

# 000 001 002 003 004 005 LEARNING HIERARCHICAL POLYNOMIALS OF MULTI- 006 PLE NONLINEAR FEATURES 007 008 009

010 **Anonymous authors**  
 011 Paper under double-blind review  
 012  
 013  
 014  
 015  
 016  
 017  
 018  
 019

## ABSTRACT

020 In deep learning theory, a critical question is to understand how neural networks  
 021 learn hierarchical features. In this work, we study the learning of hierarchical  
 022 polynomials of *multiple nonlinear features* using three-layer neural networks. We  
 023 examine a broad class of functions of the form  $f^* = g^* \circ p$ , where  $p : \mathbb{R}^d \rightarrow$   
 024  $\mathbb{R}^r$  represents multiple quadratic features with  $r \ll d$  and  $g^* : \mathbb{R}^r \rightarrow \mathbb{R}$  is a  
 025 polynomial of degree  $p$ . This can be viewed as a nonlinear generalization of the  
 026 multi-index model (Damian et al., 2022), and also an expansion upon previous  
 027 work that focused only on a single nonlinear feature, i.e.  $r = 1$  (Nichani et al.,  
 028 2023; Wang et al., 2023).  
 029

030 Our primary contribution shows that a three-layer neural network trained via lay-  
 031 erwise gradient descent suffices for

- 032 • complete recovery of the space spanned by the nonlinear features
- 033 • efficient learning of the target function  $f^* = g^* \circ p$  or transfer learning of
- 034  $f = g \circ p$  with a different link function

035 within  $\tilde{\mathcal{O}}(d^4)$  samples and polynomial time. For such hierarchical targets, our  
 036 result substantially improves the sample complexity  $\Theta(d^{2p})$  of the kernel methods,  
 037 demonstrating the power of efficient feature learning. It is important to highlight  
 038 that our results leverage novel techniques and thus manage to go beyond all prior  
 039 settings such as single-index and multi-index models as well as models depending  
 040 just on one nonlinear feature, contributing to a more comprehensive understanding  
 041 of feature learning in deep learning.

## 042 1 INTRODUCTION

043 Deep neural networks have achieved remarkable empirical success across numerous domains of  
 044 artificial intelligence (Krizhevsky et al., 2012; He et al., 2016). This success can be largely attributed  
 045 to their ability to extract latent features from real-world data and decompose complex targets into  
 046 hierarchical representations, which improves test accuracy (He et al., 2016) and allows efficient  
 047 transfer learning (Devlin, 2018). These feature learning capabilities are widely regarded as a core  
 048 strength of neural networks over non-adaptive approaches such as kernel methods (Wei et al., 2020;  
 049 Bai and Lee, 2020).

050 Despite these empirical achievements, the feature learning capabilities of neural networks are less  
 051 well understood from a theoretical point of view. Previous work on feature learning has shown that  
 052 two-layer neural networks can learn *multiple* linear features of the input (Damian et al., 2022), that  
 053 is, multi-index models. However, the two-layer architecture inherently limits the network’s ability  
 054 to represent and learn nonlinear features (Daniely, 2017). Given that many real-world scenarios  
 055 involve diverse and nonlinear features, recent studies have shifted focus to investigating the learning  
 056 of *nonlinear* features using deeper neural networks. Safran and Lee (2022); Ren et al. (2023);  
 057 Nichani et al. (2023); Wang et al. (2023) have demonstrated that three-layer networks, when trained  
 058 via gradient descent, can efficiently learn *hierarchical targets* of the form  $h = g \circ p$ , where  $p$   
 059 represents certain types of features such as the norm  $|\mathbf{x}|$  or a quadratic form  $\mathbf{x}^\top \mathbf{A} \mathbf{x}$ . However,  
 060 these studies are limited to relatively simple hierarchical functions and mainly focus on targets of  
 061 a single feature. It remains unclear whether neural networks can efficiently learn a wider range of  
 062 hierarchical functions, particularly those that depend on multiple nonlinear features. This leads us  
 063 to the following central question:

054     *Can neural networks adaptively identify **multiple nonlinear features** from the hierarchical targets*  
 055     *by gradient descent, thereby allowing an efficient learning for such targets?*

### 056     1.1 MAIN CONTRIBUTIONS

058     In this paper, we provide strong theoretical evidence that three-layer neural networks have the ability  
 059     to learn multiple hidden nonlinear features. Specifically, we study the problem of learning any  
 060     hierarchical polynomial with multiple quadratic features using a three-layer network trained via  
 061     layer-wise gradient descent. Our main contributions are summarized as follows:

- 063     • **A Novel Analytic Framework for Multi-Nonlinear Feature Learning.** We demonstrate  
 064       that when the target function belongs to a broad class of the form  $f^* = g^* \circ p$ , where  
 065        $p : \mathbb{R}^d \rightarrow \mathbb{R}^r$  represents  $r$  quadratic (nonlinear) features and  $g^*$  is a link function, the  
 066       first step of gradient descent efficiently learns and recovers the space spanned by these  
 067       nonlinear features  $p$  within only  $\tilde{\mathcal{O}}(d^4)$  samples. **We remark that our proof techniques**  
 068       **are also applicable to general nonlinear features.** The core technical novelty is that we  
 069       develop a novel and general universality argument (Lemma 1) that bridges multi nonlinear  
 070       feature models to multi-index models, which allows for an accurate reconstruction of the  
 071       features through a simple linear transformation on the learned representations with small  
 072       approximation error (Proposition 1)
- 073     • **Improved Sample Complexity and Efficient Transfer Learning.** Leveraging the learned  
 074       features in the first GD step, we prove that when the link function  $g^*$  is a polynomial of  
 075       degree  $p$ , the gradient descent on the outer layer can achieves a vanishing generalization  
 076       error with a small outer width and at most  $\mathcal{O}(r^{\mathcal{O}(p)})$  additional training samples, removing  
 077       the dependence on  $d$  (Theorem 1). This significantly improves upon the sample complexity  
 078       of kernel methods, which require  $\Theta(d^{2p})$  samples. Moreover, our analysis enables efficient  
 079       transfer learning for any other target function of the form  $f = g \circ p$  with a different link  
 080       function  $g$ , which also only requires  $\mathcal{O}(r^{\mathcal{O}(p)})$  additional samples.

### 081     1.2 RELATED WORKS

082     **Kernel Methods.** Earlier research links the behavior of gradient descent (GD) on the entire net-  
 083     work to its linear approximation near the initialization. In this scenario, neural networks act as  
 084     kernels, known as the Neural Tangent Kernel (NTK). This connection bridges neural network anal-  
 085     ysis with established kernel theory and offers initial learning guarantees for neural networks (Jacot  
 086     et al., 2018; Soltanolkotabi et al., 2018; Du et al., 2018; Chizat et al., 2019; Arora et al., 2019).  
 087     However, kernel theory fails to explain the superior empirical achievements of neural networks over  
 088     kernel methods (Arora et al., 2019; Lee et al., 2020; E et al., 2020). Networks in the kernel regime  
 089     **fail to learn features** (Yang and Hu, 2021), not adaptable to hierarchical structures of real world  
 090     targets. Ghorbani et al. (2021) proves that for uniformly distributed data on the sphere, the NTK  
 091     method requires  $\tilde{\Omega}(d^k)$  samples to learn any polynomials of degree  $k$  in  $d$  dimensions, which is  
 092     impractical when  $k$  is large. Thus, a central question is how neural networks can detect and capture  
 093     the underlying hierarchies in the target functions, which allows for a better generalization behavior  
 094     versus kernel methods.

095     **Learning Linear Features.** Recent studies have demonstrated neural networks' capability to  
 096     learn hierarchical functions of linear features more efficiently than kernel methods. Specifically,  
 097     Bietti et al. (2022); Ba et al. (2022) establish the efficient learning of single-index models, i.e.,  
 098      $f^*(\mathbf{x}) = g(\langle \mathbf{u}, \mathbf{x} \rangle)$ . Furthermore, recent works Damian et al. (2022); Abbe et al. (2023); Dandi  
 099     et al. (2023a); Bietti et al. (2023) further demonstrate that for isotropic data, two-layer or three-layer  
 100     neural networks can effectively learn multi-index models of the form  $f^*(\mathbf{x}) = g(\mathbf{U}\mathbf{x})$ . These studies  
 101     adopt certain modified training algorithms, such as layer-wise training. With sufficient feature learn-  
 102     ing, these networks can learn low-rank polynomials with a benign sample complexity of  $\mathcal{O}(d^{\mathcal{O}(1)})$ ,  
 103     which does not scale with the degree of the polynomial  $g$ . Empirically, fully connected networks  
 104     trained via gradient descent on image classification tasks also capture low-rank features (Lee et al.,  
 105     2007; Radhakrishnan et al., 2022). More recently, the learning of single-index and multi-index  
 106     models is analyzed with more advanced algorithm framework or specified data structure. Mousavi-  
 107     Hosseini et al. (2024) considers learning general multi-index models with two-layer neural networks  
 108     through a mean-field Langevin dynamics, Dandi et al. (2024b); Lee et al. (2024) goes beyond the  
 109     traditional Correlational Statistical Query (CSQ) setting and consider algorithms that reuse samples

for feature learning. Mousavi-Hosseini et al. (2023); Ba et al. (2023); Wang et al. (2024) considers learning linear features with structured data (such as data with a spiked covariance) rather than the commonly considered isotropic one. Cui et al. (2024); Dandi et al. (2024a) study the spectral structure revealed in the learned features with one huge gradient step through a spiked random feature model to understand the mechanism of feature learning in neural networks.

**Learning Nonlinear Features.** Previous studies indicate that neural networks can effectively learn specific hierarchies of nonlinear features. Safran and Lee (2022) shows that GD can efficiently learn functions such as  $1_{\|\mathbf{x}\| \geq \lambda}$  with a three-layer network. Ren et al. (2023) demonstrates that  $\text{ReLU}(1 - \|\mathbf{x}\|)$  can be learned by a multi-layer mean-field network. Moniri et al. (2024) studies the nonlinear feature learning capabilities of two-layer neural networks with one step of gradient descent. Allen-Zhu and Li (2019; 2020) explore learning target functions of the form  $p + \alpha g \circ p$  with  $p$  being the underlying feature through a three-layer residual network, though they either need  $\alpha = o_d(1)$  or cannot reach vanishing error. More recent works have addressed a broader class of nonlinear features compared with the previous research and demonstrate that three-layer neural networks can learn these hidden features efficiently. Specifically, Nichani et al. (2023) demonstrates that a three-layer network trained with layer-wise GD algorithm effectively learns  $g \circ p$  for a quadratic feature  $p(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$  with an improved sample complexity of  $\tilde{\Theta}(d^4)$ . As a follow-up, Wang et al. (2023) further demonstrates that such a network can in fact efficiently learn  $g \circ p$  for  $p$  within a broad subclass of degree  $k$  polynomials and optimizes the sample complexity to  $\tilde{\mathcal{O}}(d^k)$ .

However, all of these studies focus on a *single nonlinear feature*, limiting their applicability to scenarios involving multiple features. Our work addresses this gap by establishing the first theoretical guarantee for efficiently learning hierarchical polynomials of *multiple nonlinear features*, which significantly broadens the learnable function class and advances towards a better understanding of feature learning.

## 2 PRELIMINARIES

### 2.1 NOTATIONS

We use bold letters to denote vectors and matrices. For a vector  $\mathbf{v}$ , we denote its Euclidean norm by  $\|\mathbf{v}\|_2$ . For a matrix  $\mathbf{A}$ , we denote its operator and Frobenius norm as  $\|\mathbf{A}\|_2$  and  $\|\mathbf{A}\|_F$ , respectively. For any positive integer  $n$ , we denote  $[n] = \{1, 2, \dots, n\}$ . Moreover, for any indexes  $i$  and  $j$ , we denote  $\delta_{ij} = 1$  if  $i = j$  and 0 otherwise. We use  $\mathcal{O}$ ,  $\Theta$  and  $\Omega$  to hide absolute constants. In addition, we denote  $f \lesssim g$  when there exists some positive absolute constant  $C$  with  $f \leq Cg$ . We use  $\tilde{\mathcal{O}}$ ,  $\tilde{\Theta}$  and  $\tilde{\Omega}$  to ignore logarithmic terms. For a function  $f : \mathcal{X} \rightarrow \mathbb{R}$  and a distribution  $v$  on  $\mathcal{X}$ , we denote  $\|f\|_{L^p(\mathcal{X}, v)} = (\mathbb{E}_{\mathbf{x} \sim v} [|f(\mathbf{x})|^p])^{1/p}$ . When the domain is clear from context, we write  $\|f\|_{L^p}$  for simplicity. Finally, we write  $\mathbb{E}_{\mathbf{x}}$  as the shorthand for  $\mathbb{E}_{\mathbf{x} \sim v}$  sometimes.

### 2.2 PROBLEM SETUP

**Data distribution** Our aim is to learn the target function  $f^* : \mathcal{X} \rightarrow \mathbb{R}$ , with  $\mathcal{X} \subseteq \mathbb{R}^d$  being the input space. Throughout the paper, we assume  $\mathcal{X} = \mathbb{S}^{d-1}(\sqrt{d})$ , that is, the sphere with radius  $\sqrt{d}$  in  $d$  dimensions. Also, we consider the data distribution to be the uniform distribution on the sphere, i.e.,  $\mathbf{x} \sim \text{Unif}(\mathcal{X})$ , and we draw two independent datasets  $\mathcal{D}_1, \mathcal{D}_2$ , each with  $n_1$  and  $n_2$  i.i.d. samples, respectively. Thus, we draw  $n_1 + n_2$  samples in total.

**Target function** For the target function  $f^* : \mathbb{R}^d \rightarrow \mathbb{R}$ , we assume they are hierarchical functions of  $r$  quadratic features

$$f^*(\mathbf{x}) = g^*(\mathbf{p}(\mathbf{x})) = g^*(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \mathbf{x}^\top \mathbf{A}_2 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}).$$

This structure represents a broad class of functions where  $\mathbf{p}(\mathbf{x}) = [\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \mathbf{x}^\top \mathbf{A}_2 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}]^\top$  represents  $r$  quadratic features, and  $g^* : \mathbb{R}^r \rightarrow \mathbb{R}$  is a link function. Here we consider the case  $r \ll d$ . To simplify our analysis while maintaining generality, we make the following assumptions:

**Assumption 1** (Orthogonal quadratic features). *For any  $i, j \in [r]$ , we suppose*

$$\mathbb{E}_{\mathbf{x}} [\mathbf{x}^\top \mathbf{A}_i \mathbf{x}] = 0, \quad \mathbb{E}_{\mathbf{x}} [(\mathbf{x}^\top \mathbf{A}_i \mathbf{x})(\mathbf{x}^\top \mathbf{A}_j \mathbf{x})] = \delta_{ij} \quad \text{and} \quad \|\mathbf{A}_i\|_{op} \leq \frac{\kappa_1}{\sqrt{d}}.$$

Here we assume  $\kappa_1 = \text{poly}(\log d)$ .

The first assumption is equivalent to  $\text{tr}(\mathbf{A}_i) = 0$  for any  $i \in [r]$ . For  $\mathbf{A}_i$  such that  $\text{tr}(\mathbf{A}_i) \neq 0$ , we could simply subtract the mean of the feature to  $\mathbf{A}'_i = \mathbf{A}_i - (\text{tr}(\mathbf{A}_i)/d) \cdot \mathbf{I}$  so

$$\mathbf{x}^\top \mathbf{A}'_i \mathbf{x} = \mathbf{x}^\top (\mathbf{A}_i - (\text{tr}(\mathbf{A}_i)/d) \cdot \mathbf{I}) \mathbf{x} = \mathbf{x}^\top \mathbf{A}_i \mathbf{x} - \text{tr}(\mathbf{A}_i).$$

The second assumption on the feature orthonormality can be attained via linear transformation on the features, preserving the overall function class. The third assumption on the operator norm bound ensures that the features are balanced, which is common in the non-linear feature learning literature (Nichani et al., 2023; Wang et al., 2023). Moreover, we note that when the entries of  $\mathbf{A}_i$  are sampled i.i.d., the assumption is satisfied with high probability by standard random matrix arguments.

**Assumption 2** (Well-conditioned link function). *For the link function  $g^*$ , we assume  $g^*$  is a degree- $p$  polynomial with  $\mathbb{E}_{\mathbf{z}} [g^*(\mathbf{z})^2] = \Theta(1)$ , where  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_r)$  and  $p \in \mathbb{N}$  is a constant. Moreover, we assume the expected Hessian  $\mathbf{H} = \mathbb{E}_{\mathbf{z}} [\nabla^2 g^*(\mathbf{z})] \in \mathbb{R}^{r \times r}$  is well-conditioned, i.e., there exists a constant  $C_H$  such that  $\lambda_{\min}(\mathbf{H}) \geq \frac{C_H}{\sqrt{r}}$ .*

This assumption ensures the link function adequately emphasizes all  $r$  features, preventing degeneracy to a lower-dimensional subspace. The second-moment condition is achievable through simple normalization.

**Assumption 3** (Preprocessed target function). *For the entire target function  $f^*$ , we assume  $\mathcal{P}_0(f^*) = \mathbb{E}_{\mathbf{x}}[f^*(\mathbf{x})] = 0$  and  $\|\mathcal{P}_2(f^*)\|_{L^2} \leq \kappa_2/\sqrt{d}$ . Here  $\mathcal{P}_k$  is the projection onto the function space of degree  $k$  spherical harmonics on the sphere  $\mathbb{S}^{d-1}(\sqrt{d})$ , and  $\kappa_2$  satisfies  $\kappa_2 = \text{poly}(r, \log d)$ .*

We will give a rigorous definition of  $\mathcal{P}_k$  in Section 2.3.1. This assumption is analogous to a pre-processing procedure conducted in Damian et al. (2022), which subtracts out the mean and linear component of the features from the target. The zero-mean condition ensures the network focuses on learning the function's variability rather than a constant offset. While Nichani et al. (2023); Wang et al. (2023) assume the link function  $g$  has non-zero linear component, we rather assume  $g$  has a *nearly zero linear component*, which prevents the target function from being dominated by a single linear combination of the quadratic features and keeps the learned representation space from collapsing to the one-dimensional space of that certain linear combination. This is an essential difference between single-feature and multi-feature learning, because our assumptions ensure that the network genuinely learns to *represent and distinguish all  $r$  features* rather than conflate them, while assumptions in Nichani et al. (2023); Wang et al. (2023) represent a *degenerate case* that neural network may only learn the dominant linear combination of the  $r$  features. We provide examples and counterexamples as follows.

**Remark 1.** These assumptions accommodate a wide range of target functions. For instance,  $f^*(\mathbf{x}) = \frac{1}{\sqrt{r}} \sum_{k=1}^r (\mathbf{x}^\top \mathbf{A}_k \mathbf{x})^2 - \sqrt{r}$  satisfies Assumption 3 with  $\kappa_2 \lesssim \sqrt{r}\kappa_1$  for any  $\{\mathbf{a}_k\}_{k \in [r]}$  under Assumption 1. Moreover, for diagonal  $\mathbf{A}_k$  with  $\mathbf{A}_k = \text{diag}(\mathbf{a}_k)$ , where  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_r$  are orthogonal zero-sum vectors with entries  $a_{k,i} \in \{\pm c/\sqrt{d}\}$ , we can achieve  $\kappa_2 = 0$ . Here  $c = \Theta(1)$  is a normalizing constant. Notably, linear combinations of features like  $f(\mathbf{x}) = \frac{1}{\sqrt{r}} \sum_{k=1}^r (\mathbf{x}^\top \mathbf{A}_k \mathbf{x})$  violate our assumptions, since it represents a degenerate case with  $\|\mathcal{P}_2(f)\|_{L^2} = \|f\|_{L^2} = \Theta(1)$ .

**Three-layer neural network** We adopt a standard three-layer neural network for learning the target functions. Let  $m_1, m_2$  be the two hidden layer widths, and  $\sigma_1, \sigma_2$  be two activation functions. Our learner is a three-layer neural network parameterized by  $\theta = (\mathbf{a}, \mathbf{W}, \mathbf{b}, \mathbf{V})$ , where  $\mathbf{a} \in \mathbb{R}^{m_1}$ ,  $\mathbf{W} \in \mathbb{R}^{m_1 \times m_2}$ ,  $\mathbf{b} \in \mathbb{R}^{m_1}$ , and  $\mathbf{V} \in \mathbb{R}^{m_2 \times d}$ . The network  $f(\mathbf{x}; \theta)$  is defined as

$$f(\mathbf{x}; \theta) = \frac{1}{m_1} \sum_{i=1}^{m_1} a_i \sigma_1(\langle \mathbf{w}_i, \sigma_2(\mathbf{V}\mathbf{x}) \rangle + b_i) = \frac{1}{m_1} \sum_{i=1}^{m_1} a_i \sigma_1(\langle \mathbf{w}_i, \mathbf{h}^{(0)}(\mathbf{x}) \rangle + b_i). \quad (1)$$

Here,  $\mathbf{w}_i \in \mathbb{R}^{m_2}$  is the  $i$ -th row of  $\mathbf{W}$ , and  $\mathbf{h}^{(0)}(\mathbf{x}) := \sigma_2(\mathbf{V}\mathbf{x}) \in \mathbb{R}^{m_2}$  is the random feature embedding lying in the innermost layer. We initialize each row of  $\mathbf{V}$  to be drawn uniformly on the sphere of radius  $\sqrt{d}$ , i.e.,  $\mathbf{v}_i^{(0)} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ . For  $\mathbf{a}$ ,  $\mathbf{b}$  and  $\mathbf{W}$ , we use a symmetric initialization so that  $f(\mathbf{x}; \theta^{(0)}) = 0$  (Chizat et al., 2019). Explicitly, we assume that  $m_1$  is an even number and for any  $j \in [m_1/2]$ , we initialize the parameters as

$$a_j^{(0)} = -a_{m_1-j}^{(0)} \sim \text{Unif}(\{-1, 1\}), \quad \mathbf{w}_j^{(0)} = \mathbf{w}_{m_1-j}^{(0)} \sim \mathcal{N}(0, \epsilon \mathbf{I}_{m_2}), \quad \text{and } b_j^{(0)} = b_{m_1-j}^{(0)} = 0.$$

216 Here  $\epsilon > 0$  is a hyperparameter to control the magnitude of the initial neurons. Different from  
 217 Nichani et al. (2023) where the weights  $w_j$  are initialized at zeros, we require a *random initialization*, which enables the learned weights to capture the multiple features in all directions instead of  
 218 converging to a specific direction like the previous results for learning a single feature.  
 219

220 For the activation functions  $\sigma_1$  and  $\sigma_2$ , we have the following assumptions:  
 221

222 **Assumption 4** (Activation Function). *We take the outer activation function  $\sigma_1$  and the inner activation function  $\sigma_2$  as*  
 223

$$\sigma_1(t) = \begin{cases} 2|t| - 1, & |t| \geq 1, \\ t^2, & |t| < 1. \end{cases} \quad \text{and} \quad \sigma_2(t) = \sum_{i=2}^{\infty} c_i Q_i(t), \quad (2)$$

224 where  $Q_i(t)$  is the  $i$ -th degree Gegenbauer polynomial in the  $d$ -dimensional space. Moreover, we  
 225 assume there exist constants  $C_\sigma, \alpha_\sigma$  such that  $|\sigma_2(t)| \leq C_\sigma$  for  $|t| \leq d$ , and  $\mathbb{E}_{\mathbf{x}} [\sigma_2^k(\mathbf{x}^\top \mathbf{1}_d)] \leq$   
 226  $d^{-k} C_k$  for  $k = 2, 4$ . We assume  $c_2 = \Theta(1)$ , and  $C_2, C_4$  and  $\{c_i\}_{i=2}^{\infty}$  are all constants independent  
 227 of  $n, d, m_1$  and  $m_2$ .

