varepsilon_1} \left(\frac{\Delta}{\varepsilon_1}\right)^{2k/(k+1)}\right).$$

748 Thus,

$$\begin{aligned} \epsilon_{\text{approx}} &\leq \sup_{\|\mathbf{z}\| \leq r_z} \left| \sum_{j \in S} a_j^* \sigma(\langle \mathbf{v}_j, \mathbf{z} \rangle + b_j) - \mathcal{T}\hat{h}(\mathbf{z}) \right| + \left| \mathcal{T}\hat{h}(\mathbf{z}) - h(\mathbf{z}) \right| \\ &\leq \tilde{\mathcal{O}}\left(\frac{r_b}{\varepsilon_1} \left(\frac{\Delta}{\varepsilon_1}\right)^{\frac{2k}{k+1}} (\sqrt{\zeta} + \frac{r_b}{\zeta^{(k-1)/2} \alpha N})\right) + \tilde{\mathcal{O}}\left(\varepsilon_1 \left(\frac{\Delta}{\varepsilon_1}\right)^{-\frac{2}{k+1}} + \varepsilon_1 \left(\frac{\Delta}{\varepsilon_1}\right)^{\frac{2k}{k+1}} \left(\frac{\varepsilon_1}{r_b}\right)^k\right), \end{aligned}$$

where we bounded the first term via Proposition 31 with $q = 1$, and the second term via Lemma 30. Additionally, we have

$$|a_j^*| \leq \frac{r_b}{\varepsilon_1 \alpha N} \left(\frac{\Delta}{\varepsilon_1}\right)^{\frac{2k}{k+1}},$$

749 for all $j \in [N]$. To obtain $\text{AR}(\mathbf{a}^*, \mathbf{W}, \mathbf{b}) \leq \text{AR}^* + \epsilon$, we must choose $\Delta = \tilde{\Theta}(\varepsilon_1 (\varepsilon_1 / \sqrt{\tilde{\epsilon}})^{(k+1)/2})$.
 750 Next, we choose $r_b = \tilde{\Theta}(\varepsilon_1 (\varepsilon_1 / \sqrt{\tilde{\epsilon}})^{(k+1)/k})$. This combination ensures $|\mathcal{T}\hat{h}(\mathbf{z}) - h(\mathbf{z})| \lesssim \sqrt{\tilde{\epsilon}}$. To
 751 make sure $|\hat{\mathcal{T}}\hat{h}(\mathbf{z}) - \mathcal{T}\hat{h}(\mathbf{z})| \lesssim \sqrt{\tilde{\epsilon}}$, we should let

$$\zeta \leq \tilde{\mathcal{O}}\left(\tilde{\epsilon} \left(\frac{\tilde{\epsilon}}{\varepsilon_1^2}\right)^{k+1+1/k}\right), \quad \text{and} \quad N = \tilde{\Theta}\left(\frac{1}{\zeta^{(k-1)/2} \alpha} \left(\frac{\varepsilon_1}{\sqrt{\tilde{\epsilon}}}\right)^{k+3+2/k}\right).$$

The above guarantee that $\epsilon_{\text{approx}} \lesssim \sqrt{\tilde{\epsilon}}$ and consequently $\mathcal{E}_1 + \sqrt{\mathcal{E}_3 \mathcal{E}_1} \lesssim \epsilon$. Note that the above choices imply $|a_j^*| \leq \tilde{r}_a / |S|$ for all $i \in S$ with $\tilde{r}_a = \tilde{\mathcal{O}}((\varepsilon_1 / \sqrt{\tilde{\epsilon}})^{k+1+1/k})$. From Lemma 22 with $q = 1$, we have $\mathcal{E}_2 \lesssim \tilde{r}_a^2 \varepsilon_1^2 \zeta$. Therefore, if we let

$$\zeta = \tilde{\Theta}\left(\left(\frac{\tilde{\epsilon}}{\varepsilon_1^2}\right)^{k+2+1/k}\right),$$

752 we have $\mathcal{E}_2 \lesssim \tilde{\epsilon}$ and consequently $\mathcal{E}_2 + \sqrt{\mathcal{E}_3 \mathcal{E}_1} \lesssim \epsilon$. This concludes the proof of the α -DFL case.

753 Next, we consider the case of (α, β) -SFL. Note that the error $|\mathcal{T}\hat{h}(\mathbf{z}) - h(\mathbf{z})|$ remains unchanged.