228 We remark the outer activation  $\sigma_1$  is a slightly modified version of the absolute value function  
 229  $|t|$ , smoothed around the origin. The assumptions on  $\sigma_2$  are based on the Gegenbauer expansion,  
 230 often considered in the spherical analysis (introduced in Section 2.3.2). Compared to standard inner  
 231 activations, we remove the constant term ( $Q_0(t) = 1$ ) and the linear term ( $Q_1(t) = t/d$ ) to focus  
 232 on learning nonlinear features without low-order interference. Importantly, these assumptions on  
 233 activation functions maintain significant generality. The assumptions on magnitude and moments  
 234 are satisfied by many common activation functions with appropriate scaling. The core assumption  
 235 in the Gegenbauer expansion is the non-zero component of  $Q_2$ , i.e.,  $c_2 = \Theta(1)$ , which we rely on  
 236 for a subspace recovery of the  $r$  quadratic features while other assumptions are made to simplify  
 237 our analysis since other components in inner activation will lead to useless noises or biases in the  
 238 weights after training. Moreover, if we consider higher degree nonlinear features such as degree  $q$   
 239 polynomials, we expect that  $\sigma_2$  has sufficient emphasis on  $Q_q$  for efficient feature learning.  
 240

241 **Remark 2.**  $\sigma_2(t) = Q_2(t) = \frac{t^2-d}{d(d-1)}$  is an example of the inner activation function.  
 242

243 **Training Algorithm** Following Nichani et al. (2023), our network is trained via layer-wise gra-  
 244 dient descent with sample splitting. Throughout the training process, we freeze the innermost layer  
 245 weights  $\mathbf{V}$ . In the first stage, the second layer weights  $\mathbf{W}$  are trained for one step with a specified  
 246 learning rate  $\eta_1$  and weight decay  $\lambda_1$ . In the second stage, we reinitialize the bias  $\mathbf{b}$  and train the  
 247 outer layer weights  $\mathbf{a}$  for  $T - 1$  steps.  
 248

249 **Transfer Learning** We remark that our algorithm allows transfer learning of a different target  
 250 function  $f$  that shares the same features of the original target:  
 251

$$f^*(\mathbf{x}) \rightarrow f(\mathbf{x}) = g(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \mathbf{x}^\top \mathbf{A}_2 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}) \quad (\text{transferred target})$$

252 In this case, we switch the target function from  $f^* = g^*(\mathbf{p})$  to  $f = g(\mathbf{p})$  in the second training  
 253 stage. For the loss function, we use the standard squared loss:  
 254

$$\hat{\mathcal{L}}^{(1)}(\theta) = \frac{1}{n_1} \sum_{\mathbf{x} \in \mathcal{D}_1} (f(\mathbf{x}; \theta) - f^*(\mathbf{x}))^2, \quad \hat{\mathcal{L}}^{(2)}(\theta) = \begin{cases} \frac{1}{n_2} \sum_{\mathbf{x} \in \mathcal{D}_2} (f(\mathbf{x}; \theta) - f^*(\mathbf{x}))^2 & (\text{original}), \\ \frac{1}{n_2} \sum_{\mathbf{x} \in \mathcal{D}_2} (f(\mathbf{x}; \theta) - f(\mathbf{x}))^2 & (\text{transferred}). \end{cases}$$

255 This layer-wise training approach, combined with the ability to perform transfer learning, provides a  
 256 powerful framework for learning and adapting to hierarchical functions with hidden features (Kulka-  
 257 rni and Karande, 2017; Damian et al., 2022; Nichani et al., 2023). The pseudocode for the entire  
 258 training procedure is presented in Algorithm 1.  
 259

### 2.3 TECHNICAL BACKGROUND: ANALYSIS OVER THE SPHERE

260 We briefly introduce spherical harmonics and Gegenbauer polynomials, which forms the foundation  
 261 of our analysis over the sphere  $\mathbb{S}^{d-1}(\sqrt{d})$ . For more details, see Appendix B.5.  
 262

#### 2.3.1 SPHERICAL HARMONICS

263 Let  $\tau_{d-1}$  be the uniform distribution on  $\mathbb{S}^{d-1}(\sqrt{d})$ . Consider functions in  $L^2(\mathbb{S}^{d-1}(\sqrt{d}), \tau_{d-1})$ , with  
 264 scalar product and norm denoted as  $\langle \cdot, \cdot \rangle_{L^2}$  and  $\|\cdot\|_{L^2}$ . For  $\ell \in \mathbb{Z}_{\geq 0}$ , let  $V_{d,\ell}$  be the linear space of  
 265

---

270   **Algorithm 1** Layer-wise training algorithm  
271   **Input:** Learning rates  $\eta_1, \eta_2$ , weight decay  $\lambda_1, \lambda_2$ , parameter  $\epsilon$ , number of steps  $T$   
272   1 **initialize**  $\mathbf{a}, \mathbf{b}, \mathbf{W}$  and  $\mathbf{V}$ .  
273   2 **train**  $\mathbf{W}$  on dataset  $\mathcal{D}_1$   
274   3     $\mathbf{W}^{(1)} \leftarrow \mathbf{W}^{(0)} - \eta_1 [\nabla_{\mathbf{W}} \hat{\mathcal{L}}^{(1)}(\theta) + \lambda_1 \mathbf{W}^{(0)}]$   
275   4 **end**  
276   5 **re-initialize**  
277   6     $b_i^{(1)} \sim \text{Unif}([-3, 3]), i \in [m_1]$   
278   7     $\mathbf{a}^{(1)}, \mathbf{V}^{(1)} \leftarrow \mathbf{a}^{(0)}, \mathbf{V}^{(0)}$   
279   8     $\theta^{(1)} \leftarrow (\mathbf{a}^{(1)}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{V}^{(0)})$   
280   7 **end**  
281   8 **train**  $\mathbf{a}$  on dataset  $\mathcal{D}_2$   
282   9    **for**  $t = 2$  to  $T$  **do**  
283   10     $\mathbf{a}^{(t)} \leftarrow \mathbf{a}^{(t-1)} - \eta_2 [\nabla_{\mathbf{a}} \hat{\mathcal{L}}^{(2)}(\theta^{(t-1)}) + \lambda_2 \mathbf{a}^{(t-1)}]$   
284   11    **end**  
285   12 **end**  
286 13 **return** Prediction function  $f(\cdot; \theta^{(T)})$ :  $\mathbf{x} \rightarrow \frac{1}{m_1} \langle \mathbf{a}^{(T)}, \sigma_1(\mathbf{W}^{(1)} \mathbf{h}^{(0)}(\mathbf{x}) + \mathbf{b}^{(1)}) \rangle$

---

288  
289 homogeneous harmonic polynomials of degree  $\ell$  restricted on  $\mathbb{S}^{d-1}(\sqrt{d})$ . The set  $\{V_{d,\ell}\}_{\ell \geq 0}$  forms  
290 an orthogonal basis of the  $L^2$  space, with dimension  $\dim(V_{d,\ell}) = \Theta(d^\ell)$ . For each  $\ell \in \mathbb{Z}_{\geq 0}$ , the  
291 spherical harmonics  $\{Y_{\ell,j}\}_{j \in [B(d,\ell)]}$  form an orthonormal basis of  $V_{d,\ell}$ . Moreover, we denote by  $\mathcal{P}_k$   
292 the orthogonal projections to  $V_{d,k}$ , which can be written as

293  
294   
$$\mathcal{P}_k(f)(\mathbf{x}) = \sum_{\ell=1}^{B(d,k)} \langle f, Y_{k,\ell} \rangle_{L^2} Y_{k,\ell}(\mathbf{x}).$$
  
295  
296

297 We also define  $\mathcal{P}_{\leq \ell} \equiv \sum_{k=0}^{\ell} \mathcal{P}_k$ ,  $\mathcal{P}_{>\ell} \equiv \mathbf{I} - \mathcal{P}_{\leq \ell}$ ,  $\mathcal{P}_{<\ell} \equiv \mathcal{P}_{\leq \ell-1}$ , and  $\mathcal{P}_{\geq \ell} \equiv \mathcal{P}_{>\ell-1}$ .

### 2.3.2 GEGENBAUER POLYNOMIALS

300 Corresponding to the degree  $\ell$  spherical harmonics in the  $d$ -dimension space, the  $\ell$ -th Gegenbauer  
301 polynomial  $Q_\ell : [-d, d] \rightarrow \mathbb{R}$  is a polynomial of degree  $\ell$ . The set  $\{Q_\ell\}_{\ell \geq 0}$  forms an orthogonal ba-  
302 sis on  $L^2([-d, d], \tilde{\tau}_{d-1})$ , where  $\tilde{\tau}_{d-1}$  is the distribution of  $\sqrt{d} \langle \mathbf{x}, \mathbf{e}_1 \rangle$  when  $\mathbf{x} \sim \tau_{d-1}$ . In particular,  
303 these polynomials are normalized so that  $Q_\ell(d) = 1$ . We present the explicit forms of Gegenbauer  
304 polynomials of degree no more than 2:

305   
$$Q_0(t) = 1, \quad Q_1(t) = \frac{t}{d}, \quad \text{and} \quad Q_2(t) = \frac{t^2 - d}{d(d-1)}.$$
  
306  
307

308 Gegenbauer polynomials are directly related to spherical harmonics, leading to a number of elegant  
309 properties. We provide further details on these properties in Appendix B.5.

## 3 MAIN RESULTS

310 The following is our main theorem, which bounds the population absolute loss of Algorithm 1:

311 **Theorem 1.** Suppose  $n_1, m_2 = \tilde{\Omega}(d^4)$ . Let  $\hat{\theta}$  be the output of Algorithm 1 after  $T =$   
312  $\text{poly}(n_1, n_2, m_1, m_2, d)$  steps. Then, there exists a set of hyper-parameters  $(\epsilon, \eta_1, \eta_2, \lambda_1, \lambda_2)$  such  
313 that, with high probability over the initialization of parameters and draws of  $\mathcal{D}_1, \mathcal{D}_2$ , we have

314  
315   
$$\mathbb{E}_{\mathbf{x}} \left[ |f(\mathbf{x}; \hat{\theta}) - f^*(\mathbf{x})| \right] = \tilde{\mathcal{O}} \left( \underbrace{\sqrt{\frac{r^p \kappa_2^{2p}}{\min(m_1, n_2)}}}_{\text{Complexity of } g^*} + \underbrace{\sqrt{\frac{d^6 r^{p+1}}{m_2}} + \sqrt{\frac{d^2 r^{p+1}}{n_1}} + \frac{r^{p+2}}{d^{1/6}}}_{\text{Feature Learning Error}} \right).$$
  
316  
317  
318  
319  
320  
321

322 Moreover, for any other degree  $p$  polynomial  $g : \mathbb{R}^r \rightarrow \mathbb{R}$  with  $\|g\|_{L^2} \lesssim 1$ , by substituting the target  
323 function  $f^* = g^* \circ \mathbf{p}$  by  $f = g \circ \mathbf{p}$  in the second training stage, we can achieve the same result for  
learning the new target function.

324 The full proof is provided in Appendix E.1. To interpret the results, we provide the following  
 325 discussion of Theorem 1.

327 **Feature learning error** This term quantifies the requirements on the first-stage sample complexity  
 328 and the inner width to sufficiently capture the non-linear features. Given  $d \gg r$ , if the width  
 329  $m_2 = \tilde{\Omega}(d^6 r^{p+1})$  and the sample size  $n_1 = \tilde{\Omega}(d^4 + d^2 r^{p+1})$ , we can fully capture the underlying  
 330 feature information and approximate any degree  $p$  polynomials of the features. We will demon-  
 331 strate how Algorithm 1 learns these features through the learned representations in Proposition 1  
 332 and express hierarchical polynomials in Proposition 2.

334 **Complexity of  $g^*$**  This term is the second-stage sample (and width) complexity given that the  $r$   
 335 features have been fully captured in the first stage. Moreover, for a sufficiently preprocessed target  
 336 function, i.e.,  $\kappa_2 = \mathcal{O}(1)$ , we achieve the standard results of  $\tilde{\mathcal{O}}(r^p)$  complexity in learning a degree-  
 337  $p$  polynomial in the  $r$ -dimensional space in the kernel regime.

339 **Transfer learning** Leveraging the two-stage structure of training, we can learn a different target  
 340 function in the second stage that shares the same features with the original target. This also sup-  
 341 ports the fact that we have fully captured the information of the  $r$  nonlinear features in the first  
 342 stage, making it possible for the efficient learning with a different polynomial head  $g$ . Moreover,  
 343 by viewing the first stage as a pre-training process with  $\tilde{\Omega}(d^4 + d^2 r^{p+1})$  samples, only additional  
 344  $\tilde{\mathcal{O}}(r^p \kappa_2^{2p})$  samples are required to learn any degree  $p$  polynomial of the features, which gets rid of  
 345 the polynomial dependence on the ambient dimension of  $d$ .

347 **Comparison with previous works** Compared with the sample complexity of  $\tilde{\Omega}(d^2 r + dr^p)$  in  
 348 Damian et al. (2022) for learning multi-index models, we have a similar polynomial dependence on  
 349  $r$ , and the dependence on  $d$  increases from  $d^2$  to  $d^4$  because of the increased complexity of quadratic  
 350 features rather than linear ones. Moreover, our approach significantly improves upon the  $\Theta(d^{2p})$   
 351 sample complexity required by kernel methods to learn degree  $p$  polynomials of quadratic features  
 352 (i.e., degree  $2p$  polynomials of the input). Crucially, our polynomial dependence on  $d$  in the overall  
 353 sample complexity is independent of the degree  $p$  of the link function  $g$ .

355 **Near optimality of the sample complexity** We remark that our sample complexity of  $\tilde{\mathcal{O}}(d^4)$  is  
 356 nearly optimal with respect to  $d$  for all algorithms that use one step of gradient descent for feature  
 357 learning. Our assumptions on the target functions imply that the leap index<sup>1</sup> of our target functions  
 358 are basically 4 (more specifically, the second order information of  $g \circ p$ , where  $p$  are quadratic  
 359 features), and we also utilize  $\mathcal{P}_4(f)$  for recovering the subspace of the  $r$  quadratic features, which  
 360 will be discussed in details in Section 4. Dandi et al. (2023b) indicates that  $\Omega(d^4)$  samples are  
 361 required for an efficient learning of terms in  $\mathcal{P}_4(f^*)$ , which substantiates the near optimality of our  
 362 result.

## 363 4 PROOF ROADMAP OF THEOREM 1

365 The proof of Theorem 1 unfolds in two training stages. First, by a novel universality argument  
 366 (Lemma 1), we show that after the first training stage, with sufficient training samples, the net-  
 367 work learns to fully extract out the hidden features  $p$  (Proposition 1). Next, we show that during  
 368 the second stage, the network is capable of expressing the link function with a mild outer width  
 369  $m_1$  (Proposition 2). We conclude the proof through standard Rademacher complexity analysis to  
 370 quantify the generalization error of the second-stage model (detailed in Appendix E.1).

### 371 4.1 STAGE 1: LEARNING THE FEATURES

372 We provide a brief analysis on the learned representations after the first training stage. Denote  
 373  $\mathbf{w}_j = \epsilon^{-1} \mathbf{w}_j^{(0)} \sim \mathcal{N}(0, \mathbf{I}_{m_2})$ . According to Algorithm 1, by setting  $\epsilon$  sufficiently small, after

375  
 376  
 377 <sup>1</sup>The leap index of a target function  $f^*$  is the first integer  $\ell$  that  $\mathcal{P}_\ell f^* \neq 0$ . Our assumptions imply a  
 diminishing  $\mathcal{P}_{<4}(f^*)$  and a non-degenerate  $\mathcal{P}_4(f^*)$  as  $d \rightarrow \infty$ .

378 one-step gradient descent on  $\mathbf{W}$ , we know for each  $j \in [m_1]$ ,

$$\begin{aligned} 380 \quad \eta_1 \nabla_{\mathbf{w}_j^{(0)}} \mathcal{L}(\theta^{(0)}) &= -\eta_1 \frac{a_j^{(0)}}{m_1} \cdot \frac{1}{n_1} \sum_{\mathbf{x} \in \mathcal{D}_1} f^*(\mathbf{x}_i) \mathbf{h}^{(0)}(\mathbf{x}_i) \sigma'_1 \left( \langle \epsilon \mathbf{w}_j, \mathbf{h}^{(0)}(\mathbf{x}_i) \rangle \right) \\ 381 \quad &\xrightarrow{\epsilon \rightarrow 0} -\frac{2\epsilon\eta_1}{m_1} a_j^{(0)} \cdot \frac{1}{n_1} \sum_{\mathbf{x} \in \mathcal{D}_1} f^*(\mathbf{x}_i) \mathbf{h}^{(0)}(\mathbf{x}_i) \mathbf{h}^{(0)}(\mathbf{x}_i)^\top \mathbf{w}_j. \\ 382 \end{aligned}$$

383 By taking  $\eta_1 = \frac{m_1}{2\epsilon m_2} \cdot \eta$  for some  $\eta > 0$  to be chosen later and  $\lambda_1 = \eta_1^{-1}$ , we have

$$388 \quad \mathbf{w}_j^{(1)} = \mathbf{w}_j^{(0)} - \eta_1 \left[ \nabla_{\mathbf{w}_j^{(0)}} \mathcal{L}(\theta^{(0)}) + \lambda_1 \mathbf{w}_j^{(0)} \right] = \frac{\eta a_j^{(0)}}{m_2} \cdot \frac{1}{n_1} \sum_{\mathbf{x} \in \mathcal{D}_1} f^*(\mathbf{x}_i) \mathbf{h}^{(0)}(\mathbf{x}_i) \mathbf{h}^{(0)}(\mathbf{x}_i)^\top \mathbf{w}_j.$$

389 Then for any second-stage training sample  $\mathbf{x}' \in \mathcal{D}_2$ , the inner-layer representation becomes

$$\begin{aligned} 392 \quad \left\langle \mathbf{w}_j^{(1)}, \sigma_2(\mathbf{V}\mathbf{x}') \right\rangle &= \frac{\eta a_j^{(0)}}{m_2} \left\langle \frac{1}{n_1} \sum_{\mathbf{x} \in \mathcal{D}_1} f^*(\mathbf{x}_i) \mathbf{h}^{(0)}(\mathbf{x}_i) \mathbf{h}^{(0)}(\mathbf{x}_i)^\top \mathbf{w}_j, \mathbf{h}^{(0)}(\mathbf{x}') \right\rangle \\ 393 \quad &= \eta a_j^{(0)} \cdot \underbrace{\left\langle \mathbf{w}_j, \frac{1}{n_1 m_2} \sum_{\mathbf{x} \in \mathcal{D}_1} f^*(\mathbf{x}_i) \langle \mathbf{h}^{(0)}(\mathbf{x}_i), \mathbf{h}^{(0)}(\mathbf{x}') \rangle \mathbf{h}^{(0)}(\mathbf{x}_i) \right\rangle}_{\mathbf{h}^{(1)}(\mathbf{x}')}. \\ 394 \end{aligned}$$

400 Our main contribution in this part is that the first-step trained presentations representations  $\mathbf{h}^{(1)}(\mathbf{x})$   
401 approximately spans the space of the target features  $(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \mathbf{x}^\top \mathbf{A}_2 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x})$ . Thus, the  
402 target features  $\mathbf{p}(\mathbf{x})$  can be reconstructed through a linear transformation from the learned representations  
403  $\mathbf{h}^{(1)}(\mathbf{x})$ , which is formalized in the following proposition.

404 **Proposition 1** (Reconstruct the feature). *Suppose  $m_2, n_1 = \tilde{\Omega}(d^4)$ . With high probability jointly on  
405  $\mathbf{V}$ ,  $\mathcal{D}_1$  and  $\mathcal{D}_2$ , there exists a matrix  $\mathbf{B}^* \in \mathbb{R}^{r \times m_2}$  such that for any  $\mathbf{x} \in \mathcal{D}_2$ , we have*

$$407 \quad \left\| \mathbf{B}^* \mathbf{h}^{(1)}(\mathbf{x}) - \mathbf{p}(\mathbf{x}) \right\|_2 = \tilde{\mathcal{O}} \left( \frac{d^3 r}{\sqrt{m_2}} + \frac{dr}{\sqrt{n_1}} + \frac{r^{\frac{p+5}{2}}}{d^{1/6}} \right). \quad (3)$$

410 The proof is provided in Appendix C.3. We summarize the main idea of the proof as follows.

411 **Universality of features** The foundation of the proof lies in the universality result that the joint  
412 distribution of the multiple features  $\mathbf{p}$  is approximately multivariate standard Gaussian:

$$414 \quad (\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \mathbf{x}^\top \mathbf{A}_2 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}) \stackrel{d}{\approx} \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r), \quad d \gg r.$$

416 It is worth mentioning that we provide a general universality theory that quantifies the difference  
417 between the distribution of any  $r$ -dimensional function (not limited in quadratic forms) and the  
418  $r$ -dimensional Gaussian distribution, which is presented in Lemma 1.

419 **Lemma 1** (Universality of vector-valued functions). *Suppose  $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  is an  $d$ -dimensional  
420 standard Gaussian variable. If a function  $\mathbf{p} : \mathbb{R}^d \rightarrow \mathbb{R}^r$  satisfies  $\mathbb{E}_{\mathbf{X}}[\mathbf{p}(\mathbf{X})] = \mathbf{0}_r$  and  
421  $\text{Cov}(\mathbf{p}(\mathbf{X}), \mathbf{p}(\mathbf{X})) = \mathbf{I}_r$ , then we have*

$$422 \quad W_1(\text{Law}(\mathbf{p}(\mathbf{X})), \mathcal{N}(\mathbf{0}, \mathbf{I}_r)) \leq \frac{4}{\sqrt{\pi}} \left( \sum_{i=1}^r \mathbb{E} \left[ \|\nabla p_i(\mathbf{X})\|_2^4 \right]^{1/4} \right) \left( \sum_{j=1}^r \mathbb{E} \left[ \|\nabla^2 p_j(\mathbf{X})\|_{\text{op}}^4 \right]^{1/4} \right).$$

425 Here  $\mathbf{p}(\mathbf{x}) = [p_1(\mathbf{x}), p_2(\mathbf{x}), \dots, p_r(\mathbf{x})]^\top$  and  $W_1$  denotes the Wasserstein-1 distance.

427 The proof is provided in Appendix B.2. This lemma extends the previous universality results of  
428 univariate Gaussian approximation theory (Chatterjee, 2007) to the multivariate version and could  
429 be of independent interest for the field of high dimensional probability theory. As a corollary, when  
430 we take  $\mathbf{p}$  to be  $r$  quadratic features satisfying Assumption 1, we ensure the  $W_1$  distance is bounded  
431 by  $\tilde{\mathcal{O}}(r^2/\sqrt{d})$  (see Lemma 16 in the appendix for the formal statement). This approximation error  
finally contributes to third term in the error bound of Proposition 1 (Equation (3)).

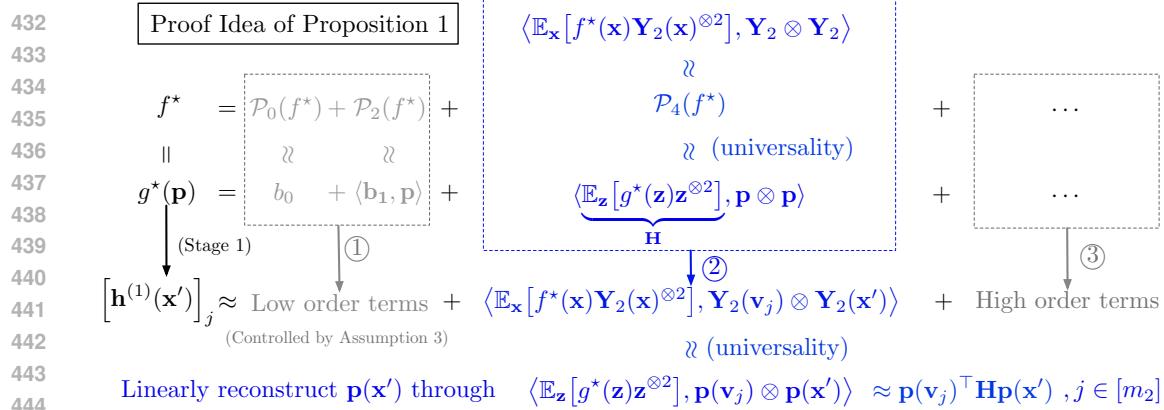


Figure 1: The proof idea of Proposition 1. Block 1 characterizes the constant and linear terms of  $g^*$ , which is approximately equivalent to the low-order terms  $\mathcal{P}_{<4}(f^*)$  by our universality theory and results into biases in the learned weights  $\mathbf{h}^{(1)}(\mathbf{x}')$  after Stage 1. This bias is vanishing with  $d \rightarrow \infty$  by our assumptions on  $\mathcal{P}_0(f^*)$  and  $\mathcal{P}_2(f^*)$ . Block 2 describes the second-order information of  $g^*$  (approximately  $\mathcal{P}_4(f^*)$ ), which is of the greatest importance and captured by the quadratic component  $c_2 Q_2(\cdot)$  in the inner activation  $\sigma_2(\cdot)$  and converted into quantities spanned by the  $r$  quadratic features  $\mathbf{p}$ . Block 3 represents the remaining terms of  $f^*$ , which leads to high-order nuisance in the learned weights, but still dominated by the second term due to Assumption 2 when  $d$  is large, which enables us to utilize the terms in blue (resulted from Block 2) to reconstruct the features efficiently.

**Utilizing the second-order information of  $g^*$**  Lemma 1 establishes a crucial link between our model and the multi-index model studied by Damian et al. (2022). This connection allows us to simplify the analysis on non-linear features and utilize the second-order information of the link function  $g^*$  to fully recover the feature space. In the context of multi-index models where  $f^*(\mathbf{x}) = g^*(\mathbf{p}(\mathbf{x}))$  with  $\mathbf{p}(\mathbf{x}) = \mathbf{U}\mathbf{x}$ , it has been shown that for a prepossessed target with a non-degenerate expected Hessian  $\mathbf{H} = \mathbb{E}_{\mathbf{z}}[\nabla^2 g^*(\mathbf{z})]$ , the learned representations, dominated by the degree 2 component of  $f^*$  which takes form  $\mathbb{E}_{\mathbf{x}}[f^*(\mathbf{x})\mathbf{x}^{\otimes 2}] \approx \mathbf{U}^T \mathbf{H} \mathbf{U}$ , are spanned by  $\{\mathbf{u}_i \otimes \mathbf{u}_j\}_{i,j \in [r]}$ . Extending this to our setting with quadratic features and applying the universality argument from Lemma 1, we demonstrate that the degree 4 component of our  $f^*$ , namely  $\mathbb{E}_{\mathbf{x}}[f^*(\mathbf{x})\mathbf{Y}_2(\mathbf{x})^{\otimes 2}]$ , is approximately spanned by the quantities  $\{\mathbf{A}_i \otimes \mathbf{A}_j\}_{i,j \in [r]}$ , which is formalized in Proposition 3 in Appendix C.1. Here  $\mathbf{Y}_2(\mathbf{x})$  represents the tensorized quadratic spherical harmonics. Under Assumption 3, it turns out that after the first step of GD (Stage 1 of Algorithm 1), the learned representations are dominated by this degree 4 component (Proposition 4 in Appendix C.2). This domination enables efficient recovery of the "span" of the hidden features  $\mathbf{p}$ . For a visual representation of our proof strategy, we also present our main idea of the proof in Figure 1. Remarkably, we find that the reconstruction matrix admits a surprisingly simple form of  $\mathbf{B}^* \propto \mathbf{H}^{-1}[\mathbf{p}(\mathbf{v}_1), \mathbf{p}(\mathbf{v}_2), \dots, \mathbf{p}(\mathbf{v}_{m_2})]$ . We provide empirical support for the effectiveness of this reconstruction through experiments in Section A.

#### 4.2 STAGE 2: LEARNING THE LINK FUNCTION

By the deduction above, after the first training stage, the model becomes a random-feature model (Rahimi and Recht, 2007):

$$f(\mathbf{x}'; \theta) = \frac{1}{m_1} \sum_{j=1}^{m_1} a_j \sigma_1 \left( \eta a_j^{(0)} \langle \mathbf{w}_j, \mathbf{h}^{(1)}(\mathbf{x}') \rangle + b_j^{(1)} \right). \quad (4)$$

Here  $\theta = (\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{V})$ , with  $\mathbf{a} = [a_1, a_2, \dots, a_{m_1}]^\top \in \mathbb{R}^{m_1}$  being the trainable parameters in the second stage. Leveraging the construction in Proposition 1, we can construct a corresponding weight vector  $\mathbf{a}$  in the outer layer to express the polynomial  $g(\mathbf{B}^* \mathbf{h}^{(1)}(\mathbf{x})) \approx g(\mathbf{p}(\mathbf{x}))$ .

**Proposition 2** (Expressivity of the second-stage model). *Suppose  $g$  is a degree  $p$  polynomial with  $\|g\|_{L^2} \lesssim 1$ . Then there exists a learning rate  $\eta$  such that, with high probability over  $\mathcal{D}_1, \mathcal{D}_2, \mathbf{W}$  and  $\mathbf{V}$ , there exists  $\mathbf{a}^* \in \mathbb{R}^{m_1}$  such that the parameter  $\theta^* = (\mathbf{a}^*, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{V})$  achieves a small*

486 empirical loss:

$$488 \quad \frac{1}{n_2} \sum_{\mathbf{x} \in \mathcal{D}_2} (f(\mathbf{x}; \theta^*) - g(\mathbf{p}(\mathbf{x})))^2 = \tilde{\mathcal{O}}\left(\frac{\|\mathbf{a}^*\|_2^2}{m_1^2} + \frac{d^6 r^{p+1}}{m_2} + \frac{d^2 r^{p+1}}{n_1} + \frac{r^{2p+4}}{d^{1/3}}\right).$$

491 Here  $\mathbf{a}^*$  satisfies  $\|\mathbf{a}^*\|_2^2 = \tilde{\mathcal{O}}(m_1 r^p \kappa_2^{2p})$ .

493 The proof is provided in Appendix D.1. We provide following discussions.

494 **Error propagation** To explain the increased polynomial dependence on  $r$ , we remark that the  
495 approximation error in Proposition 1 gets multiplied by the averaged Lipschitz smoothness of the  
496 link function  $g$ , which is upper bounded by  $\mathcal{O}(r^{\frac{p-1}{2}})$ . This product is then squared due to the use of  
497 squared loss in Proposition 2.

498 **Reduced complexity of  $\mathbf{a}$**  Moreover, we remark that the complexity of  $\mathbf{a}$ , i.e.,  $\|\mathbf{a}\|_2$ , gets rid of  
499 the polynomial dependence on  $d$ , which is greatly reduced compared with a naive random-feature  
500 model that requires  $\|\mathbf{a}\|_2^2 = \Theta(m_1 d^{2p})$ . This directly saves the second-stage sample complexity  $n_1$   
501 and the outer width  $m_1$ , since  $n_1, m_2 = \Theta(m_1^{-1} \|\mathbf{a}^*\|_2^2)$  is required for efficient approximation and  
502 generalization (Ghorbani et al., 2021). We also examine this reduced dependency by comparing our  
503 model with a naive random feature model in learning hierarchical target functions in Section A.

505 **Arbitrariness of  $g$**  Thanks to the two-stage architecture and the sufficient learning of the features,  
506 the choice on the link function  $g$  can be an arbitrary degree  $p$  polynomial, not limited to the truth  
507 target  $g^*$ . This allows us to conduct transfer learning tasks in Stage 2 of Algorithm 1.

508 Finally, by standard Rademacher complexity analysis on the random feature model presented in  
509 Appendix E.1, we conclude our proof.

## 5 CONCLUSIONS AND DISCUSSIONS

512 **Comparison with Nichani et al. (2023); Wang et al. (2023)** As discussed under assumptions 3  
513 and the initialization of our neural networks, our work differs significantly in the *targets of interests*,  
514 the *parametrization of neural networks*, the *mathematical strategies*, and the *intuitions behind the results*. Our assumptions ensure a nearly zero linear component and a non-degenerate second  
515 order term of the link function  $g$  which significantly contrasts the assumptions posed in Nichani  
516 et al. (2023); Wang et al. (2023) that emphasize the linear component. Our random initialization  
517 (rather than a deterministic initialization used in the aforementioned two works) in the weights of  
518 the three-layer neural networks allows the learned weights to capture multiple features in all direc-  
519 tions simultaneously after training rather than converge to a single direction. We develop a novel  
520 universality result to relate multiple nonlinear features to multivariate Gaussian, while these two  
521 works adopt existing result of the approximate Stein’s lemma which only applies to single nonlinear  
522 feature. Most importantly, subspace recovery is completely different from and also significantly  
523 harder than single feature recovery considered in Nichani et al. (2023); Wang et al. (2023).

524 **Conclusions** In this work, we have shown the provable capabilities of three-layer networks in ef-  
525 ficiently learning targets of multiple quadratic features. Leveraging a novel universality result, we  
526 have shown that one gradient step suffices for a full recovery of the subspace spanned by multiple  
527 quadratic features. In addition, leveraging the learned features, we have demonstrated the transfer  
528 learning capabilities of this three-layer neural network with a constant polynomial sample com-  
529 plexity guarantee. To the best of our knowledge, this is the first theoretical result of efficiently learning  
530 such a board target function class of multiple nonlinear features with neural networks. We have  
531 made a great improvement on the sample complexity by highlighting feature learning compared to  
532 kernel methods.

533 **Future works** First, it may be possible that the sample complexity bound of  $\tilde{\mathcal{O}}(d^4)$  could be  
534 improved to the information-theoretic optimal sample complexity  $\mathcal{O}(d^2)$  in learning general hier-  
535 archical polynomials of quadratic features. We think that this result may be achieved when we  
536 consider more advanced algorithms that utilize the samples more thoroughly such as using multiple  
537 steps of GD, which could be a great future extension of our work. Moreover, our methodology is  
538 not inherently limited to quadratic features. The principles shown in Figure 1 and techniques devel-  
539 oped here give a foundation for understanding the learning of even more complex function classes.  
Another natural future direction of our work is to understand whether and when our results can be  
generalized to learning multiple high-degree features.

540 REFERENCES  
541

- 542 Emmanuel Abbe, Enric Boix Adserà, and Theodor Misiakiewicz. Sgd learning on neural networks:  
543 leap complexity and saddle-to-saddle dynamics. In *The Thirty Sixth Annual Conference on Learn-*  
544 *ing Theory*, pages 2552–2623. PMLR, 2023.
- 545 Zeyuan Allen-Zhu and Yuanzhi Li. What can resnet learn efficiently, going beyond kernels? *Ad-*  
546 *vances in Neural Information Processing Systems*, 32, 2019.
- 547 Zeyuan Allen-Zhu and Yuanzhi Li. Backward feature correction: How deep learning performs deep  
548 learning. *arXiv preprint arXiv:2001.04413*, 2020.
- 549 Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On  
550 exact computation with an infinitely wide neural net. *Advances in neural information processing*  
551 *systems*, 32, 2019.
- 552 Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg Yang. High-  
553 dimensional asymptotics of feature learning: How one gradient step improves the representation.  
554 *Advances in Neural Information Processing Systems*, 35:37932–37946, 2022.
- 555 Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, and Denny Wu. Learning in the  
556 presence of low-dimensional structure: A spiked random matrix perspective. In A. Oh,  
557 T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neu-*  
558 *ral Information Processing Systems*, volume 36, pages 17420–17449. Curran Associates, Inc.,  
559 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/38a1671ab0747b6ffe4d1c6ef117a3a9-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/38a1671ab0747b6ffe4d1c6ef117a3a9-Paper-Conference.pdf).
- 560 Yu Bai and Jason D. Lee. Beyond linearization: On quadratic and higher-order approximation of  
561 wide neural networks, 2020.
- 562 Alberto Bietti, Joan Bruna, Clayton Sanford, and Min Jae Song. Learning single-index models with  
563 shallow neural networks. *Advances in Neural Information Processing Systems*, 35:9768–9783,  
564 2022.
- 565 Alberto Bietti, Joan Bruna, and Lucas Pillaud-Vivien. On learning gaussian multi-index models  
566 with gradient flow. *arXiv preprint arXiv:2310.19793*, 2023.
- 567 Sourav Chatterjee. Fluctuations of eigenvalues and second order poincaré inequalities, 2007. URL  
568 <https://arxiv.org/abs/0705.1224>.
- 569 Lenaic Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programing.  
570 *Advances in neural information processing systems*, 32, 2019.
- 571 Hugo Cui, Luca Pesce, Yatin Dandi, Florent Krzakala, Yue M. Lu, Lenka Zdeborová, and Bruno  
572 Loureiro. Asymptotics of feature learning in two-layer networks after one gradient-step, 2024.  
573 URL <https://arxiv.org/abs/2402.04980>.
- 574 Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi. Neural networks can learn representations  
575 with gradient descent. In *Conference on Learning Theory*, pages 5413–5452. PMLR, 2022.
- 576 Yatin Dandi, Florent Krzakala, Bruno Loureiro, Luca Pesce, and Ludovic Stephan. Learning two-  
577 layer neural networks, one (giant) step at a time. *arXiv preprint arXiv:2305.18270*, 2023a.
- 578 Yatin Dandi, Florent Krzakala, Bruno Loureiro, Luca Pesce, and Ludovic Stephan. How two-layer  
579 neural networks learn, one (giant) step at a time, 2023b. URL <https://arxiv.org/abs/2305.18270>.
- 580 Yatin Dandi, Luca Pesce, Hugo Cui, Florent Krzakala, Yue M. Lu, and Bruno Loureiro. A ran-  
581 dom matrix theory perspective on the spectrum of learned features and asymptotic generalization  
582 capabilities, 2024a. URL <https://arxiv.org/abs/2410.18938>.
- 583 Yatin Dandi, Emanuele Troiani, Luca Arnaboldi, Luca Pesce, Lenka Zdeborová, and Florent Krza-  
584 kala. The benefits of reusing batches for gradient descent in two-layer networks: Breaking the  
585 curse of information and leap exponents, 2024b. URL <https://arxiv.org/abs/2402.03220>.

- 594 Amit Daniely. Depth separation for neural networks, 2017. URL <https://arxiv.org/abs/1702.08489>.
- 595
- 596
- 597 Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding.
- 598 *arXiv preprint arXiv:1810.04805*, 2018.
- 599
- 600 Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes
- 601 over-parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2018.
- 602
- 603 Weinan E, Chao Ma, and Lei Wu. A comparative analysis of optimization and generalization prop-
- 604 erties of two-layer neural network and random feature models under gradient descent dynamics.
- 605 *Science China Mathematics*, 63(7):1235–1258, jan 2020. doi: 10.1007/s11425-019-1628-5. URL
- 606 <https://doi.org/10.1007%2Fs11425-019-1628-5>.
- 607
- 608 Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Linearized two-layers
- 609 neural networks in high dimension. *The Annals of Statistics*, 49(2):1029 – 1054, 2021.
- 610
- 611 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
- 612 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages
- 613 770–778, 2016.
- 614
- 615 Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and gen-
- 616 eralization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- 617
- 618 Tom H. Koornwinder. *Dual Addition Formulas Associated with Dual Product Formulas*, page
- 619 373–392. WORLD SCIENTIFIC, January 2018. ISBN 9789813228887. doi: 10.1142/
- 620 9789813228887\_0019. URL [http://dx.doi.org/10.1142/9789813228887\\_0019](http://dx.doi.org/10.1142/9789813228887_0019).
- 621
- 622 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet classification with deep convo-
- 623 lutional neural networks. In *Advances in Neural Information Processing Systems*, 2012.
- 624
- 625 Mandar Kulkarni and Shirish Karande. Layer-wise training of deep networks using kernel similarity,
- 626 2017. URL <https://arxiv.org/abs/1703.07115>.
- 627
- 628 Honglak Lee, Chaitanya Ekanadham, and Andrew Ng. Sparse deep belief net model for visual area
- 629 v2. volume Vol 20, 01 2007.
- 630
- 631 Jaehoon Lee, Samuel Schoenholz, Jeffrey Pennington, Ben Adlam, Lechao Xiao, Roman Novak,
- 632 and Jascha Sohl-Dickstein. Finite versus infinite neural networks: an empirical study. *Advances*
- 633 in *Neural Information Processing Systems*, 33:15156–15172, 2020.
- 634
- 635 Jason D. Lee, Kazusato Oko, Taiji Suzuki, and Denny Wu. Neural network learns low-dimensional
- 636 polynomials with sgd near the information-theoretic limit, 2024. URL <https://arxiv.org/abs/2406.01581>.
- 637
- 638 Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Learning with invariances in random
- 639 features and kernel models. In *Conference on Learning Theory*, pages 3351–3418. PMLR, 2021.
- 640
- 641 Behrad Moniri, Donghwan Lee, Hamed Hassani, and Edgar Dobriban. A theory of non-linear feature
- 642 learning with one gradient step in two-layer neural networks, 2024. URL <https://arxiv.org/abs/2310.07891>.
- 643
- 644 Alireza Mousavi-Hosseini, Denny Wu, Taiji Suzuki, and Murat A. Erdogdu. Gradient-based feature
- 645 learning under structured data, 2023. URL <https://arxiv.org/abs/2309.03843>.
- 646
- 647 Alireza Mousavi-Hosseini, Denny Wu, and Murat A. Erdogdu. Learning multi-index models with
- 648 neural networks via mean-field langevin dynamics, 2024. URL <https://arxiv.org/abs/2408.07254>.
- 649
- 650 Eshaan Nichani, Alex Damian, and Jason D Lee. Provable guarantees for nonlinear feature learning
- 651 in three-layer neural networks. *arXiv preprint arXiv:2305.06986*, 2023.
- 652
- 653 Giuseppe Da Prato and Luciano Tubaro. Wick powers in stochastic pdes: an introduction. 2007.
- 654 URL <https://api.semanticscholar.org/CorpusID:55493217>.

- 648 Adityanarayanan Radhakrishnan, Daniel Beaglehole, Parthe Pandit, and Mikhail Belkin. Feature  
 649 learning in neural networks and kernel machines that recursively learn features. *arXiv preprint*  
 650 *arXiv:2212.13881*, 2022.
- 651 Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines.  
 652 In J. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural In-*  
 653 *formation Processing Systems*, volume 20. Curran Associates, Inc., 2007. URL  
 654 [https://proceedings.neurips.cc/paper\\_files/paper/2007/file/013a006f03dbc5392effeb8f18fda755-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2007/file/013a006f03dbc5392effeb8f18fda755-Paper.pdf).
- 655
- 656 Yunwei Ren, Mo Zhou, and Rong Ge. Depth separation with multilayer mean-field networks. *arXiv*  
 657 *preprint arXiv:2304.01063*, 2023.
- 658
- 659 Nathan Ross. Fundamentals of stein’s method. 2011.
- 660
- 661 Itay Safran and Jason Lee. Optimization-based separations for neural networks. In *Conference on*  
 662 *Learning Theory*, pages 3–64. PMLR, 2022.
- 663
- 664 Mahdi Soltanolkotabi, Adel Javanmard, and Jason D Lee. Theoretical insights into the optimization  
 665 landscape of over-parameterized shallow neural networks. *IEEE Transactions on Information*  
 666 *Theory*, 65(2):742–769, 2018.
- 667
- 668 Ramon van Handel. Probability in high dimensions. 2016. URL <https://web.math.princeton.edu/~rvan/APC550.pdf>.
- 669
- 670 Zhichao Wang, Denny Wu, and Zhou Fan. Nonlinear spiked covariance matrices and signal propa-  
 671 gation in deep neural networks, 2024. URL <https://arxiv.org/abs/2402.10127>.
- 672
- 673 Zihao Wang, Eshaan Nichani, and Jason D. Lee. Learning hierarchical polynomials with three-layer  
 674 neural networks, 2023.
- 675
- 676 Colin Wei, Jason D. Lee, Qiang Liu, and Tengyu Ma. Regularization matters: Generalization and  
 677 optimization of neural nets v.s. their induced kernel, 2020.
- 678
- 679 Greg Yang and Edward J Hu. Tensor programs iv: Feature learning in infinite-width neural networks.  
 680 In *International Conference on Machine Learning*, pages 11727–11737. PMLR, 2021.
- 681
- 682
- 683
- 684
- 685
- 686
- 687
- 688
- 689
- 690
- 691
- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701

# Appendix

---

756    **A NUMERICAL EXPERIMENTS**

758    We empirically verify Theorem 1 and Proposition 1. We consider learning functions with  $r = 3$   
 759    quadratic features. Regarding the target function, we choose the target functions to be of the form

$$760 \quad f_{d,p}^*(\mathbf{x}) = \frac{f_{d,p}(\mathbf{x}) - \mathbb{E}[f_{d,p}(\mathbf{x})]}{\sqrt{\text{Var}[f_{d,p}(\mathbf{x})]}}, \quad \text{with } f_{d,p}(\mathbf{x}) = \sum_{i=1}^r (\mathbf{x}^\top \mathbf{A}_i \mathbf{x})^p, \quad p \in \mathbb{N}. \quad (5)$$

764    For the underlying features, we take  $\mathbf{p}(\mathbf{x}) = [\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \mathbf{x}^\top \mathbf{A}_2 \mathbf{x}, \mathbf{x}^\top \mathbf{A}_3 \mathbf{x}]^\top$  with  $\mathbf{A}_k = \text{diag}(c \cdot \mathbf{a}_k)$ ,  
 765    and  $c > 0$  is a normalizing constant. To ensure the orthogonality of the features and  $\text{tr}(\mathbf{A}_k) = 0$ ,  
 766    we choose the ambient dimension  $d$  to be divisible by 4 and take  $\mathbf{a}_k$  to be

$$767 \quad \mathbf{a}_1 = \text{Vec}([\mathbf{1}, \mathbf{1}, -\mathbf{1}, -\mathbf{1}]), \quad \mathbf{a}_2 = \text{Vec}([\mathbf{1}, -\mathbf{1}, \mathbf{1}, -\mathbf{1}]), \quad \text{and } \mathbf{a}_3 = \text{Vec}([\mathbf{1}, -\mathbf{1}, -\mathbf{1}, \mathbf{1}]).$$

769    Here  $\mathbf{1}$  is a vector of ones in  $d/4$  dimensions, and  $c = \sqrt{\frac{d+2}{2d^2}}$  to ensure that  $\mathbb{E}_{\mathbf{x}}[(\mathbf{x}^\top \mathbf{A}_k \mathbf{x})^2] = 1$   
 770    for each  $k = 1, 2, 3$ .

772    For the network architecture, we choose  $\sigma_1$  as per (2) and  $\sigma_2 = Q_2$ , with network sizes set to  
 773     $m_1 = 10000$  and  $m_2 = 20000$ . We compare our proposed model (4) (given by Algorithm 1) against  
 774    the naive random-feature model defined as

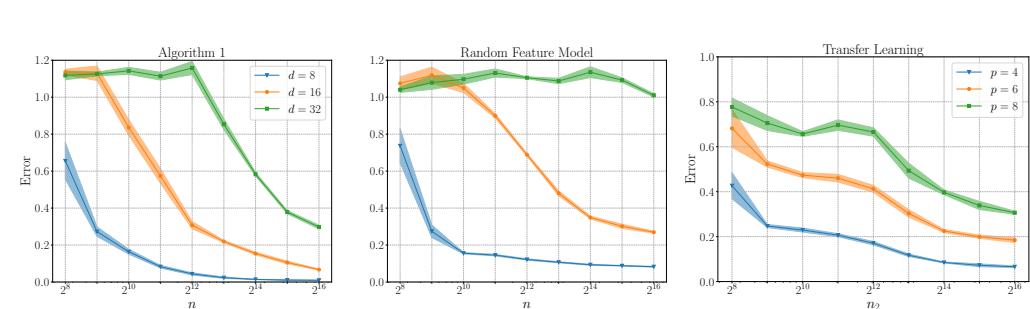
$$776 \quad f^{\text{RF}}(\mathbf{x}'; \theta) = \frac{1}{m_1} \sum_{j=1}^{m_1} a_j \sigma_1\left(\eta a_j^{(0)} \langle \mathbf{w}_j, \mathbf{h}^{(0)}(\mathbf{x}') \rangle + b_j^{(1)}\right), \quad (6)$$

779    where  $\mathbf{a}$  is the only trainable parameter throughout the training process. Our experiments involve  
 780    learning  $f_{d,p}^*$  with  $p = 4$  and  $d \in \{8, 16, 32\}$ . To examine our model’s transfer learning capabilities,  
 781    we also train the model on an initial target function  $f_{d,2}^*$  with  $d = 16$  and  $n_1 = 2^{16}$  in the first stage,  
 782    then transfer to targets  $f_{d,p}^*$  with  $p = 4, 6, 8$ . For each task, we explore a range of sample sizes from  
 783     $2^8$  to  $2^{16}$ . The results of these experiments are presented in Figure 2.

785    **Improved sample complexity and Polynomial dependence on  $d$**  The left panel of Figure 2  
 786    demonstrates that our model outperforms the naive random-feature model across all dimensions.  
 787    As the dimension  $d$  increases, both models show larger test errors, but our model exhibits less sensi-  
 788    tivity to  $d$ . This aligns with our theoretical analysis in Theorem 1 that the sample complexity of  
 789    kernel methods should be  $\Omega(d^{2p-4})$  times greater than that of our model. Moreover, we redraw  
 790    Figure 2 by plotting the test error against  $\log_d n$ . As shown in Figure 3, the loss curves for our  
 791    model (Algorithm 1) align closely for different values of  $d$ , indicating that it achieves low error rates  
 792    with only  $\tilde{\mathcal{O}}(d^4)$  samples. In stark contrast, the naive random feature model exhibits significant  
 793    separation between curves for different  $d$  values, requiring more than  $\tilde{\mathcal{O}}(d^4)$  samples to achieve  
 794    comparable error rates. This graphical evidence powerfully demonstrates how our approach elimi-  
 795    nates the dependence on dimension  $\Theta(d^{2p})$  presented in kernel methods, resulting in substantially  
 796    improved sample complexity in high-dimensional settings.

797    **Efficient transfer learning** The right panel of Figure 2 showcases our algorithm’s strong transfer  
 798    learning capabilities. our algorithm successfully learns all three transferred target functions with  
 799    benign second-stage sample complexity. Notably, as the degree  $p$  increases, the test error grows  
 800    no faster than  $r^p$ , which is significantly slower than  $d^{2p}$ . This supports our theoretical result that  
 801    the second-stage sample complexity depends on the number of features  $r$  rather than the ambient  
 802    dimension  $d$ , underscoring our model’s strong transfer learning capabilities.