754 However, this time we invoke Proposition 32 for controlling $|\hat{\mathcal{T}}\hat{h}(\mathbf{z}) - \mathcal{T}\hat{h}(\mathbf{z})|$. Therefore,

$$\begin{aligned} \epsilon_{\text{approx}} &\leq \sup_{\|\mathbf{z}\| \leq r_z} \left| \sum_{j \in S} a_j^* \sigma(\langle \mathbf{v}_j, \mathbf{z} \rangle + b_j) - \mathcal{T}\hat{h}(\mathbf{z}) \right| + \left| \mathcal{T}\hat{h}(\mathbf{z}) - h(\mathbf{z}) \right| \\ &\leq \tilde{\mathcal{O}}\left(\frac{r_b^2}{\beta \varepsilon_1} \left(\frac{\Delta}{\varepsilon_1}\right)^{\frac{2k}{k+1}} \sqrt{\frac{1}{\alpha N}}\right) + \tilde{\mathcal{O}}\left(\varepsilon_1 \left(\frac{\Delta}{\varepsilon_1}\right)^{-\frac{2}{k+1}} + \varepsilon_1 \left(\frac{\Delta}{\varepsilon_1}\right)^{\frac{2k}{k+1}} \left(\frac{\varepsilon_1}{r_b}\right)^k\right). \end{aligned}$$

755 Since the second term is unchanged, we have the same choices of $\Delta = \tilde{\Theta}(\varepsilon_1 (\varepsilon_1 / \sqrt{\tilde{\epsilon}})^{(k+1)/2})$ and
 756 $r_b = \tilde{\Theta}(\varepsilon_1 (\varepsilon_1 / \sqrt{\tilde{\epsilon}})^{1+1/k})$ as in the α -DFL case. However, for the finite-width discretization, we
 757 should choose

$$N = \tilde{\Theta}\left(\frac{1}{\alpha \beta^2} \left(\frac{\varepsilon_1^2}{\tilde{\epsilon}}\right)^{k+3+2/k}\right). \quad (\text{D.7})$$

Moreover, Proposition 32 implies $|a_j^*| \leq \tilde{r}_a/|S|$ with $\tilde{r}_a = \tilde{O}((\varepsilon_1/\sqrt{\tilde{\varepsilon}})^{k+1+1/k}/\beta)$. As a result, to get $\mathcal{E}_2 \lesssim \tilde{r}_a^2 \varepsilon_1^2 \zeta \leq \tilde{\varepsilon}$ from Lemma 22 with $q = 1$, we let

$$\zeta = \tilde{\Theta}\left(\beta^2 \left(\frac{\tilde{\varepsilon}}{\varepsilon_1^2}\right)^{k/2+2+1/(2k)}\right),$$

758 completing the proof.

759

□

760 E Numerical Experiments

761 We perform the following small-scale experiment to support intuitions from our theory. We consider
 762 a single-index setting, where the teacher non-linearity is given by either ReLU, tanh, or $\text{He2}(z) =$
 763 $(z^2 - 1)/\sqrt{2}$ which is the normalized second Hermite polynomial. The student network has $N = 100$
 764 neurons, and the input is sampled from $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$ with $d = 100$. We implement adversarial
 765 training in the following manner. At each iteration, we sample a new batch of i.i.d. training examples.
 766 We estimate the adversarial perturbations on this batch by performing 5 steps of signed projected
 767 gradient ascent, with a stepsize of 0.1. We then perform a gradient descent step on the perturbed
 768 batch. To estimate the robust test risk, we fix a test set of 10,000 i.i.d. samples, and use 20 iterations
 769 to estimate the adversarial perturbation. Because of the online nature of the algorithm, the total
 770 number of samples used is the batch size times the number of iterations taken.

771 The first row of Figure 1 compares the performance of three different approaches. Full AD training
 772 refers to adversarially training all layers from random initialization, where first layer weights are
 773 initialized uniformly on the sphere \mathbb{S}^{d-1} , second layer weights are initialized i.i.d. from $\mathcal{N}(0, 1/N^2)$,
 774 and biases are initialized i.i.d. from $\mathcal{N}(0, 1)$. In the two other approaches, we initialize all first layer
 775 weights to the target direction \mathbf{u} . In one approach we fix this direction and do not train it, while in the
 776 other approach we allow the training of first layer weights from this initialization. As can be seen
 777 from Figure 1, there is a considerable improvement in initializing from \mathbf{u} , which is consistent with
 778 our theory that this direction provides a Bayes optimal projection for robust learning.

779 In the typical setting where we do not have knowledge of \mathbf{u} , we consider the following alternative. We
 780 first perform standard training on the network, i.e. assume $\varepsilon = 0$ (denoted in Figure 1 by SD training).
 781 We can then either fix the first layer weights to these directions, or further train them adversarially
 782 from this initialization. Note that for a fair comparison with the full AD method, we provide the
 783 same random bias and second layer weight initializations across all methods at the beginning of
 784 the adversarial training stage. Even though this approach is not perfect at estimating the unknown
 785 direction, it still provides a considerable benefit over adversarial training of all layer from random
 786 initialization, as demonstrated in the second row of Figure 1.