803    **Accurate reconstruction of quadratic features** To further demonstrate our model’s feature learn-  
 804    ing capabilities, we extract the learned features  $\mathbf{h}^{(1)}$  after the first training stage of Algorithm 1, us-  
 805    ing  $f_{16,2}$  as the target. We then reconstruct these features using a linear transformation  $\mathbf{B}^* \in \mathbb{R}^{r \times m_2}$ ,  
 806    as described in Proposition 1. We examine how reconstruction accuracy changes with first-stage  
 807    sample sizes. Figure 4 shows the correlation between true and reconstructed features for each sam-  
 808    ple size. As  $n_1$  increases, all features are better approximated simultaneously. Notably,  $d^4$  samples  
 809    prove sufficient to reconstruct the features with high accuracy, supporting our model’s effective fea-  
 810    ture learning ability.



(a) Comparison between Algorithm 1 and the random-feature model (b) Performance of transfer learning

Figure 2: For the left panel, Algorithm 1 uses two equally sized datasets, while the random feature model uses the full dataset. For the right panel, we conduct transfer learning with  $n_1 = 2^{16}$  pretraining samples and plot the dependence on  $n_2$ . The figure reports the mean and normalized standard error of the test error using 10,000 fresh samples, based on 5 independent experimental instances.

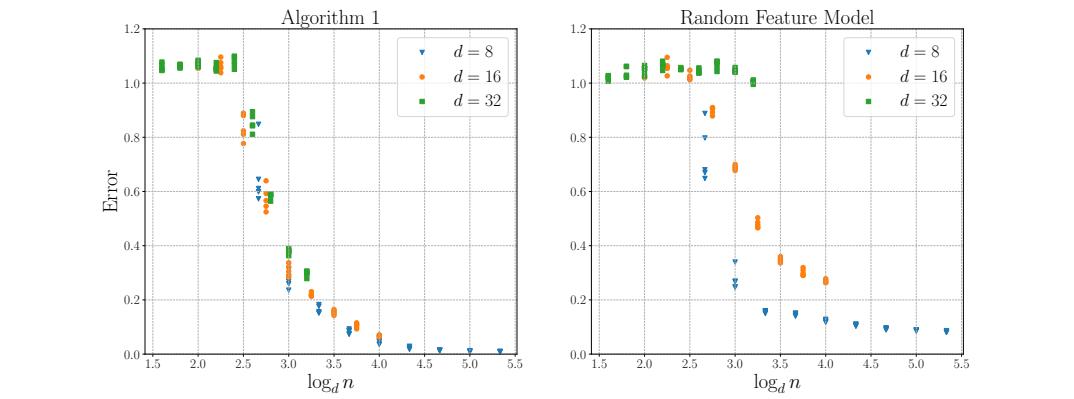


Figure 3: Test error of Algorithm 1 and the naive random feature models with x-axis being the relative sample complexity ( $\log_d n$ ). We plot the test error of 5 independent instances for each  $d \in \{8, 16, 32\}$ .

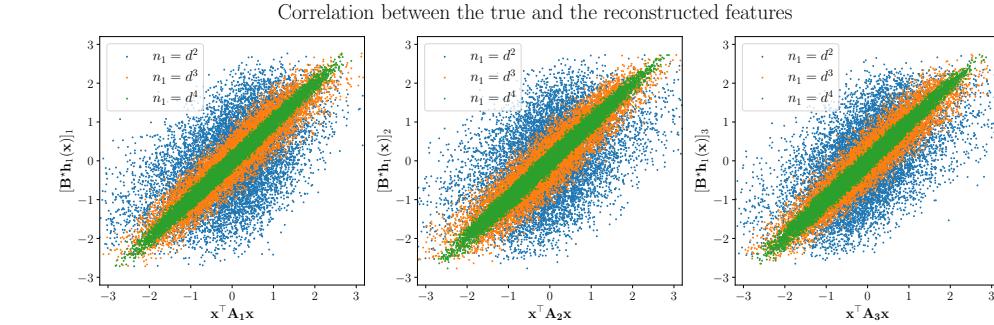


Figure 4: The linear correlation between the three true features and their corresponding reconstructed features for varying first-stage sample sizes  $n_1$ . The reconstructed features are standardized to match the variance of the true features. For  $i = 1, 2, 3$ , the  $i$ -th scatter plot represents 10,000 test sample points of  $([B^* h^{(1)}(x)]_i, x^\top A_i x)$  for  $n_1 \in \{d^2, d^3, d^4\}$ , where  $d = 16$ .

864    **B TECHNICAL BACKGROUND**  
 865

866    **B.1 ASYMPTOTIC NOTATION**

867    Throughout the proof we will let  $C$  be a fixed but sufficiently large constant.

868    **Definition 1** (high probability events). Let  $\iota = C \log(dn_1n_2m_1m_2)$ . We say that an event happens  
 869    with high probability if it happens with probability at least  $1 - \text{poly}(d, n_1, n_2, m_1, m_2)e^{-\iota}$ .

870  
 871    **Example 1.** If  $z \sim N(0, 1)$  then  $|z| \leq \sqrt{2\iota}$  with high probability.

872    Note that high probability events are closed under union bounds over sets of size  
 873     $\text{poly}(d, n_1, n_2, m_1, m_2)$ , such as  $\mathcal{D}_1$ ,  $\mathcal{D}_2$  and  $\{\mathbf{w}_j\}_{j \in [m_1]}$ . We will also assume throughout the  
 874    paper that  $\iota \leq C^{-1}d$ .

875    **B.2 MULTIVARIATE GAUSSIAN APPROXIMATION**  
 876

877    In this section, we assume that  $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  and aim to establish an upper bound of Wasserstein  
 878    distance between the distribution of  $\mathbf{p}(\mathbf{X})$  and the standard  $r$ -dimensional Gaussian distribution,  
 879    i.e., Lemma 1.

880    To prove Lemma 1, we introduce Stein’s method (Ross, 2011) for multivariate Gaussian approxima-  
 881    tion. We will use the following additional notations.

- 882  
 883    •  $\mathcal{G}f(\mathbf{x}) := \int_0^\infty \mathbb{E}_{\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [f(e^{-t}\mathbf{x} + \sqrt{1-e^{-2t}}\mathbf{Z}) - f(\mathbf{Z})] dt$  denotes the potential op-  
 884    erator of  $f$ .  
 885    •  $\mathcal{J}(\mathbf{p}) := [\nabla p_1, \nabla p_2, \dots, \nabla p_r]^\top \in \mathbb{R}^{r \times n}$  denotes the Jacobian matrix of  $\mathbf{p}$ .

886  
 887    Now we state the supporting lemmas to prove Lemma 1.

888    **Lemma 2** (Corollary 9.12 in van Handel (2016)). *For any probability measure  $\mu$  in  $\mathbb{R}^r$ , we have*

889  
 890    
$$W_1(\mu, \mathcal{N}(\mathbf{0}, \mathbf{I}_r)) \leq \sup_{\|\nabla g\| \leq 1, \|\nabla^2 g\| \leq \sqrt{\frac{2}{\pi}}} \mathbb{E}_{\mathbf{Y} \sim \mu} [\Delta g(\mathbf{Y}) - \langle \nabla g(\mathbf{Y}), \mathbf{Y} \rangle].$$

891  
 892    **Lemma 3** (Lemma 9.21 in van Handel (2016)). *Suppose  $\mathbf{X} = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$  is an  
 893     $d$ -dimensional standard Gaussian variable. Then for any functions  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  and  $h : \mathbb{R}^d \rightarrow \mathbb{R}$ ,  
 894    we have*

895  
 896    
$$\text{Cov}(g(\mathbf{X}), h(\mathbf{X})) = \mathbb{E}_{\mathbf{X}} [\langle \nabla g(\mathbf{X}), \nabla \mathcal{G}h(\mathbf{X}) \rangle]$$

897    With the lemmas above, we begin our proof of Lemma 1.

898  
 899    *Proof of Lemma 1.* By invoking Lemma 2 with  $\mu = \text{Law}(\mathbf{p})$  and  $\mathbf{Y} = \mathbf{p}(\mathbf{X})$ , for any  $g(\mathbf{y}) : \mathbb{R}^r \rightarrow \mathbb{R}$  with  $\|\nabla g\| \leq 1$  and  $\|\nabla^2 g\| \leq \sqrt{\frac{2}{\pi}}$ , we aim to bound

900  
 901    
$$\underbrace{\mathbb{E}_{\mathbf{X}} [\Delta g(\mathbf{p}(\mathbf{X})) - \langle \nabla g(\mathbf{p}(\mathbf{X})), \mathbf{p}(\mathbf{X}) \rangle]}_{\blacklozenge} = \sum_{i=1}^r \mathbb{E}_{\mathbf{X}} \left[ \frac{\partial^2 g}{\partial y_i^2} \Big|_{\mathbf{y}=\mathbf{p}(\mathbf{X})} - p_i(\mathbf{X}) \frac{\partial g}{\partial y_i} \Big|_{\mathbf{y}=\mathbf{p}(\mathbf{X})} \right].$$

902    Since for any  $i \in [r]$ ,  $\mathbb{E}[p_i(\mathbf{X})] = 0$ , we have

903  
 904    
$$\begin{aligned} \mathbb{E} \left[ p_i(\mathbf{X}) \frac{\partial g}{\partial y_i} \Big|_{\mathbf{y}=\mathbf{p}(\mathbf{X})} \right] &= \text{Cov} \left( p_i(\mathbf{X}), \frac{\partial g}{\partial y_i} \Big|_{\mathbf{y}=\mathbf{p}(\mathbf{X})} \right) \\ 905 &= \mathbb{E} \left[ \left\langle \nabla_{\mathbf{x}} \frac{\partial g}{\partial y_i} \Big|_{\mathbf{y}=\mathbf{p}(\mathbf{X})}, \nabla_{\mathbf{x}} \mathcal{G}p_i(\mathbf{X}) \right\rangle \right] \\ 906 &= \mathbb{E} \left[ \left\langle \sum_{j=1}^r \frac{\partial^2 g}{\partial y_i \partial y_j} \Big|_{\mathbf{y}=\mathbf{p}(\mathbf{X})}, \nabla_{\mathbf{x}} p_j(\mathbf{X}), \nabla_{\mathbf{x}} \mathcal{G}p_i(\mathbf{X}) \right\rangle \right] \\ 907 &= \sum_{j=1}^r \mathbb{E} \left[ \frac{\partial^2 g}{\partial y_i \partial y_j} \Big|_{\mathbf{y}=\mathbf{p}(\mathbf{X})} \langle \nabla_{\mathbf{x}} p_j(\mathbf{X}), \nabla_{\mathbf{x}} \mathcal{G}p_i(\mathbf{X}) \rangle \right], \end{aligned}$$

918 where the second equality follows from Lemma 3 and we obtain the third equality by the chain rule.  
919 Thus, we have

$$920 \quad \spadesuit = \mathbb{E} [\langle \nabla^2 g(\mathbf{p}(\mathbf{X})), \mathbf{I}_r - \mathcal{J}(\mathbf{p}(\mathbf{X}))\mathcal{J}(\mathcal{G}\mathbf{p}(\mathbf{X}))^\top \rangle]. \quad (7)$$

922 For a special case, for any  $i, j \in [r]$ , we take  $g(\mathbf{y}) = y_i y_j$  in (7), obtaining that

$$924 \quad \mathbb{E} [\langle \nabla_{\mathbf{x}} p_j(\mathbf{X}), \nabla_{\mathbf{x}} \mathcal{G}p_i(\mathbf{X}) \rangle] = \begin{cases} \mathbb{E} [2p_j(\mathbf{X})p_i(\mathbf{X})] = 0, & i \neq j, \\ \mathbb{E} [2p_i^2(\mathbf{X})] - 1 = 1, & i = j. \end{cases}$$

926 Thus,  $\mathbb{E} [\mathbf{I}_r - \mathcal{J}(\mathbf{p}(\mathbf{X}))\mathcal{J}(\mathcal{G}\mathbf{p}(\mathbf{X}))^\top] = \mathbf{0}_{r \times r}$ . Since  $\|\nabla^2 g\| \leq \sqrt{\frac{2}{\pi}}$ , we have  $|\langle \nabla^2 g \rangle_{i,j}| \leq \sqrt{\frac{2}{\pi}}$   
927 for any  $i, j \in [r]$ . We can therefore estimate

$$\begin{aligned} 929 \quad W_1(\text{Law}(\mathbf{p}(\mathbf{X})), \mathcal{N}(\mathbf{0}, \mathbf{I}_r)) &\leq \sqrt{\frac{2}{\pi}} \sum_{i,j \in [r]} \mathbb{E} [|\delta_{i,j} - \langle \nabla_{\mathbf{x}} p_j(\mathbf{X}), \nabla_{\mathbf{x}} \mathcal{G}p_i(\mathbf{X}) \rangle|] \\ 930 \quad &\leq \sqrt{\frac{2}{\pi}} \sum_{i,j \in [r]} \text{Var}[\langle \nabla_{\mathbf{x}} p_j(\mathbf{X}), \nabla_{\mathbf{x}} \mathcal{G}p_i(\mathbf{X}) \rangle]^{1/2} \\ 931 \quad &\leq \sqrt{\frac{2}{\pi}} \sum_{i,j \in [r]} \mathbb{E} [\|\nabla_{\mathbf{x}} \langle \nabla_{\mathbf{x}} p_j(\mathbf{X}), \nabla_{\mathbf{x}} \mathcal{G}p_i(\mathbf{X}) \rangle\|^2]^{1/2}, \end{aligned}$$

938 where we invoke Poincaré inequality in the last inequality. For any  $i, j \in [r]$ , we have

$$\begin{aligned} 940 \quad &\mathbb{E} [\|\nabla_{\mathbf{x}} \langle \nabla_{\mathbf{x}} p_j(\mathbf{X}), \nabla_{\mathbf{x}} \mathcal{G}p_i(\mathbf{X}) \rangle\|^2] \\ 941 \quad &= \mathbb{E} [\|\nabla_{\mathbf{x}}^2 p_j(\mathbf{X}) \nabla \mathcal{G}p_i(\mathbf{X}) + \nabla_{\mathbf{x}} p_j(\mathbf{X}) \nabla^2 \mathcal{G}p_i(\mathbf{X})\|^2] \\ 942 \quad &\leq 2 \mathbb{E} [\|\nabla_{\mathbf{x}}^2 p_j(\mathbf{X}) \nabla \mathcal{G}p_i(\mathbf{X})\|^2] + 2 \mathbb{E} [\|\nabla_{\mathbf{x}} p_j(\mathbf{X}) \nabla^2 \mathcal{G}p_i(\mathbf{X})\|^2] \\ 943 \quad &\leq 2 \mathbb{E} [\|\nabla^2 p_j\|^4]^{1/2} \mathbb{E} [\|\nabla \mathcal{G}p_i\|^4]^{1/2} + 2 \mathbb{E} [\|\nabla p_j\|^4]^{1/2} \mathbb{E} [\|\nabla^2 \mathcal{G}p_i\|^4]^{1/2} \\ 944 \quad &\leq 2 \mathbb{E} [\|\nabla^2 p_j\|^4]^{1/2} \mathbb{E} [\|\nabla p_i\|^4]^{1/2} + 2 \mathbb{E} [\|\nabla p_j\|^4]^{1/2} \mathbb{E} [\|\nabla^2 p_i\|^4]^{1/2}. \end{aligned}$$

949 The last inequality follows from the inequality in Page 308 in van Handel (2016). By adding up all  
950 the terms along  $i$  and  $j$ , we have

$$\begin{aligned} 951 \quad &W_1(\text{Law}(\mathbf{p}(\mathbf{X})), \mathcal{N}(\mathbf{0}, \mathbf{I}_r)) \\ 952 \quad &\leq \sqrt{\frac{2}{\pi}} \sum_{i,j \in [r]} \sqrt{2 \mathbb{E} [\|\nabla^2 p_j\|^4]^{1/2} \mathbb{E} [\|\nabla p_i\|^4]^{1/2} + 2 \mathbb{E} [\|\nabla p_j\|^4]^{1/2} \mathbb{E} [\|\nabla^2 p_i\|^4]^{1/2}} \\ 953 \quad &\leq \frac{2}{\sqrt{\pi}} \sum_{i,j \in [r]} \left( \mathbb{E} [\|\nabla^2 p_j\|^4]^{1/4} \mathbb{E} [\|\nabla p_i\|^4]^{1/4} + \mathbb{E} [\|\nabla p_j\|^4]^{1/4} \mathbb{E} [\|\nabla^2 p_i\|^4]^{1/4} \right) \\ 954 \quad &= \frac{4}{\sqrt{\pi}} \left( \sum_{i=1}^r \mathbb{E} [\|\nabla p_i\|^4]^{1/4} \right) \left( \sum_{j=1}^r \mathbb{E} [\|\nabla^2 p_j\|^4]^{1/4} \right). \end{aligned}$$

962 We complete our proof. □

### 964 B.3 HYPERCONTRACTIVITY OF POLYNOMIALS

965 The following Lemma is cited from Mei et al. (2021) and is designed for uniform distribution on the  
966 sphere in  $d$  dimension.

967 **Lemma 4.** For any  $\ell \in \mathbb{N}$  and  $f \in L^2(\mathbb{S}^{d-1})$  to be a degree  $\ell$  polynomial, for any  $q \geq 2$ , we have

$$968 \quad \left( \mathbb{E}_{\mathbf{z} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))} [f(\mathbf{z})^q] \right)^{2/q} \leq (q-1)^\ell \mathbb{E}_{\mathbf{z} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))} [f(\mathbf{z})^2].$$

971 We remark that the results above are also multiplicative.

972    **Lemma 5.** For any  $\ell \in \mathbb{N}$  and  $f \in L^2((\mathbb{S}^{d-1})^k)$  to be a degree  $\ell$  polynomial in the components of  
 973    each  $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_k$ , for any  $q \geq 2$ , we have  
 974

$$975 \quad \left( \mathbb{E}_{\mathbf{z} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))^k} [f(\mathbf{z})^q] \right)^{2/q} \leq (q-1)^{k\ell} \mathbb{E}_{\mathbf{z} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))^k} [f(\mathbf{z})^2].$$

977    Here  $\mathbf{z} = \text{Vec}([\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_k])$ .  
 978

979    For the case where the input distribution is standard Gaussian in  $d$  dimension (denoted as  $\gamma$ ), we  
 980    have the next Lemma from Theorem 4.3, Prato and Tubaro (2007).  
 981

982    **Lemma 6.** For any  $\ell \in \mathbb{N}$  and  $f \in L^2(\gamma)$  to be a degree  $\ell$  polynomial, for any  $q \geq 2$ , we have  
 983

$$983 \quad \mathbb{E}_{\mathbf{z} \sim \gamma} [f(\mathbf{z})^q] \leq \mathcal{O}_{q,\ell}(1) (\mathbb{E}_{\mathbf{z} \sim \gamma} [f(\mathbf{z})^2])^{q/2}.$$

985    where we use  $\mathcal{O}_{q,\ell}(1)$  to denote some universal constant that only depends on  $q, \ell$ .  
 986

987    Moreover, we introduce lemmas to control the deviation of random variables which polynomially  
 988    depend on some Gaussian random variables. We will use a slightly modified version of Lemma 30  
 989    from Damian et al. (2022).

990    **Lemma 7.** Let  $g$  be a polynomial of degree  $p$  and  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . Then there exists an absolute  
 991    positive constant  $C_p$  depending only on  $p$  such that for any  $\delta > 1$ ,

$$992 \quad \mathbb{P} [|g(\mathbf{x}) - \mathbb{E}[g(\mathbf{x})]| \geq \delta \sqrt{\text{Var}(g(\mathbf{x}))}] \leq 2 \exp(-C_p \delta^{2/p}).$$

994    We also have the spherical version of Lemma 7.  
 995

996    **Lemma 8.** Let  $g$  be a polynomial of degree  $p$  and  $\mathbf{x} \sim \mathbb{S}^{d-1}(\sqrt{d})$ . Then there exists an absolute  
 997    positive constant  $C_p$  depending only on  $p$  such that for any  $\delta > 1$ ,

$$998 \quad \mathbb{P} [|g(\mathbf{x}) - \mathbb{E}[g(\mathbf{x})]| \geq \delta \sqrt{\text{Var}(g(\mathbf{x}))}] \leq 2 \exp(-C_p \delta^{2/p}).$$

1001    Thus, for a degree- $p$  polynomial  $g$ , we have  $g(\mathbf{x}) \lesssim \delta^{p/2} \|g\|_{L^2}$  with high probability.  
 1002

#### 1003    B.4 MOMENTS AND FACTORIZATION OF POLYNOMIALS

1004    In this section, we present formulae for calculating moments of Gaussian or spherical variables,  
 1005    cited from Damian et al. (2022).

1006    **Lemma 9** (Expectations of Gaussian tensors). For  $\mathbf{w} \in \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$  and  $k \in \mathbb{N}$ , we have  
 1007

$$1008 \quad \mathbb{E}_{\mathbf{w}} [\mathbf{w}^{\otimes 2k}] = (2k-1)!! \text{Sym}(\mathbf{I}_d^{\otimes k})$$

1009    Here  $\text{Sym}(\mathbf{T})$  is the symmetrization of a  $k$ -tensor  $\mathbf{T} \in (\mathbb{R}^d)^{\otimes k}$  across all  $k$  axes.  
 1010

1011    Leveraging this calculation, we can factorize any polynomial  $g$  into inner products between high-  
 1012    order tensors and bound the Frobenius norm of the tensors.

1013    **Lemma 10.** (Lemma 21 in Damian et al. (2022)) Given Let  $g : \mathbb{R}^r \rightarrow \mathbb{R}$  be an degree- $p$  polynomial.  
 1014    Then there exists  $\mathbf{T}_0, \mathbf{T}_1, \dots, \mathbf{T}_p$  such that  
 1015

$$1016 \quad g(\mathbf{z}) = \sum_{k=0}^p \langle \mathbf{T}_k, \mathbf{z}^{\otimes k} \rangle \quad \text{with} \quad \|\mathbf{T}_k\|_{\text{F}} \lesssim \|g\|_{L^2} r^{\frac{p-k}{4}}, \quad k = 0, 1, \dots, p.$$

1019    Here  $\|g\|_{L^2} = \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [g^2(\mathbf{z})]$ .  
 1020

1021    As a corollary, we then have  $\nabla g(\mathbf{z}) = \sum_{k=1}^p k \mathbf{T}_k (\mathbf{z}^{\otimes k-1})$  and  
 1022

$$1023 \quad \|\nabla g(\mathbf{z})\| \leq \sum_{k=1}^p k \|\mathbf{T}_k\|_{\text{F}} \|\mathbf{z}\|^{k-1} \lesssim \|g\|_{L^2} \sum_{k=1}^p k r^{\frac{p-k}{4}} \|\mathbf{z}\|^{k-1}. \quad (8)$$

1025    For a spherical variable  $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$  we can also compute its moments.

1026 **Lemma 11** (Expectations of Spherical tensors). *For  $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$  and  $k \in \mathbb{N}$ , we have*

$$1028 \quad \mathbb{E}_{\mathbf{z}} [\mathbf{z}^{\otimes 2k}] = d^k \cdot \frac{\mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)} [\mathbf{w}^{\otimes 2k}]}{\mathbb{E}_{v \sim \chi(d)} [v^{2k}]},$$

1030 where  $\chi(d)$  represents the chi-distribution with the degree of freedom being  $d$ , and its moments can  
1031 be computed as  
1032

$$1033 \quad \mathbb{E}_{v \sim \chi(d)} [v^{2k}] = \prod_{j=0}^{k-1} (d + 2j) = \Theta(d^k).$$

1036 As an example, the moments of spherical quadratic forms  $\mathbf{x}^\top \mathbf{A} \mathbf{x}$  can be computed explicitly as

$$1039 \quad \mathbb{E}_{\mathbf{x}} [\mathbf{x}^\top \mathbf{A} \mathbf{x}] = \text{tr}(\mathbf{A}), \text{ and } \mathbb{E}_{\mathbf{x}} [(\mathbf{x}^\top \mathbf{A} \mathbf{x})(\mathbf{x}^\top \mathbf{B} \mathbf{x})] = \frac{d}{d+2} \cdot (\text{tr}(\mathbf{A})\text{tr}(\mathbf{B}) + 2\langle \mathbf{A}, \mathbf{B} \rangle).$$

1041 Thus, to satisfy Assumption 1, we require  $\text{tr}(\mathbf{A}_k) = 0$ ,  $\|\mathbf{A}_k\|_F = \sqrt{(d+2)/(2d)}$  and  $\langle \mathbf{A}_k, \mathbf{A}_\ell \rangle = 0$  for any  $k, \ell \in [r]$ .  
1042

## 1043 B.5 SPHERICAL HARMONICS AND GEGENBAUER POLYNOMIALS

1044 We introduce some facts of spherical harmonics and Gegenbauer polynomials, with the first four  
1045 properties from [Ghorbani et al. \(2021\)](#) and the last one from [Koornwinder \(2018\)](#).  
1046

1047 1. For  $\mathbf{x}, \mathbf{y} \in \mathbb{S}^{d-1}(\sqrt{d})$ ,

$$1049 \quad |Q_j(\langle \mathbf{x}, \mathbf{y} \rangle)| \leq Q_j(d) = 1. \quad (9)$$

1050 2. For  $\mathbf{x}, \mathbf{y} \in \mathbb{S}^{d-1}(\sqrt{d})$ ,

$$1052 \quad \langle Q_j(\langle \mathbf{x}, \cdot \rangle), Q_k(\langle \mathbf{y}, \cdot \rangle) \rangle_{L^2} = \frac{1}{B(d, k)} \delta_{jk} Q_k(\langle \mathbf{x}, \mathbf{y} \rangle). \quad (10)$$

1054 Here  $B(d, k)$  denotes the dimension of subspace of degree  $k$  spherical harmonics  
1055

$$1056 \quad B(d, k) := \dim(V_{d,k}) = \frac{2k+d-2}{k} \binom{k+d-3}{k-1} = \Theta(d^k).$$

1059 3. For  $\mathbf{x}, \mathbf{y} \in \mathbb{S}^{d-1}(\sqrt{d})$ ,

$$1061 \quad Q_k(\langle \mathbf{x}, \mathbf{y} \rangle) = \frac{1}{B(d, k)} \sum_{i=1}^{B(d, k)} Y_{k,i}(\mathbf{x}) Y_{k,i}(\mathbf{y}). \quad (11)$$

1064 4. For any  $k \in \mathbb{N}_{\geq 1}$ ,

$$1066 \quad \frac{t}{d} Q_k(t) = \frac{k}{2k+d-2} Q_{k-1}(t) + \frac{k+d-2}{2k+d-2} Q_{k+1}(t). \quad (12)$$

1068 5. For any  $i, j \in \mathbb{N}$ ,

$$1070 \quad Q_i(t) Q_j(t) = \sum_{k=0}^{\min(i,j)} b_{i+j-2k}^{(i,j)} \binom{i}{k} \binom{j}{k} k! Q_{i+j-2k}(t). \quad (13)$$

1073 Here, we have

$$1075 \quad b_{i+j-2k}^{(i,j)} = \frac{2(i+j-2k)+d-2}{d-2} \cdot \frac{((d-2)/2)_k ((d-2)/2)_{i-k} ((d-2)/2)_{j-k} (d-2)_{i+j-k}}{(d-2)_i (d-2)_j (d/2)_{i+j-k}}.$$

1077 We note that  $(z)_k = z(z+1)\cdots(z+k-1) = \Gamma(z+k)/\Gamma(z)$  is the Pochhammer symbol. Given  
1078 any  $i$  and  $j$ , we have  $d^k b_{i+j-2k}^{(i,j)} \rightarrow 1$  when  $d \rightarrow \infty$ . We derive a quantitative bound on the scale of  
1079  $b_{i+j-2k}^{(i,j)}$  in Lemma 12.

1080 **Lemma 12.** For any  $i, j \geq k \geq 0$ , denote

$$1081 \quad c_{i+j-2k}^{(i,j)} = \frac{((d-2)/2)_k ((d-2)/2)_{i-k} ((d-2)/2)_{j-k} (d-2)_{i+j-k}}{(d-2)_i (d-2)_j (d/2)_{i+j-k}}.$$

1082 Then, when  $d \geq 4$ , it holds that  $c_{i+j-2k}^{(i,j)} \leq \frac{1}{(d-2)_k}$ .

1083 *Proof of Lemma 12.* Note that when  $d \geq 4$ , i.e.,  $d-2 \geq d/2$ ,

$$\begin{aligned} 1084 \quad \frac{c_{(i+1)+j-2k}^{(i+1,j)}}{c_{i+j-2k}^{(i,j)}} &= \frac{\left(\frac{d-2}{2} + i - k\right)(d-2+i+j-k)}{(d-2+i)\left(\frac{d}{2} + i + j - k\right)} \quad (\text{monotone decreasing with } j) \\ 1085 \quad &\leq \frac{\left(\frac{d-2}{2} + i - k\right)(d-2+i)}{(d-2+i)\left(\frac{d}{2} + i\right)} \\ 1086 \quad &< 1. \end{aligned}$$

1087 Thus, we have  $c_{(i+1)+j-2k}^{(i+1,j)} \leq c_{i+j-2k}^{(i,j)}$ . Similarly, we have  $c_{i+(j+1)-2k}^{(i,j+1)} \leq c_{i+j-2k}^{(i,j)}$ . Consequently, for any  $i, j \geq k$ , we have

$$\begin{aligned} 1088 \quad c_{i+j-2k}^{(i,j)} &\leq c_{k+k-2k}^{(k,k)} \\ 1089 \quad &= \frac{((d-2)/2)_k ((d-2)/2)_0 ((d-2)/2)_0 (d-2)_k}{(d-2)_k (d-2)_k (d/2)_k} \\ 1090 \quad &= \frac{((d-2)/2)_k}{(d-2)_k (d/2)_k} \\ 1091 \quad &= \frac{(d-2)/2}{(d-2)_k (d/2 + k - 1)} \\ 1092 \quad &\leq \frac{1}{(d-2)_k}. \end{aligned}$$

1093 The proof is complete.  $\square$

## C APPROXIMATION THEORY OF THE INNER LAYER

1110 Since we focus on the first training stage throughout this section, we denote  $n = n_1$  for notation simplicity when the context is clear, and let the training set be  $\mathcal{D}_1 = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ .

### C.1 ASYMPTOTIC ANALYSIS OF THE LEARNED FEATURE

1111 In this subsection, we analyse the learned feature  $\mathbf{h}^{(1)}(\mathbf{x}')$  in the asymptotic way, i.e.,  $m_2, n \rightarrow \infty$ . Note that we can rewrite the learned feature as

$$\begin{aligned} 1112 \quad \mathbf{h}^{(1)}(\mathbf{x}') &= \frac{1}{nm_2} \sum_{i=1}^n f^*(\mathbf{x}_i) \langle \mathbf{h}^{(0)}(\mathbf{x}_i), \mathbf{h}^{(0)}(\mathbf{x}') \rangle \mathbf{h}^{(0)}(\mathbf{x}_i) \\ 1113 \quad &= \frac{1}{n} \sum_{i=1}^n f^*(\mathbf{x}_i) K_{m_2}^{(0)}(\mathbf{x}, \mathbf{x}') \mathbf{h}^{(0)}(\mathbf{x}_i). \end{aligned}$$

1114 where the initial kernel  $K_{m_2}^{(0)}(\mathbf{x}, \mathbf{x}')$  is defined as

$$1115 \quad K_{m_2}^{(0)}(\mathbf{x}, \mathbf{x}') = \frac{1}{m_2} \langle \sigma_2(\mathbf{V}\mathbf{x}), \sigma_2(\mathbf{V}\mathbf{x}') \rangle \approx \mathbb{E}_{\mathbf{v}} [\sigma_2(\mathbf{v}^\top \mathbf{x}) \sigma_2(\mathbf{v}^\top \mathbf{x}')].$$

1116 In this case, we have for any  $j \in [m_2]$ ,

$$1117 \quad [\mathbf{h}^{(1)}(\mathbf{x}')]_j = \frac{1}{n} \sum_{i=1}^n K_{m_2}^{(0)}(\mathbf{x}_i, \mathbf{x}') \sigma_2(\mathbf{v}_j^\top \mathbf{x}_i) \xrightarrow{m_2, n \rightarrow \infty} \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x}) K^{(0)}(\mathbf{x}, \mathbf{x}') \sigma_2(\mathbf{v}_j^\top \mathbf{x})].$$

1118 Here the infinite-inner-width kernel  $K^{(0)}$  is defined as

$$1119 \quad K^{(0)}(\mathbf{x}, \mathbf{x}') := \mathbb{E}_{\mathbf{v}} [\sigma_2(\mathbf{v}^\top \mathbf{x}) \sigma_2(\mathbf{v}^\top \mathbf{x}')] = \sum_{i=2}^{\infty} \frac{c_i^2}{B(d, i)} Q_i(\mathbf{x}^\top \mathbf{x}').$$

Recall that  $Q_2(t) = \frac{t^2-d}{d(d-1)}$ , so  $Q_2(\mathbf{v}_j^\top \mathbf{x}) = \langle \mathbf{x}\mathbf{x}^\top - \mathbf{I}, \mathbf{v}_j\mathbf{v}_j^\top - \mathbf{I} \rangle / (d(d-1))$ . Let's focus on the contribution of the quadratic term  $Q_2$  in  $K^{(0)}(\mathbf{x}, \mathbf{x}')$  and  $\sigma_2(\mathbf{v}_j^\top \mathbf{x})$ , which is

$$\begin{aligned} & \frac{c_2^3}{B(d, 2)} \cdot \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x}) Q_2(\mathbf{x}^\top \mathbf{x}') Q_2(\mathbf{v}_j^\top \mathbf{x})] \\ &= \frac{c_2^3}{B(d, 2)d^2(d-1)^2} \cdot \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x}) \langle \mathbf{x}\mathbf{x}^\top - \mathbf{I}, \mathbf{v}_j\mathbf{v}_j^\top - \mathbf{I} \rangle \langle \mathbf{x}\mathbf{x}^\top - \mathbf{I}, \mathbf{x}'\mathbf{x}'^\top - \mathbf{I} \rangle] \\ &\approx \frac{1}{d^6} \langle \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x})(\mathbf{x}\mathbf{x}^\top - \mathbf{I})^{\otimes 2}], (\mathbf{v}_j\mathbf{v}_j^\top - \mathbf{I}) \otimes (\mathbf{x}'\mathbf{x}'^\top - \mathbf{I}) \rangle. \end{aligned}$$

The following proposition provides an approximation of the tensor  $\mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x})(\mathbf{x}\mathbf{x}^\top - \mathbf{I})^{\otimes 2}]$ , which lays the foundation of our feature reconstruction theory.

**Proposition 3.** Consider two linear operators  $T$  and  $T^*$  that map  $\mathbb{R}^{d \times d}$  to  $\mathbb{R}^{d \times d}$  and satisfy

$$T(\mathbf{W}) = \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x}) \langle \mathbf{W}, \mathbf{x}\mathbf{x}^\top - \mathbf{I} \rangle (\mathbf{x}\mathbf{x}^\top - \mathbf{I})], \text{ and} \quad (14)$$

$$T^*(\mathbf{W}) = \sum_{k=1}^r \frac{1}{\|\mathbf{A}_k\|_F^2} \langle \mathbf{W}, \mathbf{A}_k \rangle \sum_{j=1}^r \mathbf{H}_{k,j} \mathbf{A}_j \quad (15)$$

for any  $\mathbf{W} \in \mathbb{R}^{d \times d}$ , where  $\mathbf{H}$  is the expected Hessian matrix  $\mathbf{H} = \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)} [\nabla^2 g^*(\mathbf{z})]$ . Then for any  $\mathbf{W} \in \mathbb{R}^{d \times d}$ , we have

$$\|T(\mathbf{W}) - T^*(\mathbf{W})\|_F \lesssim d^{-1/6} L r^2 \kappa_1 \log^2 d \cdot \|\mathbf{W}\|_F.$$

Here  $L = \tilde{\mathcal{O}}(r^{\frac{p-1}{2}})$  is the Lipschitz constant of  $g^*$  that holds with high probability.

The proof is provided in Appendix C.1.1. This proposition shows that, when  $\mathbf{H}$  is well-conditioned and  $d \gg r$ ,  $T$  can fully recover the space spanned by  $\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_r$ , which enables us to reconstruct the features efficiently. Specifically, when taking  $\mathbf{W}_k = \sum_{j=1}^r [\mathbf{H}^{-1}]_{k,j} \mathbf{A}_j$  for any  $k \in [r]$ , we have  $T(\mathbf{W}_k) \approx T^*(\mathbf{W}_k) = \mathbf{A}_k$ .

Now we consider the construction of  $\mathbf{B}^*$ . If we set  $\mathbf{B}^* = \frac{1}{m_2} [\mathbf{p}_0(\mathbf{v}_1), \dots, \mathbf{p}_0(\mathbf{v}_{m_2})]$  for some vector-valued function  $\mathbf{p}_0 : \mathbb{R}^d \rightarrow \mathbb{R}^r$  and denote  $\mathbf{p}_0(\mathbf{v}) = [p_{0,1}(\mathbf{v}), \dots, p_{0,r}(\mathbf{v})]^\top$ , we directly have for any  $k \in [r]$ ,

$$\begin{aligned} [\mathbf{B}^* \mathbf{h}^{(1)}(\mathbf{x}') ]_k &\approx \frac{1}{m_2} \sum_{j=1}^{m_2} \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x}) K^{(0)}(\mathbf{x}, \mathbf{x}') \sigma_2(\mathbf{v}_j^\top \mathbf{x}) p_{0,k}(\mathbf{v}_j)] \\ &\approx \frac{1}{d^6} \mathbb{E}_{\mathbf{v}} [\langle \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x})(\mathbf{x}\mathbf{x}^\top - \mathbf{I})^{\otimes 2}], p_{0,k}(\mathbf{v})(\mathbf{v}\mathbf{v}^\top - \mathbf{I}) \otimes (\mathbf{x}'\mathbf{x}'^\top - \mathbf{I}) \rangle] \\ &\approx \frac{1}{d^6} \langle T^*(\mathbb{E}_{\mathbf{v}} [p_{0,k}(\mathbf{v})(\mathbf{v}\mathbf{v}^\top - \mathbf{I})]), \mathbf{x}'\mathbf{x}'^\top - \mathbf{I} \rangle. \end{aligned}$$

Thus, it suffices to solve

$$T^*(\mathbb{E}_{\mathbf{v}} [p_{0,k}(\mathbf{v})(\mathbf{v}\mathbf{v}^\top - \mathbf{I})]) \propto \mathbf{A}_k, \quad k = 1, 2, \dots, r,$$

which is equivalent to solving

$$\mathbb{E}_{\mathbf{v}} [p_{0,k}(\mathbf{v})(\mathbf{v}\mathbf{v}^\top - \mathbf{I})] \propto \mathbf{W}_k = \sum_{j=1}^r [\mathbf{H}^{-1}]_{k,j} \mathbf{A}_j.$$

Since we have  $\mathbb{E}_{\mathbf{v}} [(\mathbf{v}^\top \mathbf{A}_k \mathbf{v})(\mathbf{v}\mathbf{v}^\top - \mathbf{I})] \propto \mathbf{A}_k$ , we can explicitly construct  $p_{0,k}(\mathbf{v})$  as

$$p_{0,k}(\mathbf{v}) \propto \sum_{j=1}^r [\mathbf{H}^{-1}]_{k,j} \mathbf{v}^\top \mathbf{A}_j \mathbf{v}, \text{ i.e., } \mathbf{p}_0(\mathbf{v}) \propto \mathbf{H}^{-1} \mathbf{p}(\mathbf{v}).$$

Thus, with a well conditioned  $\mathbf{H}$ , we can fully reconstruct the features.

1188 C.1.1 PROOF OF PROPOSITION 3  
1189

1190 To prove Proposition 3, it suffices to prove that the approximation error

$$1191 R(\mathbf{W}, \mathbf{V}) = \langle T(\mathbf{W}) - T^*(\mathbf{W}), \mathbf{V} \rangle \lesssim d^{-1/6} L r^2 \kappa_1 \log^2 d$$

1192 holds for any test matrix  $\mathbf{V}$  with  $\|\mathbf{V}\|_F = 1$ . We rely on the following three lemmas.

1193 **Lemma 13** (Bound  $R(\mathbf{A}_i, \mathbf{A}_j)$ ). For any  $i, j \in [r]$ , we have

$$1195 \left| \mathbb{E}_{\mathbf{x}} [g^*(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}) \langle \mathbf{A}_i, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle \langle \mathbf{A}_j, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle] - \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)} [\nabla^2 g(\mathbf{z})]_{i,j} \right| \leq \frac{L r^2 \kappa_1 \log^2 d}{\sqrt{d}}.$$

1198 Here  $L = C_g R r^{\frac{p-1}{2}}$  is the Lipschitz constant of  $g^*$  that holds with high probability.

1199 Following the proof above, we have the following more general lemma.

1200 **Lemma 14** (Bound  $R(\mathbf{A}_i, \mathbf{B})$ ). For any matrix  $\mathbf{B} \in \mathbb{R}^{d \times d}$  satisfying  $\mathbb{E}[\mathbf{x}^\top \mathbf{B} \mathbf{x}] = 0$ ,

1201  $\mathbb{E}[(\mathbf{x}^\top \mathbf{B} \mathbf{x})^2] = 1$  and  $\langle \mathbf{B}, \mathbf{A}_i \rangle = 0$  for any  $i = 1, 2, \dots, r$ , we have

$$1204 |\mathbb{E}[g^*(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}) \langle \mathbf{A}_i, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle \langle \mathbf{B}, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle]| \lesssim d^{-1/4} L r^2 \kappa_1 \log^2 d.$$

1205 **Lemma 15** (Bound  $R(\mathbf{B}_1, \mathbf{B}_2)$ ). For any two matrices  $\mathbf{B}_1, \mathbf{B}_2 \in \mathbb{R}^{d \times d}$  satisfying  $\mathbb{E}[\mathbf{x}^\top \mathbf{B}_j \mathbf{x}] = 0$ ,

1206  $\mathbb{E}[(\mathbf{x}^\top \mathbf{B}_j \mathbf{x})^2] = 1$  and  $\langle \mathbf{B}_j, \mathbf{A}_i \rangle = 0$  for any  $j = 1, 2$  and  $i = 1, 2, \dots, r$ , we have

$$1207 |\mathbb{E}[g^*(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}) \langle \mathbf{B}_1, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle \langle \mathbf{B}_2, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle]| \lesssim d^{-1/6} L r^2 \kappa_1 \log^2 d.$$

1208 Here  $L \lesssim r r^{\frac{p-1}{2}}$  is a Lipschitz constant satisfying  $\|\nabla g^*(\mathbf{p}(\mathbf{x}))\|_2 \leq L$  with high probability.

1209 The proof of the three lemmas is provided in Appendix C.1.2. With the lemmas above, we begin our  
1210 proof of Proposition 3.

1211 *Proof of Proposition 3.* Given any  $\mathbf{W} \in \mathbb{R}^{d \times d}$ , we assume  $\|\mathbf{W}\|_F = 1$  without loss of generality.  
1212 Let's decompose  $\mathbf{W}$  as

$$1217 \mathbf{W} = \sum_{k=1}^r \lambda_k \mathbf{A}_k + \lambda_{r+1} \frac{\mathbf{I}_d}{\sqrt{d}} + \lambda_{r+2} \mathbf{B}, \text{ where } \langle \mathbf{B}, \mathbf{A}_k \rangle = \langle \mathbf{B}, \mathbf{I}_d \rangle = 0, \quad k \in [r].$$

1218 Here the coefficients  $\{\lambda_k\}_{k=1}^{r+2}$  satisfy  $\sum_{k=1}^{r+2} \lambda_k^2 \lesssim 1$ , so  $\sum_{k=1}^{r+2} |\lambda_k| \lesssim \sqrt{r}$ . Since  $T^*(\mathbf{B}) = T(\mathbf{I}) = T^*(\mathbf{I}) = \mathbf{0}_{d \times d}$ , we have

$$1222 \begin{aligned} \|T(\mathbf{W}) - T^*(\mathbf{W})\|_F &\leq \sum_{k=1}^r |\lambda_k| \|T(\mathbf{A}_k) - T^*(\mathbf{A}_k)\|_F \\ 1223 &\quad + |\lambda_{r+1}| \left\| T\left(\frac{\mathbf{I}_d}{\sqrt{d}}\right) - T^*\left(\frac{\mathbf{I}_d}{\sqrt{d}}\right) \right\|_F + |\lambda_{r+2}| \|T(\mathbf{B}) - T^*(\mathbf{B})\|_F \\ 1224 &= \sum_{k=1}^r |\lambda_k| \|T(\mathbf{A}_k) - T^*(\mathbf{A}_k)\|_F + |\lambda_{r+2}| \|T(\mathbf{B})\|_F. \end{aligned}$$

1225 Since both  $T(\mathbf{A}_k)$  and  $T^*(\mathbf{A}_k)$  are traceless, by Lemma 13 and 14, we have for any  $k = 1, 2, \dots, r$ ,

$$\begin{aligned} 1226 \|T(\mathbf{A}_k) - T^*(\mathbf{A}_k)\|_F &= \max_{\|\mathbf{V}\|=1, \text{tr}(\mathbf{V})=0} \langle T(\mathbf{A}_k) - T^*(\mathbf{A}_k), \mathbf{V} \rangle \\ 1227 &\lesssim \sqrt{r} \cdot d^{-1/2} L r^2 \kappa_1 \log^2 d + d^{-1/4} L r^2 \kappa_1 \log^2 d \\ 1228 &\lesssim d^{-1/4} L r^2 \kappa_1 \log^2 d. \end{aligned}$$

1229 This is because we can decompose  $\mathbf{V} = \sum_{k=1}^r c_k \mathbf{A}_k + c_{r+1} \mathbf{B}'$  with  $\langle \mathbf{B}', \mathbf{A}_k \rangle = 0$  and apply the  
1230 two lemmas to obtain the results above. Similarly, by Lemma 14 and 15, we have

$$\begin{aligned} 1231 \|T(\mathbf{B})\|_F &= \max_{\|\mathbf{V}\|=1, \text{tr}(\mathbf{V})=0} \langle T(\mathbf{B}), \mathbf{V} \rangle \\ 1232 &\lesssim \sqrt{r} \cdot d^{-1/4} L r^2 \kappa_1 \log^2 d + d^{-1/6} L r^2 \kappa_1 \log^2 d \\ 1233 &\lesssim d^{-1/6} L r^2 \kappa_1 \log^2 d. \end{aligned}$$

1242 Thus, we have  
 1243

$$\begin{aligned} \|T(\mathbf{W}) - T^*(\mathbf{W})\|_F &\lesssim \sum_{k=1}^{r+2} |\lambda_k| d^{-1/4} L r^2 \log^2 d + |\lambda_{r+2}| d^{-1/6} L r^2 \kappa_1 \log^2 d \\ &\lesssim d^{-1/6} L r^2 \kappa_1 \log^2 d. \end{aligned}$$

1248 Here we invoke  $\sum_{k=1}^{r+2} |\lambda_k| \lesssim \sqrt{r} = o_d(1)$  in the last inequality. The proof is complete.  $\square$   
 1249

### 1250 C.1.2 OMITTED PROOFS IN APPENDIX C.1.1

1252 The following lemmas lay the foundation for our approximation process.  
 1253

1254 **Lemma 16.** Suppose Assumption 1 holds. Then the Wasserstein-1 distance between the distribution  
 1255 of  $(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x})$  and standard Gaussian  $\mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)$  can be bounded by

$$W_1((\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}), \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)) \lesssim \frac{r^2 \kappa_1}{\sqrt{d}}. \quad (16)$$

1258 Moreover, for any orthogonal unit vectors  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_s \in \mathbb{R}^d$ , we have a similar bound of  
 1259

$$W_1((\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}, \mathbf{u}_1^\top \mathbf{x}, \mathbf{u}_2^\top \mathbf{x}, \dots, \mathbf{u}_s^\top \mathbf{x}), \mathcal{N}(\mathbf{0}_{r+s}, \mathbf{I}_{r+s})) \lesssim \frac{(r+s)^2 \kappa_1}{\sqrt{d}}. \quad (17)$$

1264 *Proof of Lemma 16.* For a fixed matrix  $\mathbf{A}_i$ , define the function  $f_i(\mathbf{z}) = d \frac{\mathbf{z}^\top \mathbf{A}_i \mathbf{z}}{\|\mathbf{z}\|^2}$  and let  $\mathbf{x} = \frac{\mathbf{z} \sqrt{d}}{\|\mathbf{z}\|}$ .  
 1265 Observe that when  $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$ , we have  $\mathbf{x} \sim \text{Unif}(\mathcal{S}^{d-1}(\sqrt{d}))$ . Therefore  $[f_i(\mathbf{z})]_{i \in [r]}$  is equal in  
 1266 distribution to  $[\mathbf{x}^\top \mathbf{A}_i \mathbf{x}]_{i \in [r]}$ . We have for any  $i \in [r]$ ,

$$\nabla f_i(\mathbf{z}) = 2d \left( \frac{\mathbf{A}_i \mathbf{z}}{\|\mathbf{z}\|^2} - \frac{\mathbf{z}^\top \mathbf{A}_i \mathbf{z} \cdot \mathbf{z}}{\|\mathbf{z}\|^4} \right)$$

1271 and

$$\nabla^2 f_i(\mathbf{z}) = 2d \left( \frac{\mathbf{A}_i}{\|\mathbf{z}\|^2} - \frac{2\mathbf{A}_i \mathbf{z} \mathbf{z}^\top}{\|\mathbf{z}\|^4} - \frac{2\mathbf{z} \mathbf{z}^\top \mathbf{A}_i}{\|\mathbf{z}\|^4} - 2 \frac{\mathbf{z}^\top \mathbf{A}_i \mathbf{z}}{\|\mathbf{z}\|^4} \mathbf{I} + 4 \frac{\mathbf{z}^\top \mathbf{A}_i \mathbf{z} \mathbf{z} \mathbf{z}^\top}{\|\mathbf{z}\|^6} \right).$$

1276 Thus, we have

$$\|\nabla f_i(\mathbf{z})\| \leq 2d \left( \frac{\|\mathbf{A}_i \mathbf{z}\|}{\|\mathbf{z}\|^2} + \frac{|\mathbf{z}^\top \mathbf{A}_i \mathbf{z}|}{\|\mathbf{z}\|^3} \right) \leq \frac{\sqrt{d}}{\|\mathbf{z}\|} \cdot \|\mathbf{A}_i \mathbf{x}\| + \frac{|\mathbf{x}^\top \mathbf{A}_i \mathbf{x}|}{\|\mathbf{z}\|}.$$

1280 and

$$\|\nabla^2 f_i(\mathbf{z})\|_{op} \lesssim \frac{d}{\|\mathbf{z}\|^2} \|\mathbf{A}_i\|_{op}.$$

1284 Since  $\|\mathbf{z}\|^2$  is distributed as a chi-squared random variable with  $d$  degrees of freedom, and thus

$$\mathbb{E} [\|\mathbf{z}\|^{-2k}] = \frac{1}{\prod_{j=1}^k (d-2j)}.$$

1289 Therefore, we have

$$\mathbb{E} [\|\nabla^2 f_i(\mathbf{z})\|_{op}^4]^{1/4} \lesssim d \|\mathbf{A}_i\|_{op} \mathbb{E} [\|\mathbf{z}\|^{-8}]^{1/4} \lesssim \|\mathbf{A}_i\|_{op}.$$

1293 Then, using the fact that  $\mathbf{x}$  and  $\|\mathbf{z}\|$  are independent,

$$\mathbb{E} [\|\nabla f_i(\mathbf{z})\|^4]^{1/4} \lesssim \sqrt{d} \mathbb{E} [\|\mathbf{z}\|^{-4}]^{1/4} \mathbb{E} [\|\mathbf{A}_i \mathbf{x}\|^4]^{1/4} + \mathbb{E} [\|\mathbf{z}\|^{-4}]^{1/4} \mathbb{E} [(\mathbf{x}^\top \mathbf{A}_i \mathbf{x})^4]^{1/4} \lesssim 1.$$

1296 Thus by Lemma 1 we have

$$\begin{aligned} & W_1((\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}), \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)) \\ &= W_1(f_1(\mathbf{z}), \dots, f_r(\mathbf{z}), \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)) \\ &\lesssim \frac{4}{\sqrt{\pi}} \left( \sum_{i=1}^r \mathbb{E} [\|\nabla f_i(\mathbf{z})\|^4]^{1/4} \right) \left( \sum_{j=1}^r \mathbb{E} [\|\nabla^2 f_j(\mathbf{z})\|^4]^{1/4} \right) \\ &\lesssim \frac{r^2 \kappa_1}{\sqrt{d}}. \end{aligned}$$

1306 Now let's focus on the function  $g_j(\mathbf{z}) = \sqrt{d} \frac{\mathbf{u}_j^\top \mathbf{z}}{\|\mathbf{z}\|}$ . It holds that

$$\nabla g_j(\mathbf{z}) = \sqrt{d} \left( \frac{\mathbf{u}_j}{\|\mathbf{z}\|} - \frac{\mathbf{u}_j^\top \mathbf{z} \cdot \mathbf{z}}{\|\mathbf{z}\|^3} \right)$$

1310 and

$$\nabla^2 g_j(\mathbf{z}) = \sqrt{d} \left( -\frac{\mathbf{u}_j \mathbf{z}^\top + \mathbf{z} \mathbf{u}_j^\top}{\|\mathbf{z}\|^3} + 3 \frac{\mathbf{u}_j^\top \mathbf{z} \cdot \mathbf{z} \mathbf{z}^\top}{\|\mathbf{z}\|^5} \right).$$

1314 Thus, we have

$$\|\nabla g_j(\mathbf{z})\| \leq \frac{2\sqrt{d}}{\|\mathbf{z}\|} \quad \text{and} \quad \|\nabla^2 g_j(\mathbf{z})\|_{op} \leq \frac{5\sqrt{d}}{\|\mathbf{z}\|^2},$$

1318 which directly gives rise to

$$\mathbb{E} [\|\nabla g_j(\mathbf{z})\|^4]^{1/4} \lesssim 1 \quad \text{and} \quad \mathbb{E} [\|\nabla^2 g_j(\mathbf{z})\|_{op}^4]^{1/4} \lesssim \frac{1}{\sqrt{d}}.$$

1321 Again by Lemma 1, we have

$$\begin{aligned} & W_1((\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}, \mathbf{u}_1^\top \mathbf{x}, \mathbf{u}_2^\top \mathbf{x}, \dots, \mathbf{u}_s^\top \mathbf{x}), \mathcal{N}(\mathbf{0}_{r+s}, \mathbf{I}_{r+s})) \\ &= W_1(f_1(\mathbf{z}), \dots, f_r(\mathbf{z}), g_1(\mathbf{z}), \dots, g_s(\mathbf{z}), \mathcal{N}(\mathbf{0}_{r+s}, \mathbf{I}_{r+s})) \\ &\lesssim \frac{4}{\sqrt{\pi}} \left( \sum_{i=1}^r \mathbb{E} [\|\nabla f_i(\mathbf{z})\|^4]^{1/4} \sum_{i=1}^s \mathbb{E} [\|\nabla g_i(\mathbf{z})\|^4]^{1/4} \right) \\ &\quad \cdot \left( \sum_{j=1}^r \mathbb{E} [\|\nabla^2 f_j(\mathbf{z})\|^4]^{1/4} + \sum_{j=1}^s \mathbb{E} [\|\nabla^2 g_j(\mathbf{z})\|^4]^{1/4} \right) \\ &\lesssim \frac{(r+s)^2 \kappa_1}{\sqrt{d}}. \end{aligned}$$

1332 The proof is complete.  $\square$

1334 With the lemma above, we begin our proof of Lemma 13.

1337 *Proof of Lemma 13.* For  $\mathbf{z} \in \mathbb{R}^d$ , define  $H(\mathbf{z}) = g(\mathbf{z}) z_i z_j$ . Then by Stein's Lemma, we have

$$\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)} [H(\mathbf{z})] = \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)} [\nabla^2 g(\mathbf{z})]_{i,j} + \delta_{i,j} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)} [g(\mathbf{z})].$$

1339 Moreover, let  $R > 0$  be a truncation radius and we define  $\bar{H}(\mathbf{z}) = H(\text{clip}(\mathbf{z}, R))$ . Here the clipping 1340 function is defined as

$$\text{clip}(\mathbf{z}, R)_i = \max(\min(z_i, R), -R), \quad i = 1, 2, \dots, r.$$

1343 By (8), we know  $g(\text{clip}(\mathbf{z}, R))$  is  $\mathcal{O}(Rr^{\frac{p-1}{2}})$ -Lipschitz continuous, so  $\bar{H}$  has a Lipschitz constant of 1344  $\mathcal{O}(R^3 r^{\frac{p-1}{2}})$ . Thus, by Lemma 16, we have

$$|\mathbb{E}_{\mathbf{x}} [g(\text{clip}(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}), R)] - \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)} [g(\text{clip}(\mathbf{z}, R))]| \lesssim \frac{Rr^{\frac{p-1}{2}} \cdot r^2 \kappa_1}{\sqrt{d}},$$

$$|\mathbb{E} [\bar{H}(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x})] - \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)} [\bar{H}(\mathbf{z})]| \lesssim \frac{R^3 r^{\frac{p-1}{2}} \cdot r^2 \kappa_1}{\sqrt{d}}.$$

1350 Since  $\mathbb{E}_{\mathbf{x}}[(\mathbf{x}^\top \mathbf{A}_k \mathbf{x})^2] = 1$  for any  $k \in [r]$ , by Lemma 8, choosing  $R = C \log d$  for an appropriate  
 1351 constant  $C$  can ensure that

$$\begin{aligned} 1353 \quad & |\mathbb{E} [\bar{H}(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}) - H(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x})]| \leq \frac{1}{d}, \\ 1354 \quad & |\mathbb{E}_{\mathbf{x}} [g(\text{clip}(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}), R) - g(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x})]| \leq \frac{1}{d}, \\ 1355 \quad & |\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)} [H(\mathbf{z}) - \bar{H}(\mathbf{z})]| \leq \frac{1}{d}, \\ 1356 \quad & |\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)} [g(\text{clip}(\mathbf{z}, R)) - g(\mathbf{z})]| \leq \frac{1}{d}. \end{aligned}$$

1360 Altogether, we have

$$1362 \quad \left| \mathbb{E} [g^*(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}) \langle \mathbf{A}_i, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle \langle \mathbf{A}_j, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle] - \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)} [\nabla^2 g(\mathbf{z})]_{i,j} \right| \leq \frac{L r^2 \kappa_1 \log^2 d}{\sqrt{d}}.$$

1364 Here  $L = C_g R r^{\frac{p-1}{2}}$  for some constant  $C_g > 0$  is the Lipschitz constant of  $g^*$  that holds with high  
 1365 probability. The proof is complete.  $\square$   
 1366

1367 Following the above proof and replacing  $\mathbf{A}_i$  and  $\mathbf{A}_j$  by any other traceless matrices  $\mathbf{B}_1$  and  $\mathbf{B}_2$  that  
 1368 are orthogonal to all  $\mathbf{A}_k$ , we directly have the following corollary:

1370 **Corollary 1.** *For any two matrices  $\mathbf{B}_1, \mathbf{B}_2 \in \mathbb{R}^{d \times d}$  satisfying  $\mathbb{E} [\mathbf{x}^\top \mathbf{B}_j \mathbf{x}] = 0$  and  $\langle \mathbf{B}_j, \mathbf{A}_i \rangle = 0$   
 1371 for any  $i = 1, 2, \dots, r$  and  $j = 1, 2$ , we have*

$$\begin{aligned} 1372 \quad & |\mathbb{E} [g^*(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}) \langle \mathbf{A}_i, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle \langle \mathbf{B}_j, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle]| \\ 1373 \quad & \lesssim \left( \frac{r \kappa_1}{\sqrt{d}} + \frac{\|\mathbf{B}_j\|_{op}}{\|\mathbf{B}_j\|_F} \right) \|\mathbf{B}_j\|_F L r \log^2 d, \end{aligned}$$

1376 for any  $j = 1, 2$ , and

$$\begin{aligned} 1378 \quad & |\mathbb{E} [g^*(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}) \langle \mathbf{B}_1, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle \langle \mathbf{B}_2, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle]| \\ 1379 \quad & \lesssim \left( \frac{r \kappa_1}{\sqrt{d}} + \frac{\|\mathbf{B}_1\|_{op}}{\|\mathbf{B}_1\|_F} + \frac{\|\mathbf{B}_2\|_{op}}{\|\mathbf{B}_2\|_F} \right) \|\mathbf{B}_1\|_F \|\mathbf{B}_2\|_F L r \log^2 d. \end{aligned}$$

1382 Also, by (17) we know for any unit vector  $u \in \mathbb{R}^d$ ,  $(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \mathbf{x}^\top \mathbf{A}_2 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}, u^\top \mathbf{x})$  is ap-  
 1383 proximately Gaussian when  $d$  is sufficiently large, which gives rise to the following lemma by the  
 1384 same deduction.

1385 **Corollary 2.** *For any  $i \in [r]$  and unit vector  $\mathbf{u}_1, \mathbf{u}_2 \in \mathbb{R}^d$  and matrix  $\mathbf{B}$  satisfying the same  
 1386 requirements in Lemma 14, we have*

$$\begin{aligned} 1388 \quad & |\mathbb{E} [g^*(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}) \langle \mathbf{A}_i, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle ((\mathbf{u}_j^\top \mathbf{x})^2 - 1)]| \lesssim \frac{L r^2 \kappa_1 \log^2 d}{\sqrt{d}}, \\ 1389 \quad & |\mathbb{E} [g^*(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}) \langle \mathbf{B}, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle ((\mathbf{u}_j^\top \mathbf{x})^2 - 1)]| \lesssim \left( \frac{r \kappa_1}{\sqrt{d}} + \frac{\|\mathbf{B}_j\|_{op}}{\|\mathbf{B}_j\|_F} \right) \|\mathbf{B}_j\|_F L r \log^2 d \end{aligned}$$

1392 for any  $j = 1, 2$ , and we further have

$$1394 \quad |\mathbb{E} [g^*(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}) ((\mathbf{u}_1^\top \mathbf{x})^2 - 1)((\mathbf{u}_2^\top \mathbf{x})^2 - 1)]| \lesssim \frac{L r^2 \kappa_1 \log^2 d}{\sqrt{d}}.$$

1397 With the lemmas above, we can derive a stronger version of Corollary 1, i.e., Lemma 14, in which  
 1398 the error gets rid of the dependence on  $\|\mathbf{B}\|_{op}$ .

1400 *Proof of Lemma 14.* Let  $\tau > 1/\sqrt{d}$  be a threshold to be determined later. Decompose  $\mathbf{B}$  as follows:  
 1401

$$1402 \quad \mathbf{B} = \sum_{j=1}^d \lambda_j \mathbf{u}_j \mathbf{u}_j^\top = \sum_{|\lambda_j| > \tau} \lambda_j \left( \mathbf{u}_j \mathbf{u}_j^\top - \frac{1}{d} \mathbf{I} \right) - \sum_{k=1}^r \frac{1}{\|\mathbf{A}_k\|_F^2} \sum_{|\lambda_j| > \tau} \lambda_j \mathbf{u}_j^\top \mathbf{A}_k \mathbf{u}_j \cdot \mathbf{A}_k + \tilde{\mathbf{B}},$$

where  $\{u_j\}_{i=1}^d$  are orthogonal unit vectors and

$$\tilde{\mathbf{B}} = \sum_{|\lambda_j| \leq \tau} \lambda_j \mathbf{u}_j \mathbf{u}_j^\top + \mathbf{I} \cdot \frac{1}{d} \sum_{|\lambda_j| > \tau} \lambda_j + \sum_{k=1}^r \frac{1}{\|\mathbf{A}_k\|_F^2} \sum_{|\lambda_j| > \tau} \lambda_j \mathbf{u}_j^\top \mathbf{A}_k \mathbf{u}_j \cdot \mathbf{A}_k.$$

By construction, we have

$$\text{tr}(\tilde{\mathbf{B}}) = \sum_{|\lambda_j| \leq \tau} \lambda_j + \sum_{|\lambda_j| > \tau} \lambda_j = \sum_{j \in [d]} \lambda_j = 0.$$

Moreover, for any  $k \in [r]$ , we have

$$\langle \tilde{\mathbf{B}}, \mathbf{A}_k \rangle = \sum_{|\lambda_j| \leq \tau} \lambda_j \mathbf{u}_j^\top \mathbf{A}_k \mathbf{u}_j + \sum_{|\lambda_j| > \tau} \lambda_j \mathbf{u}_j^\top \mathbf{A}_k \mathbf{u}_j = \langle \mathbf{A}_k, \mathbf{B} \rangle = 0.$$

Therefore by Lemma 1, we have

$$\begin{aligned} & \left| \mathbb{E} \left[ g^*(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}) \langle \mathbf{A}_i, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle \langle \tilde{\mathbf{B}}, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle \right] \right| \\ & \lesssim \left( \frac{r \kappa_1}{\sqrt{d}} + \frac{\|\tilde{\mathbf{B}}\|_{op}}{\|\tilde{\mathbf{B}}\|_F} \right) \|\tilde{\mathbf{B}}\|_F L r \log^2 d. \end{aligned} \quad (18)$$

Since  $\sum_{j=1}^d \lambda_j^2 = \|\mathbf{B}\|_F^2 = (d+2)/(2d) = \mathcal{O}(1)$ , there are at most  $O(\tau^{-2})$  indices  $j$  satisfying  $|\lambda_j| > \tau$ , which gives rise to

$$\sum_{|\lambda_j| > \tau} |\lambda_j| \lesssim \sqrt{\tau^{-2} \cdot \sum_{|\lambda_j| > \tau} |\lambda_j|^2} \leq \tau^{-1}.$$

Thus, we can bound the Frobenius norm of  $\tilde{\mathbf{B}}$  by

$$\begin{aligned} \|\tilde{\mathbf{B}}\|_F^2 & \lesssim \sum_{|\lambda_j| \leq \tau} \lambda_j^2 + \frac{1}{d} \left( \sum_{\lambda_j > \tau} \lambda_j \right)^2 + \left\| \sum_{k=1}^r \frac{1}{\|\mathbf{A}_k\|_F^2} \sum_{|\lambda_j| > \tau} \lambda_j \mathbf{u}_j^\top \mathbf{A}_k \mathbf{u}_j \cdot \mathbf{A}_k \right\|_F^2 \\ & = \|\tilde{\mathbf{B}}\|_F^2 \sum_{|\lambda_j| \leq \tau} \lambda_j^2 + \frac{1}{d} \left( \sum_{\lambda_j > \tau} \lambda_j \right)^2 + \sum_{k=1}^r \left( \sum_{|\lambda_j| > \tau} \lambda_j \mathbf{u}_j^\top \mathbf{A}_k \mathbf{u}_j \right)^2 \\ & \lesssim \sum_{|\lambda_j| \leq \tau} \lambda_j^2 + \left( \frac{1}{d} + \sum_{k=1}^r \|\mathbf{A}_k\|_{op}^2 \right) \left( \sum_{|\lambda_j| > \tau} \lambda_j \right)^2 \\ & \lesssim 1 + \frac{r \kappa_1^2}{d \tau^2}. \end{aligned}$$

Thus, we have  $\|\tilde{\mathbf{B}}\|_F \lesssim 1 + \frac{\sqrt{r} \kappa_1}{\sqrt{d} \tau}$  and

$$\begin{aligned} \|\tilde{\mathbf{B}}\|_{op} & \leq \tau + \left( \frac{1}{d} + \sum_{k=1}^r \|\mathbf{A}_k\|_{op} \right) \left| \sum_{|\lambda_j| > \tau} \lambda_j u_j^\top A u_j \right| \\ & \lesssim \tau + \frac{r \kappa_1}{d \tau}. \end{aligned}$$

Thus, plugging the norm bounds into (18), we obtain that

$$\left| \mathbb{E} \left[ g^*(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}) \langle \mathbf{A}_i, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle \langle \tilde{\mathbf{B}}, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle \right] \right| \lesssim \left( \frac{r \kappa_1}{\sqrt{d}} + \frac{r^{3/2} \kappa_1^2}{d \tau} + \tau + \frac{r \kappa_1}{d \tau} \right) L r \log^2 d.$$

Next, applying Corollary 2 with  $\mathbf{u} = \mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d$ , we have

$$\left| \mathbb{E} \left[ g^*(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}) \langle \mathbf{A}_i, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle \langle \mathbf{u}_j \mathbf{u}_j^\top - \mathbf{I}/d, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle \right] \right| \lesssim \frac{L r^2 \kappa_1 \log^2 d}{\sqrt{d}}, \quad \forall j \in [d].$$

1458 Thus, we have

$$\begin{aligned}
& \left| \mathbb{E} \left[ g^*(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}) \langle \mathbf{A}_i, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle \left\langle \sum_{|\lambda_j| > \tau} \lambda_j \left( \mathbf{u}_j \mathbf{u}_j^\top - \frac{1}{d} \mathbf{I} \right), \mathbf{x} \mathbf{x}^\top - \mathbf{I} \right\rangle \right] \right| \\
& \leq \sum_{|\lambda_j| > \tau} |\lambda_j| \frac{L r^2 \kappa_1 \log^2 d}{\sqrt{d}} \\
& \lesssim \frac{L r^2 \kappa_1 \log^2 d}{\sqrt{d} \tau}.
\end{aligned}$$

1468 Besides, by Lemma 13, we have

$$\begin{aligned}
& \left| \mathbb{E} \left[ g^*(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}) \langle \mathbf{A}_i, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle \left\langle \sum_{k=1}^r \frac{1}{\|\mathbf{A}_k\|_F^2} \sum_{|\lambda_j| > \tau} \lambda_j \mathbf{u}_j^\top \mathbf{A}_k \mathbf{u}_j \cdot \mathbf{A}_k, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \right\rangle \right] \right| \\
& \lesssim \sum_{k=1}^r \left| \sum_{|\lambda_j| > \tau} \lambda_j \mathbf{u}_j^\top \mathbf{A}_k \mathbf{u}_j \right| \left( \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{o}_r, \mathbf{I}_r)} [\nabla^2 g(\mathbf{z})]_{k,i} + \frac{L r^2 \kappa_1 \log^2 d}{\sqrt{d}} \right) \\
& \lesssim \sum_{k=1}^r L \left| \sum_{|\lambda_j| > \tau} \lambda_j \mathbf{u}_j^\top \mathbf{A}_k \mathbf{u}_j \right| \\
& \lesssim \frac{L r \kappa_1}{\sqrt{d} \tau}.
\end{aligned}$$

1482 Altogether, we have

$$\begin{aligned}
& \left| \mathbb{E} [g^*(\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_r \mathbf{x}) \langle \mathbf{A}_i, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle \langle \mathbf{B}, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle] \right| \\
& \lesssim \left( \frac{r \kappa_1}{\sqrt{d}} + \frac{r^{3/2} \kappa_1^2}{d \tau} + \tau + \frac{r \kappa_1}{d \tau} \right) L r \log^2 d + \frac{L r^2 \kappa_1 \log^2 d}{\sqrt{d} \tau} + \frac{L r \kappa_1}{\sqrt{d} \tau} \\
& \lesssim d^{-1/4} L r^2 \kappa_1 \log^2 d.
\end{aligned}$$

1489 where we set  $\tau = \kappa_1 d^{-1/4}$ . The proof is complete.  $\square$

1490 Following the proof above, we can complete the proof of Lemma 15.

1493 *Proof of Lemma 15.* Similar to the proof of Lemma 14, we decompose  $\mathbf{B}_1$  and  $\mathbf{B}_2$  as follows:

$$\begin{aligned}
\mathbf{B}_i &= \sum_{j=1}^d \lambda_{i,j} \mathbf{u}_{i,j} \mathbf{u}_{i,j}^\top \\
&= \sum_{|\lambda_{i,j}| > \tau} \lambda_{i,j} \left( \mathbf{u}_{i,j} \mathbf{u}_{i,j}^\top - \frac{1}{d} \mathbf{I} \right) - \sum_{k=1}^r \frac{1}{\|\mathbf{A}_k\|_F^2} \sum_{|\lambda_{i,j}| > \tau} \lambda_{i,j} \mathbf{u}_{i,j}^\top \mathbf{A}_k \mathbf{u}_{i,j} \cdot \mathbf{A}_k + \tilde{\mathbf{B}}_i,
\end{aligned}$$

1501 where  $\{\mathbf{u}_{i,j}\}_{j=1}^d$  are orthogonal unit vectors for  $i = 1, 2$ , respectively, and

$$\tilde{\mathbf{B}}_i = \sum_{|\lambda_{i,j}| \leq \tau} \lambda_{i,j} \mathbf{u}_{i,j} \mathbf{u}_{i,j}^\top + \mathbf{I} \cdot \frac{1}{d} \sum_{|\lambda_{i,j}| > \tau} \lambda_{i,j} + \sum_{k=1}^r \frac{1}{\|\mathbf{A}_k\|_F^2} \sum_{|\lambda_{i,j}| > \tau} \lambda_{i,j} \mathbf{u}_{i,j}^\top \mathbf{A}_k \mathbf{u}_{i,j} \cdot \mathbf{A}_k.$$

1505 Then following the proof of Lemma 14, we know for any  $i = 1, 2$  and  $k = 1, 2, \dots, r$ ,

$$\text{tr}(\tilde{\mathbf{B}}_i) = \langle \tilde{\mathbf{B}}_i, \mathbf{A}_k \rangle = 0, \quad \|\tilde{\mathbf{B}}_i\|_F \lesssim 1 + \frac{\sqrt{r} \kappa_1}{\sqrt{d} \tau}, \quad \|\tilde{\mathbf{B}}_i\|_{\text{op}} \lesssim \tau + \frac{r \kappa_1}{d \tau}, \quad \sum_{|\lambda_{i,j}| > \tau} |\lambda_{i,j}| \lesssim \tau^{-1}.$$

1510 Let's denote the bi-linear operator  $\Gamma(\cdot, \cdot) : \mathbb{R}^{d \times d} \times \mathbb{R}^{d \times d} \rightarrow \mathbb{R}$  being

$$\Gamma(\mathbf{A}, \mathbf{B}) = \mathbb{E} [f^*(\mathbf{x}) \langle \mathbf{A}, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle \langle \mathbf{B}, \mathbf{x} \mathbf{x}^\top - \mathbf{I} \rangle]. \quad (19)$$

1512 By Corollary 1 and the proof of Lemma 14, we have  
 1513

$$1515 \quad \left| \Gamma(\tilde{\mathbf{B}}_1, \tilde{\mathbf{B}}_2) \right| \lesssim \left( \frac{r\kappa_1}{\sqrt{d}} \left( 1 + \frac{r\kappa_1^2}{d\tau^2} \right) + \left( \tau + \frac{r\kappa_1}{d\tau} \right) \left( 1 + \frac{\sqrt{r}}{\sqrt{d}\tau} \right) \right) Lr \log^2 d. \quad (20)$$

$$1517 \quad \left| \Gamma \left( \tilde{\mathbf{B}}_i, \sum_{|\lambda_{-i,j}| > \tau} \lambda_{-i,j} (\mathbf{u}_{-i,j} \mathbf{u}_{-i,j}^\top - \mathbf{I}/d) \right) \right| \\ 1518 \quad \lesssim \tau^{-1} \left( \frac{r\kappa_1}{\sqrt{d}} + \frac{r^{3/2}\kappa_1^2}{d\tau} + \tau + \frac{r\kappa_1}{d\tau} \right) Lr \log^2 d \quad (21)$$

$$1523 \quad \left| \Gamma \left( \tilde{\mathbf{B}}_i, \sum_{k=1}^r \frac{1}{\|\mathbf{A}_k\|_F^2} \sum_{|\lambda_{-i,j}| > \tau} \lambda_{-i,j} \mathbf{u}_{-i,j}^\top \mathbf{A}_k \mathbf{u}_{-i,j} \cdot \mathbf{A}_k \right) \right| \\ 1524 \quad \lesssim \tau^{-1} \left( \frac{r\kappa_1}{\sqrt{d}} + \frac{r^{3/2}\kappa_1^2}{d\tau} + \tau + \frac{r\kappa_1}{d\tau} \right) Lr \log^2 d. \quad (22)$$

1529 Here  $-i$  means 2 when  $i = 1$  and 1 when  $i = 2$ . Moreover, by Lemma 13, we have  
 1530

$$1532 \quad \left| \Gamma \left( \sum_{k=1}^r \frac{1}{\|\mathbf{A}_k\|_F^2} \sum_{|\lambda_{1,j}| > \tau} \lambda_{1,j} \mathbf{u}_{1,j}^\top \mathbf{A}_k \mathbf{u}_{1,j} \cdot \mathbf{A}_k, \sum_{k=1}^r \frac{1}{\|\mathbf{A}_k\|_F^2} \sum_{|\lambda_{2,j}| > \tau} \lambda_{2,j} \mathbf{u}_{2,j}^\top \mathbf{A}_k \mathbf{u}_{2,j} \cdot \mathbf{A}_k \right) \right| \\ 1533 \quad \lesssim \left( \sum_{k=1}^r \left| \sum_{|\lambda_{1,j}| > \tau} \lambda_{1,j} \mathbf{u}_{1,j}^\top \mathbf{A}_k \mathbf{u}_{1,j} \right| \right) \left( \sum_{k=1}^r \left| \sum_{|\lambda_{2,j}| > \tau} \lambda_{2,j} \mathbf{u}_{2,j}^\top \mathbf{A}_k \mathbf{u}_{2,j} \right| \right) \left( 1 + \frac{Lr^2\kappa_1 \log^2 d}{\sqrt{d}} \right) \\ 1534 \quad \lesssim \frac{r^2\kappa_1^2}{d\tau^2}, \quad (23)$$

1541 and  
 1542

$$1544 \quad \left| \Gamma \left( \sum_{k=1}^r \frac{1}{\|\mathbf{A}_k\|_F^2} \sum_{|\lambda_{i,j}| > \tau} \lambda_{i,j} \mathbf{u}_{i,j}^\top \mathbf{A}_k \mathbf{u}_{i,j} \cdot \mathbf{A}_k, \sum_{|\lambda_{-i,j}| > \tau} \lambda_{-i,j} (\mathbf{u}_{-i,j} \mathbf{u}_{-i,j}^\top - \mathbf{I}/d) \right) \right| \\ 1545 \quad \lesssim \left( \sum_{k=1}^r \left| \sum_{|\lambda_{i,j}| > \tau} \lambda_{i,j} \mathbf{u}_{i,j}^\top \mathbf{A}_k \mathbf{u}_{i,j} \right| \right) \left( \sum_{|\lambda_{-i,j}| > \tau} |\lambda_{-i,j}| \right) \frac{Lr^2\kappa_1 \log^2 d}{\sqrt{d}} \\ 1546 \quad \lesssim \frac{r\kappa_1}{\sqrt{d}\tau^2} \cdot \frac{Lr^2\kappa_1 \log^2 d}{\sqrt{d}}. \quad (24)$$

1553 Finally, we have  
 1554

$$1557 \quad \left| \Gamma \left( \sum_{|\lambda_{1,j}| > \tau} \lambda_{1,j} (\mathbf{u}_{1,j} \mathbf{u}_{1,j}^\top - \mathbf{I}/d), \sum_{|\lambda_{2,j}| > \tau} \lambda_{2,j} (\mathbf{u}_{2,j} \mathbf{u}_{2,j}^\top - \mathbf{I}/d) \right) \right| \\ 1558 \quad \lesssim \left( \sum_{|\lambda_{1,j}| > \tau} |\lambda_{1,j}| \right) \left( \sum_{|\lambda_{2,j}| > \tau} |\lambda_{2,j}| \right) \frac{Lr^2\kappa_1 \log^2 d}{\sqrt{d}}. \\ 1559 \quad \lesssim \frac{1}{\tau^2} \cdot \frac{Lr^2\kappa_1 \log^2 d}{\sqrt{d}}. \quad (25)$$

1566 Summing (20) to (25) altogether, we have  
 1567  
 1568  $|\Gamma(\mathbf{B}_1, \mathbf{B}_2)| \lesssim \left( \frac{r\kappa_1}{\sqrt{d}} \left( 1 + \frac{r\kappa_1^2}{d\tau^2} \right) + \left( \tau + \frac{r\kappa_1}{d\tau} \right) \left( 1 + \frac{\sqrt{r}}{\sqrt{d}\tau} \right) \right) Lr \log^2 d$   
 1569  
 1570  $+ \tau^{-1} \left( \frac{r\kappa_1}{\sqrt{d}} + \frac{r^{3/2}\kappa_1^2}{d\tau} + \tau + \frac{r\kappa_1}{d\tau} \right) Lr \log^2 d$   
 1571  
 1572  $+ \frac{r}{\sqrt{d}\tau} \cdot \frac{Lr^2\kappa_1 \log^2 d}{\sqrt{d}} + \frac{r^2\kappa_1^2}{d\tau^2} + \frac{r\kappa_1}{\sqrt{d}\tau^2} \cdot \frac{Lr^2\kappa_1 \log^2 d}{\sqrt{d}} + \frac{1}{\tau^2} \cdot \frac{Lr^2\kappa_1 \log^2 d}{\sqrt{d}}$   
 1573  
 1574  $\lesssim d^{-1/6} Lr^2\kappa_1 \log^2 d,$

1575 where we take  $\tau = \kappa_1 d^{-1/6}$ . The proof is complete.  $\square$   
 1576  
 1577

## 1578 C.2 BOUNDEDNESS OF THE LEARNED FEATURE

1579 In this section, we aim to upper bound the magnitude of the learned feature  $\langle \mathbf{w}, \mathbf{h}^{(1)}(\mathbf{x}') \rangle$ . Since we  
 1580 focus on the first training stage throughout this section, we denote  $n = n_1$  for notation simplicity  
 1581 when the context is clear, and let the training set be  $\mathcal{D}_1 = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ . We have the following  
 1582 proposition:

1583 **Proposition 4.** Suppose  $m_2 \geq d^4 C_\sigma^4$  and  $n \geq Ci^2 d^2$  for some sufficiently large  $C$ . With high  
 1584 probability jointly on  $\mathbf{V}$  and the training dataset  $\mathcal{D}_1$ , and with probability at least  $1 - 4n \exp(-i^2/2)$   
 1585 on  $\mathbf{w}$ , for any  $\mathbf{x}' \in \mathcal{D}_2$ , we have

$$1586 |\langle \mathbf{w}, \mathbf{h}^{(1)}(\mathbf{x}') \rangle| \lesssim \frac{i^{p+2}}{d^3} + \frac{i^{p+2}\sqrt{m_2}}{d^4\sqrt{n}} + \frac{\sqrt{m_2}i^3 \log^2(m_2 n_2)}{d^6} \cdot \left( \|\mathcal{P}_{>2}(f^*)\|_{L^2} + \sqrt{d} \|\mathcal{P}_2(f^*)\|_{L^2} \right).$$

1587 As a corollary, when  $m_2 \gtrsim d^4 i^{2p+4}$ ,  $n \gtrsim d^4 i^{2p+4}$  and  $\|\mathcal{P}_2(f^*)\|_{L^2} \lesssim \frac{\kappa_2}{\sqrt{d}}$ , we have for any  $\mathbf{x}' \in \mathcal{D}_2$ ,

$$1588 \frac{1}{\sqrt{m_2}} |\langle \mathbf{w}, \mathbf{h}^{(1)}(\mathbf{x}') \rangle| \lesssim \frac{\kappa_2 i^5}{d^6}. \quad (26)$$

1589 Thus, by taking the learning rate  $\eta = Cm_2 - 1/2\kappa_2^{-1}i^{-5}d^6$  for an appropriate constant  $C > 0$ , we  
 1590 can ensure that  $|\eta \langle \mathbf{w}, \mathbf{h}^{(1)}(\mathbf{x}') \rangle| \leq 1$  with high probability.  
 1591

1592 *Proof of Proposition 4.* Note that

$$1593 \frac{1}{m_2} \langle \mathbf{w}, \mathbf{h}^{(1)}(\mathbf{x}') \rangle = \frac{1}{m_2 n} \sum_{i=1}^n \sum_{j=1}^{m_2} w_j f^*(\mathbf{x}_i) K_{m_2}^{(0)}(\mathbf{x}_i, \mathbf{x}') \sigma_2(\mathbf{v}_j^\top \mathbf{x}_i)$$

$$1594 = \frac{1}{m_2} \sum_{j=1}^{m_2} \frac{1}{n} \sum_{i=1}^n f^*(\mathbf{x}_i) K_{m_2}^{(0)}(\mathbf{x}_i, \mathbf{x}') w_j \sigma_2(\mathbf{v}_j^\top \mathbf{x}_i).$$

1595 We do a decomposition as follows

$$1596 \frac{1}{m_2} \sum_{j=1}^{m_2} \frac{1}{n} \sum_{i=1}^n f^*(\mathbf{x}_i) K_{m_2}^{(0)}(\mathbf{x}_i, \mathbf{x}') w_j \sigma_2(\mathbf{v}_j^\top \mathbf{x}_i)$$

$$1597 = \underbrace{\frac{1}{m_2} \sum_{j=1}^{m_2} \frac{1}{n} \sum_{i=1}^n f^*(\mathbf{x}_i) \left( K_{m_2}^{(0)}(\mathbf{x}_i, \mathbf{x}') - K^{(0)}(\mathbf{x}_i, \mathbf{x}') \right) w_j \sigma_2(\mathbf{v}_j^\top \mathbf{x}_i)}_{A_1}$$

$$1598 + \underbrace{\frac{1}{m_2} \sum_{j=1}^{m_2} \frac{1}{n} \left( \sum_{i=1}^n f^*(\mathbf{x}_i) K^{(0)}(\mathbf{x}_i, \mathbf{x}') w_j \sigma_2(\mathbf{v}_j^\top \mathbf{x}_i) - \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x}) K^{(0)}(\mathbf{x}, \mathbf{x}') w_j \sigma_2(\mathbf{v}_j^\top \mathbf{x})] \right)}_{A_2}$$

$$1599 + \underbrace{\frac{1}{m_2} \sum_{j=1}^{m_2} w_j \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x}) K^{(0)}(\mathbf{x}, \mathbf{x}') \sigma_2(\mathbf{v}_j^\top \mathbf{x})]}_{A_3}.$$

1600 We consider derive an upper bound on  $A_1$ ,  $A_2$  and  $A_3$ , respectively.

1620 **Lemma 17** (Bound  $A_1$ ). Suppose  $m_2 \geq d^4 C_\sigma^4$ . With high probability jointly on  $\mathbf{V}$  and the training  
 1621 dataset  $\mathcal{D}_1$ , and with probability at least  $1 - 2n \exp(-\iota^2/2)$  on  $\mathbf{w}$ , for any  $\mathbf{x}' \in \mathcal{D}_2$ , we have  
 1622

$$1623 |A_1| \lesssim \frac{\iota^{p+2}}{m_2 d^3}. \\ 1624$$

1625 **Lemma 18** (Bound  $A_2$ ). Suppose  $m_2 \geq d^4 C_\sigma^4$  and  $n \geq C\iota^2 d^2$  for some sufficiently large  $C$ . With  
 1626 high probability on the training dataset  $\mathcal{D}_1$ , for any  $\mathbf{x}' \in \mathcal{D}_2$ , we have  
 1627

$$1628 |A_2| \lesssim \frac{\iota^{p+2}}{d^4 \sqrt{m_2 n}}. \\ 1629$$

1630 **Lemma 19** (Bound  $A_3$ ). Suppose  $m_2 \geq d^4 C_\sigma^4$  for some sufficiently large  $C$ . With high probability  
 1631 jointly on  $\mathbf{V}$  and the training dataset  $\mathcal{D}_1$ , and with probability at least  $1 - 2n_2 \exp(-\iota^2/2)$  on  $\mathbf{w}$ ,  
 1632 we have

$$1633 |A_3| \lesssim \frac{\iota^3 \log^2(m_2 n_2)}{\sqrt{m_2} d^6} \cdot \left( \|\mathcal{P}_{>2}(f)\|_{L^2} + \sqrt{d} \|\mathcal{P}_2(f)\|_{L^2} \right). \\ 1634$$

1636 Similarly, for a single point  $\mathbf{x}'$ , with high probability on  $\mathbf{V}$  and  $\mathcal{D}_1$ , with probability  $1 - 2 \exp(-\iota^2/2)$   
 1637 on  $\mathbf{w}$ , we have

$$1638 |A_3| \lesssim \frac{\iota^3 \log^2(m_2 n_2)}{\sqrt{m_2} d^6} \cdot \left( \|\mathcal{P}_{>2}(f)\|_{L^2} + \sqrt{d} \|\mathcal{P}_2(f)\|_{L^2} \right). \\ 1639$$

1641 The proof of the three lemmas are provided in Appendix C.2.1. Combining the results in the three  
 1642 lemmas above directly concludes our proof.  $\square$   
 1643

### 1644 C.2.1 OMITTED PROOFS FOR PROPOSITION 4

1645 *Proof of Lemma 17.* We can rewrite  $A_1$  as

$$1646 |A_1| = \frac{1}{n} \sum_{i=1}^n f^*(\mathbf{x}_i) \left( K_{m_2}^{(0)}(\mathbf{x}_i, \mathbf{x}') - K^{(0)}(\mathbf{x}_i, \mathbf{x}') \right) \frac{1}{m_2} \sum_{j=1}^{m_2} w_j \sigma_2(\mathbf{v}_j^\top \mathbf{x}_i)$$

1647 Since  $\frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} w_j \sigma_2(\mathbf{v}_j^\top \mathbf{x}) \sim \mathcal{N}\left(0, \frac{1}{m_2} \sum_{j=1}^{m_2} \sigma_2(\mathbf{v}_j^\top \mathbf{x})^2\right)$ , we know given any  $\mathbf{x}$  and  $\mathbf{V}$ , with  
 1648 probability at least  $1 - 2 \exp(-\iota^2/2)$  on  $\mathbf{w}$ , we have

$$1649 \left| \frac{1}{m_2} \sum_{j=1}^{m_2} w_j \sigma_2(\mathbf{v}_j^\top \mathbf{x}) \right| \leq \frac{\iota}{\sqrt{m_2}} \cdot \sqrt{\frac{1}{m_2} \sum_{j=1}^{m_2} \sigma_2(\mathbf{v}_j^\top \mathbf{x})^2}$$

1650 Moreover, by (34) in the proof of Lemma 24, we know for any  $\mathbf{x}$  and  $t > 0$ , we have

$$1651 \Pr \left[ \frac{1}{m_2} \sum_{j=1}^{m_2} \left( \sigma_2(\mathbf{v}_j^\top \mathbf{x})^2 - \mathbb{E}_{\mathbf{v}_j} [\sigma_2(\mathbf{v}_j^\top \mathbf{x})^2] \right) \geq \sqrt{\frac{t}{m_2}} \right] \\ 1652 = \Pr \left[ \left| K_{m_2}^{(0)}(\mathbf{x}, \mathbf{x}) - K^{(0)}(\mathbf{x}, \mathbf{x}) \right| \geq \sqrt{\frac{t}{m_2}} \right] \\ 1653 \leq 2 \exp \left( \frac{-t/2}{\frac{C_4}{d^4} + \frac{C_\sigma^2}{3} \sqrt{\frac{t}{m_2}}} \right).$$

1654 Altogether, when  $m_2 \geq d^4 C_\sigma^4$ , by taking  $t = C^2 \iota^2 / d^4$  for sufficiently large  $C$  and union bounding  
 1655 over the dataset  $\mathcal{D}_1$ , we can ensure that with probability at least  $1 - n \exp(-\iota)$  on  $\mathbf{V}$  and at least  
 1656  $1 - 2n \exp(-\iota^2/2)$  on  $\mathbf{w}$ , i.e., high probability on  $\mathbf{w}, \mathbf{V}$ , we have

$$1657 \left| \frac{1}{m_2} \sum_{j=1}^{m_2} w_j \sigma_2(\mathbf{v}_j^\top \mathbf{x}) \right| \leq \frac{C\iota}{\sqrt{m_2}} \sqrt{\frac{C_2}{d^2} + \frac{\iota}{\sqrt{m_2} d^2}} \leq \frac{\iota C_3}{\sqrt{m_2} d}, \quad \forall \mathbf{x} \in \mathcal{D}_1. \quad (27)$$

1674 Here  $C_3$  is a constant. We denote this joint event by  $E_1$ . On the other hand, by Lemma 24, with  
 1675 high probability on  $\mathbf{V}$ , we have for any  $\mathbf{x}_i \in \mathcal{D}_1$  and  $\mathbf{x}' \in \mathcal{D}_2$ ,

$$1677 \quad |K_{m_2}^{(0)}(\mathbf{x}_i, \mathbf{x}') - K^{(0)}(\mathbf{x}_i, \mathbf{x}')| \leq \frac{\iota}{\sqrt{m_2 d^2}}. \quad (28)$$

1679 We denote this event by  $E_2$ . Last, we truncate the range of the target function  $f$ . Denoting the  
 1680 truncation radius as  $R = (C\eta)^p$  for a sufficient large constant  $C$  and  $\eta = \log(dm_1 m_2 n_1 n_2)$   
 1681  $\Pr[|f(\mathbf{x})| \geq R] \leq 2e^{-2\eta}$  (this could be guaranteed by Lemma 4). Given  $n$  i.i.d. samples  
 1682  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \sim \mathbb{S}^{d-1}(\sqrt{d})$ , we have

$$1684 \quad \Pr[|f(\mathbf{x}_i)| \leq R, \forall i \in [n]] \geq 1 - 2ne^{-2\eta}. \quad (29)$$

1685 Thus, with high probability on the dataset  $\mathcal{D}$ , we have  $|f(\mathbf{x})| \leq \iota^p$  for any  $\mathbf{x} \in \mathcal{D}_1$ . We denote this  
 1686 event by  $E_3$ . Thus, combining (27) (32) and the truncation radius of  $f$ , we directly have

$$1688 \quad |A_1| \leq \iota^p \cdot \frac{\iota}{\sqrt{m_2 d^2}} \cdot \frac{\iota C_3}{\sqrt{m_2 d}} = \frac{C_3 \iota^{p+2}}{m_2 d^3}.$$

1690 with high probability (under events  $E_1, E_2$  and  $E_3$ ). The proof is complete.  $\square$

1693 *Proof of Lemma 18.* We rewrite  $A_2$  as

$$1695 \quad A_2 = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mathbf{x}} \left[ f^*(\mathbf{x}_i) K^{(0)}(\mathbf{x}_i, \mathbf{x}') \frac{1}{m_2} \sum_{j=1}^{m_2} w_j \sigma_2(\mathbf{v}_j^\top \mathbf{x}_i) - f^*(\mathbf{x}) K^{(0)}(\mathbf{x}, \mathbf{x}') \frac{1}{m_2} \sum_{j=1}^{m_2} w_j \sigma_2(\mathbf{v}_j^\top \mathbf{x}) \right].$$

1698 Denote  $Y(\mathbf{x}) = f^*(\mathbf{x}) K^{(0)}(\mathbf{x}, \mathbf{x}') \frac{1}{m_2} \sum_{j=1}^{m_2} w_j \sigma_2(\mathbf{v}_j^\top \mathbf{x})$ . By the proof of bounding  $A_1$ , we could  
 1699 choose the truncation radius as  $R = (C\eta)^p$  such that  $|f(\mathbf{x})| \leq R$  for all  $\mathbf{x} \in \mathcal{D}$  with high probability  
 1700  $(1 - 2ne^{-2\eta})$  on the dataset  $\mathcal{D}$ . Now we denote a truncated version of  $Y$  by  
 1701

$$1703 \quad \tilde{Y}(\mathbf{x}) = f^*(\mathbf{x}) \mathbf{1}\{f^*(\mathbf{x}) \leq R\} K^{(0)}(\mathbf{x}, \mathbf{x}') \frac{1}{m_2} \sum_{j=1}^{m_2} w_j \sigma_2(\mathbf{v}_j^\top \mathbf{x}) \mathbf{1} \left\{ \frac{1}{m_2} \sum_{j=1}^{m_2} w_j \sigma_2(\mathbf{v}_j^\top \mathbf{x}) \leq \frac{\iota C_3}{\sqrt{m_2 d}} \right\}.$$

1705 Here  $C_3$  is a constant defined in (27). Now, we decompose the concentration error as

$$1708 \quad \frac{1}{n} \sum_{i=1}^n Y(\mathbf{x}_i) - \mathbb{E}_{\mathbf{x}}[Y(\mathbf{x})] = \underbrace{\frac{1}{n} \sum_{i=1}^n (Y(\mathbf{x}_i) - \tilde{Y}(\mathbf{x}_i))}_{\mathcal{L}_0} + \underbrace{\frac{1}{n} \sum_{i=1}^n (\tilde{Y}(\mathbf{x}_i) - \mathbb{E}_{x_i}[\tilde{Y}(\mathbf{x}_i)])}_{\mathcal{L}_1} \\ 1711 \quad + \underbrace{\frac{1}{n} \sum_{i=1}^n (\mathbb{E}_{x_i}[\tilde{Y}(\mathbf{x}_i)] - \mathbb{E}_{x_i}[Y(\mathbf{x}_i)])}_{\mathcal{L}_2}.$$

1716 We know with probability at least  $1 - 2ne^{-2\eta}$  on  $\mathcal{D}$ ,  $\mathcal{L}_0 = 0$ .

1718 **Bounding  $\mathcal{L}_1$ .** We attempt to use Bernstein's type bound. First we derive a uniform upper bound  
 1719 of  $\tilde{Y}(\mathbf{x})$ . By the definition, we have

$$1721 \quad |\tilde{Y}(\mathbf{x})| \leq R \left| K^{(0)}(\mathbf{x}, \mathbf{x}') \right| \left| \frac{1}{m_2} \sum_{j=1}^{m_2} w_j \sigma_2(\mathbf{v}_j^\top \mathbf{x}) \vee \frac{\iota C_3}{\sqrt{m_2 d}} \right| \\ 1722 \quad \leq R \cdot \frac{C_2}{d^2} \frac{\iota C_3}{\sqrt{m_2 d}} \\ 1723 \quad = \frac{RC_2 C_3 \iota}{d^3 \sqrt{m_2}}.$$

Then, we bound the second moments of  $\tilde{Y}(\mathbf{x}) - \mathbb{E}_{\mathbf{x}} [\tilde{Y}(\mathbf{x})]$ , which is

$$\begin{aligned} \text{Var} [\tilde{Y}(\mathbf{x})] &\leq \mathbb{E}_{\mathbf{x}} [\tilde{Y}_k^2(\mathbf{x})] \\ &\leq r^2 \kappa_1 \mathbb{E}_{\mathbf{x}} \left[ \left( K^{(0)}(\mathbf{x}, \mathbf{x}') \right)^2 \left( \frac{1}{m_2} \sum_{j=1}^{m_2} w_j \sigma_2(\mathbf{v}_j^\top \mathbf{x}) \vee \frac{\iota C_3}{\sqrt{m_2 d}} \right)^2 \right] \\ &\leq r^2 \kappa_1 \mathbb{E}_{\mathbf{x}} \left[ \left( K^{(0)}(\mathbf{x}, \mathbf{x}') \right)^2 \right] \left( \frac{\iota C_3}{\sqrt{m_2 d}} \right)^2 \\ &\leq r^2 \kappa_1 \sum_{k=2}^{\infty} \frac{c_k^4}{B(d, k)^3} \cdot \left( \frac{\iota C_3}{\sqrt{m_2 d}} \right)^2 \\ &\leq \frac{C_4 r^2 \kappa_1 \iota^2}{m_2 d^8}. \end{aligned}$$

Here  $C_4$  is a constant. Thus, by Bernstein's inequality, we have

$$\Pr \left[ |\mathcal{L}_1| \geq \frac{R \iota}{d^4 \sqrt{m_2}} \sqrt{\frac{t}{n}} \right] \leq 2 \exp \left( \frac{-\frac{t}{2}}{C_4 + \frac{C_2 C_3}{3} \sqrt{\frac{d^2 t}{n}}} \right).$$

Thus, when  $n \geq C_\epsilon \iota^2 d^2$ , by taking  $t = \iota^2$  and  $R = (C\eta)^p \leq \iota^p$ , with high probability on  $\mathcal{D}, \mathbf{V}$  and  $\mathbf{w}$ , we have

$$|\mathcal{L}_1| \leq \frac{\iota^{p+2}}{d^4 \sqrt{m_2 n}}.$$

**Bounding  $\mathcal{L}_2$ .** It suffices to bound

$$\begin{aligned} & \left| \mathbb{E}_{\mathbf{x}} [\tilde{Y}_k(\mathbf{x})] - \mathbb{E}_{\mathbf{x}} [Y_k(\mathbf{x})] \right| \\ &\leq \mathbb{E}_{\mathbf{x}} \left[ |f^*(\mathbf{x})| \left| K^{(0)}(\mathbf{x}, \mathbf{x}') \frac{1}{m_2} \sum_{j=1}^{m_2} w_j \sigma_2(\mathbf{v}_j^\top \mathbf{x}) \right| \mathbf{1} \left\{ f^*(\mathbf{x}) > R \text{ or } \frac{1}{m_2} \sum_{j=1}^{m_2} w_j \sigma_2(\mathbf{v}_j^\top \mathbf{x}) > \frac{\iota C_3}{\sqrt{m_2 d}} \right\} \right] \\ &\leq \mathbb{E}_{\mathbf{x}} [(f^*(\mathbf{x}))^2]^{\frac{1}{2}} \Pr \left[ f^*(\mathbf{x}) > R \text{ or } \frac{1}{m_2} \sum_{j=1}^{m_2} w_j \sigma_2(\mathbf{v}_j^\top \mathbf{x}) > \frac{\iota C_3}{\sqrt{m_2 d}} \right]^{\frac{1}{4}} \mathbb{E}_{\mathbf{x}} \left[ K^{(0)}(\mathbf{x}, \mathbf{x}')^4 \right]^{\frac{1}{4}} \frac{\tau C_3}{\sqrt{m_2 d}} \\ &\lesssim \frac{(\exp(-\eta) + \exp(-\iota)) \iota C_3}{d^3 \sqrt{m_2}}. \end{aligned}$$

Taking  $\eta \geq 2 \log n + 2 \log d + \log(C_3)$  and  $\iota \geq C\eta$ , we ensure that  $\mathcal{L}_2 \leq \iota / (d^4 \sqrt{m_2 n})$ . Altogether, with high probability (event  $E_3$ ) on  $\mathcal{D}$ , we have

$$|A_2| = \left| \frac{1}{n} \sum_{i=1}^n Y(\mathbf{x}_i) - \mathbb{E}_{\mathbf{x}} [Y(\mathbf{x})] \right| \leq \frac{2\iota^{p+2}}{d^4 \sqrt{m_2 n}}.$$

The proof is complete.  $\square$

*Proof of Lemma 19.* We remember that

$$A_3 = \frac{1}{m_2} \sum_{j=1}^{m_2} w_j \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x}) K^{(0)}(\mathbf{x}, \mathbf{x}') \sigma_2(\mathbf{v}_j^\top \mathbf{x})] = \frac{1}{m_2} \sum_{j=1}^{m_2} w_j h(\mathbf{v}_j, \mathbf{x}'),$$

where  $h(\mathbf{v}, \mathbf{x}') = \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x}) K^{(0)}(\mathbf{x}, \mathbf{x}') \sigma_2(\mathbf{v}^\top \mathbf{x})]$ . To bound  $h(\mathbf{v}, \mathbf{x}')$  uniformly, we have the following lemma:

**Lemma 20.** *With high probability on  $\mathbf{V}$  and the datasets  $\mathcal{D}_1$  and  $\mathcal{D}_2$ , we have for any  $j \in [m_2]$  and  $\mathbf{x}' \in \mathcal{D}_2$ ,*

$$|h(\mathbf{v}_j, \mathbf{x}')| \lesssim \frac{\iota^2 \log^2(m_2 n_2)}{d^6} \cdot \left( \|\mathcal{P}_{>2}(f)\|_{L^2} + \sqrt{d} \|\mathcal{P}_2(f)\|_{L^2} \right)$$

The proof of Lemma 20 is deferred to the end of this section. Thus, condition on the event above, by invoking the upper bound of Gaussian tail and uniformly bounding over  $\mathbf{x}' \in \mathcal{D}_2$ , we have with probability  $1 - 2n \exp(-\iota^2/2)$  on  $\mathbf{w}$ , for any  $\mathbf{x}' \in \mathcal{D}_2$ , we have

$$|A_3| \leq \frac{\iota}{\sqrt{m_2}} \sqrt{\frac{1}{m_2} \sum_{j=1}^{m_2} h^2(\mathbf{v}_j, \mathbf{x}')} \lesssim \frac{\iota^3 \log^2(m_2 n_2)}{\sqrt{m_2} d^6} \cdot \left( \|\mathcal{P}_{>2}(f)\|_{L^2} + \sqrt{d} \|\mathcal{P}_2(f)\|_{L^2} \right).$$

Also, for a single point  $\mathbf{x}'$ , with probability  $1 - 2 \exp(-\iota^2/2)$  on  $\mathbf{w}$ , we have

$$|A_3| \lesssim \frac{\iota^3 \log^2(m_2 n_2)}{\sqrt{m_2} d^6} \cdot \left( \|\mathcal{P}_{>2}(f)\|_{L^2} + \sqrt{d} \|\mathcal{P}_2(f)\|_{L^2} \right).$$

The proof is complete.  $\square$

*Proof of Lemma 20.* Recall that the activation function  $\sigma_2$  admits a Gegenbauer expansion

$$\sigma_2(t) = \sum_{i=2}^{\infty} c_i Q_i(t).$$

Let's fix  $\mathbf{x}'$  and  $\mathbf{v}$ . Note that we can decompose  $h(\mathbf{v}, \mathbf{x}')$  as

$$\begin{aligned} h(\mathbf{v}, \mathbf{x}') &= \mathbb{E}_{\mathbf{x}} \left[ f^*(\mathbf{x}) K^{(0)}(\mathbf{x}, \mathbf{x}') \sigma_2(\mathbf{v}^\top \mathbf{x}) \right] \\ &= \mathbb{E}_{\mathbf{x}} \left[ f^*(\mathbf{x}) \sum_{i=2}^{\infty} \frac{c_i^2 Q_i(\mathbf{x}^\top \mathbf{x}')}{B(d, i)} \sum_{j=2}^{\infty} c_j Q_j(\mathbf{x}^\top \mathbf{v}) \right] \\ &= \sum_{i=2}^{\infty} \sum_{j=2}^{\infty} \mathbb{E}_{\mathbf{x}} \left[ f^*(\mathbf{x}) \frac{c_i^2}{B(d, i)^2} \langle \mathbf{Y}_i(\mathbf{x}), \mathbf{Y}_i(\mathbf{x}') \rangle \cdot \frac{c_j}{B(d, j)} \langle \mathbf{Y}_j(\mathbf{x}), \mathbf{Y}_j(\mathbf{v}) \rangle \right] \\ &= \sum_{i=2}^{\infty} \sum_{j=2}^{\infty} \frac{c_i^2 c_j}{B(d, i)^2 B(d, j)} \langle \mathbf{Y}_i(\mathbf{x}') \otimes \mathbf{Y}_j(\mathbf{v}), \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x}) \mathbf{Y}_i(\mathbf{x}) \otimes \mathbf{Y}_j(\mathbf{x})] \rangle \\ &=: \sum_{i=2}^{\infty} \sum_{j=2}^{\infty} \frac{c_i^2 c_j}{B(d, i)^2 B(d, j)} h_{i,j}(\mathbf{v}, \mathbf{x}'). \end{aligned}$$

By the definition of  $h_{i,j}(\mathbf{v}, \mathbf{x}')$ , we have

$$\mathbb{E}_{\mathbf{v}, \mathbf{x}'} [h_{i,j}^2(\mathbf{v}, \mathbf{x}')] \tag{30}$$

$$\begin{aligned} &= \mathbb{E}_{\mathbf{v}, \mathbf{x}'} \left[ \langle \mathbf{Y}_i(\mathbf{x}') \otimes \mathbf{Y}_j(\mathbf{v}), \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x}) \mathbf{Y}_i(\mathbf{x}) \otimes \mathbf{Y}_j(\mathbf{x})] \rangle^2 \right] \\ &= \langle \mathbb{E}_{\mathbf{x}'} [f^*(\mathbf{x}') \mathbf{Y}_i(\mathbf{x}') \otimes \mathbf{Y}_j(\mathbf{x}')], \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x}) \mathbf{Y}_i(\mathbf{x}) \otimes \mathbf{Y}_j(\mathbf{x})] \rangle \\ &= B(d, i) B(d, j) \mathbb{E}_{\mathbf{x}, \mathbf{x}'} [f(\mathbf{x}) f(\mathbf{x}') Q_i(\mathbf{x}^\top \mathbf{x}') Q_j(\mathbf{x}^\top \mathbf{x}')] \\ &= B(d, i) B(d, j) \mathbb{E}_{\mathbf{x}, \mathbf{x}'} \left[ f(\mathbf{x}) f(\mathbf{x}') \sum_{k=0}^{\min(i, j)} b_{i+j-2k}^{(i, j)} Q_{i+j-2k}(\mathbf{x}^\top \mathbf{x}') \right] \\ &= B(d, i) B(d, j) \mathbb{E}_{\mathbf{x}, \mathbf{x}'} \left[ f(\mathbf{x}) f(\mathbf{x}') \sum_{k=0}^{\min(i, j)} \frac{b_{i+j-2k}^{(i, j)}}{B(d, i+j-2k)} \langle \mathbf{Y}_{i+j-2k}(\mathbf{x}), \mathbf{Y}_{i+j-2k}(\mathbf{x}') \rangle \right] \\ &= B(d, i) B(d, j) \sum_{k=0}^{\min(i, j)} \frac{b_{i+j-2k}^{(i, j)}}{B(d, i+j-2k)} \|\mathbb{E}_{\mathbf{x}} [f(\mathbf{x}) \mathbf{Y}_{i+j-2k}(\mathbf{x})]\|_F^2 \\ &= B(d, i) B(d, j) \sum_{k=0}^{\min(i, j)} \frac{b_{i+j-2k}^{(i, j)}}{B(d, i+j-2k)} \|\mathcal{P}_{i+j-2k}(f)\|_{L^2}^2. \end{aligned} \tag{31}$$

1836 Since  $h_{i,j}(\mathbf{v}, \mathbf{x}')$  is a degree  $i$  polynomial of  $\mathbf{x}'$  and a degree  $j$  polynomial of  $\mathbf{v}$ , by Lemma 5, we  
 1837 have for any  $q \geq 2$ ,

$$1839 \quad \mathbb{E}_{\mathbf{v}, \mathbf{x}'} [|h_{i,j}(\mathbf{v}, \mathbf{x}')|^q]^{2/q} \leq (q-1)^{i+j} \mathbb{E}_{\mathbf{v}, \mathbf{x}'} [h_{i,j}(\mathbf{v}, \mathbf{x}')^2]$$

1840 Let  $\delta = (2e\iota \log(m_2 n_2))^{(i+j)/2}$  for some  $\iota > 1$ , taking  $q = 1 + e^{-1}\delta^{2/(i+j)}$  and Markov inequality,  
 1841 we have

$$\begin{aligned} 1843 \quad \Pr [|h_{i,j}(\mathbf{v}, \mathbf{x}')| \geq \delta \sqrt{\mathbb{E}_{\mathbf{v}, \mathbf{x}'} [h_{i,j}(\mathbf{v}, \mathbf{x}')^2]}] &\leq \frac{\mathbb{E}_{\mathbf{v}, \mathbf{x}'} [|h_{i,j}(\mathbf{v}, \mathbf{x}')|^q]}{\left(\delta \sqrt{\mathbb{E}_{\mathbf{v}, \mathbf{x}'} [h_{i,j}(\mathbf{v}, \mathbf{x}')^2]}\right)^\infty} \\ 1844 \quad &\leq (q-1)^{(i+j)q/2} \delta^{-q} \\ 1845 \quad &= \exp\left(-\frac{i+j}{2}\left(1 + \frac{2e\iota \log(m_2 n_2)}{e}\right)\right) \\ 1846 \quad &= (m_2 n_2)^{-\iota(i+j)} \exp(-(i+j)/2). \\ 1847 \end{aligned}$$

1852 Thus, with probability at least  $-(m_2 n_2)^{1-\iota(i+j)} \exp(-(i+j)/2)$ ,

$$\begin{aligned} 1853 \quad h_{i,j}(\mathbf{v}_i, \mathbf{x}') &\leq (2e\iota \log(m_2 n_2))^{(i+j)/2} \sqrt{\mathbb{E}_{\mathbf{v}, \mathbf{x}'} [h_{i,j}(\mathbf{v}, \mathbf{x}')^2]} \\ 1854 \quad &\leq (2e\iota \log(m_2 n_2))^{(i+j)/2} \sqrt{\sum_{k=0}^{[(i+j)/2]} \frac{B(d, i)B(d, j)b_{i+j-2k}^{(i,j)}}{B(d, i+j-2k)} \|\mathcal{P}_{i+j-2k}(f)\|_{L^2}^2}. \\ 1855 \quad 1856 \quad 1857 \quad 1858 \end{aligned}$$

1859 In the second inequality we invoke (31). Summing over  $i$  and  $j$  gives rise to

$$\begin{aligned} 1860 \quad &|h(\mathbf{v}, \mathbf{x}')| \\ 1861 \quad &= \left| \sum_{i=2}^{\infty} \sum_{j=2}^{\infty} \frac{c_i^2 c_j}{B(d, i)^2 B(d, j)} h_{i,j}(\mathbf{v}, \mathbf{x}') \right| \\ 1862 \quad &\leq \sum_{i=2}^{\infty} \sum_{j=2}^{\infty} \frac{c_i^2 |c_j| (2e\iota \log(m_2 n_2))^{(i+j)/2}}{B(d, i)^2 B(d, j)} \sqrt{\sum_{k=0}^{[(i+j)/2]} \frac{B(d, i)B(d, j)b_{i+j-2k}^{(i,j)}}{B(d, i+j-2k)} \|\mathcal{P}_{i+j-2k}(f)\|_{L^2}^2} \\ 1863 \quad &\leq \sqrt{\sum_{i=2}^{\infty} \sum_{j=2}^{\infty} \frac{c_i^2 |c_j| (2e\iota \log(m_2 n_2))^{(i+j)/2}}{B(d, i)^2 B(d, j)^{1/2}}} \\ 1864 \quad &\cdot \sqrt{\sum_{i=2}^{\infty} \sum_{j=2}^{\infty} \frac{c_i^2 |c_j| (2e\iota \log(m_2 n_2))^{(i+j)/2}}{B(d, i)^2 B(d, j)^{3/2}} \sum_{k=0}^{[(i+j)/2]} \frac{B(d, i)B(d, j)b_{i+j-2k}^{(i,j)}}{B(d, i+j-2k)} \|\mathcal{P}_{i+j-2k}(f)\|_{L^2}^2} \\ 1865 \quad &\lesssim \frac{\iota \log(m_2 n_2)}{d^{5/2}} \cdot \sqrt{\sum_{\ell=0}^{\infty} \frac{1}{B(d, \ell)} \sum_{\substack{2 \leq i, j \\ i+j-\ell \text{ even}}}^{\infty} \frac{c_i^2 |c_j| b_{\ell}^{(i,j)} (2e\iota \log(m_2 n_2))^{(i+j)/2}}{B(d, i)B(d, j)^{1/2}} \|\mathcal{P}_{\ell}(f)\|_{L^2}^2}. \\ 1866 \quad 1867 \quad 1868 \quad 1869 \quad 1870 \quad 1871 \quad 1872 \quad 1873 \quad 1874 \quad 1875 \quad 1876 \quad 1877 \quad 1878 \end{aligned}$$

1879 In the second inequality we invoke Cauchy inequality. Then by plugging the bound on  $b_{\ell}^{(i,j)}$  in  
 1880 Lemma 12, we have

$$\begin{aligned} 1881 \quad &|h(\mathbf{v}, \mathbf{x}')| \\ 1882 \quad &\lesssim \frac{\iota \log(m_2 n_2)}{d^{5/2}} \cdot \sqrt{\sum_{\ell=0}^{\infty} \frac{4(2\ell+d-2)}{B(d, \ell)(d-2)} \sum_{i+j-\ell=2k}^{i, j \geq \max(k, 2)} \frac{c_i^2 |c_j| (2e\iota \log(m_2 n_2))^{(i+j)/2}}{B(d, i)B(d, j)^{1/2}(d-2)_k} \binom{i}{k} \binom{j}{k} k! \|\mathcal{P}_{\ell}(f)\|_{L^2}^2} \\ 1883 \quad &\lesssim \frac{\iota \log(m_2 n_2)}{d^{5/2}} \cdot \left( \frac{\iota \log(m_2 n_2)}{d^{7/2}} \cdot \|\mathcal{P}_{>2}(f)\|_{L^2} + \frac{\iota \log(m_2 n_2)}{d^3} \cdot \|\mathcal{P}_2(f)\|_{L^2} \right) \\ 1884 \quad &= \frac{\iota^2 \log^2(m_2 n_2)}{d^6} \cdot \|\mathcal{P}_{>2}(f)\|_{L^2} + \frac{\iota^2 \log^2(m_2 n_2)}{d^{11/2}} \cdot \|\mathcal{P}_2(f)\|_{L^2}. \\ 1885 \quad 1886 \quad 1887 \quad 1888 \quad 1889 \end{aligned}$$

1890 The probability of this event is at least  
 1891

$$1892 \quad 1 - \sum_{i=2}^{\infty} \sum_{j=2}^{\infty} (m_2 n_2)^{-\iota(i+j)} \exp((i+j)/2) = 1 - \frac{m_2 n_2 (- (m_2 n_2)^{-\iota} e^{-1/2})^2}{(m_2 n_2)^{4\iota} e^2}, \\ 1893 \\ 1894$$

1895 which is a high probability event when uniformly bounding over  $\mathbf{x}' \in \mathcal{D}_2$  and  $\mathbf{v} = \mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{m_2}$ .  
 1896 The proof is complete.  $\square$   
 1897

### 1898 C.3 PROOF OF PROPOSITION 1

#### 1899 C.3.1 THE FORMAL STATEMENT OF PROPOSITION 1 AND THE COROLLARY

1900 Let's consider a formal version of Proposition 1. We remind the readers that throughout this section  
 1901 we denote  $n = n_1$  for notation simplicity, since we only focus on the first training stage.  
 1902

1903 **Proposition 5** (Reconstruct the feature). *Suppose  $m_2, n \geq Cd^4$  for some sufficiently large  $C$ . With  
 1904 high probability jointly on  $\mathbf{V}$  and the training datasets  $\mathcal{D}_1$  and  $\mathcal{D}_2$ , there exists a matrix  $\mathbf{B}^* \in$   
 1905  $\mathbb{R}^{r \times m_2}$  satisfying  $\|\mathbf{B}^*\|_{\text{op}} \lesssim \frac{d^6}{\lambda_{\min}(\mathbf{H})} \sqrt{\frac{1}{m_2}}$  such that for any  $\mathbf{x}' \in \mathcal{D}_2$ , we have  
 1906*

$$1907 \quad \left\| \mathbf{B}^* \mathbf{h}^{(1)}(\mathbf{x}') - \mathbf{p}(\mathbf{x}') \right\|_2 \lesssim \frac{\sqrt{r}}{\lambda_{\min}(\mathbf{H})} \cdot \left( \frac{\iota^{p+2} d^5}{m_2} + \frac{\iota d^3}{\sqrt{m_2}} + \frac{\iota^{p+3/2} d}{\sqrt{n}} + \frac{\iota L r^2 \kappa_1 \log^2 d}{d^{1/6}} \right). \\ 1908 \\ 1909$$

1910 Here  $L \lesssim \iota r^{\frac{p-1}{2}}$  is a Lipschitz constant satisfying  $\|\nabla g^*(\mathbf{p}(\mathbf{x}))\|_2 \leq L$  with high probability.  
 1911

1912 With the proposition above, we directly have the following result.

1913 **Corollary 3.** *Under the same assumption in Proposition 5, with high probability, we have*

$$1915 \quad \sup_{\mathbf{x} \in \mathcal{D}_2} \left| g(\mathbf{B}^* \mathbf{h}^{(1)}(\mathbf{x})) - g(\mathbf{p}(\mathbf{x})) \right| \lesssim \|g\|_{L^2} \cdot \frac{r^{p/2}}{\lambda_{\min}(\mathbf{H})} \cdot \left( \frac{\iota^{p+2} d^5}{m_2} + \frac{\iota d^3}{\sqrt{m_2}} + \frac{\iota^{p+3/2} d}{\sqrt{n}} + \frac{\iota L r^2 \kappa_1 \log^2 d}{d^{1/6}} \right). \\ 1916 \\ 1917$$

1918 We provide the main proof of Proposition 5 in Appendix C.3.2, and defer the proof of Corollary 3  
 1919 and other supporting lemmas to Appendix C.3.3.  
 1920

#### 1921 C.3.2 PROOF OF PROPOSITION 5

1922 *Proof.* Denote the target features by  $\mathbf{p}(\mathbf{v}) = [\mathbf{v}^\top \mathbf{A}_1 \mathbf{v}, \dots, \mathbf{v}^\top \mathbf{A}_r \mathbf{v}]^\top \in \mathbb{R}^r$  for any  $\mathbf{v} \in \mathbb{R}^d$ , and  
 1923 we further let  $\mathbf{P} = [\mathbf{p}(\mathbf{v}_1), \mathbf{p}(\mathbf{v}_2), \dots, \mathbf{p}(\mathbf{v}_{m_2})]^\top \in \mathbb{R}^{m_2 \times r}$ . Then for any  $\mathbf{x}' \in \mathcal{D}_2$ , we have the  
 1924 following decomposition

$$1925 \quad \frac{1}{m_2} \mathbf{P}^\top \mathbf{h}^{(1)}(\mathbf{x}') = \frac{1}{m_2 n} \sum_{i=1}^n \sum_{j=1}^{m_2} f^*(\mathbf{x}_i) K_{m_2}^{(0)}(\mathbf{x}_i, \mathbf{x}') \sigma_2(\mathbf{v}_j^\top \mathbf{x}_i) \mathbf{p}(\mathbf{v}_j) \\ 1926 \\ 1927 \quad = \underbrace{\frac{1}{n} \sum_{i=1}^n \frac{1}{m_2} \sum_{j=1}^{m_2} f^*(\mathbf{x}_i) \left( K_{m_2}^{(0)}(\mathbf{x}_i, \mathbf{x}') - K^{(0)}(\mathbf{x}_i, \mathbf{x}') \right) \sigma_2(\mathbf{v}_j^\top \mathbf{x}_i) \mathbf{p}(\mathbf{v}_j)}_{\mathbf{D}_{1,1}} \\ 1928 \\ 1929 \quad + \underbrace{\frac{1}{n} \sum_{i=1}^n f^*(\mathbf{x}_i) K^{(0)}(\mathbf{x}_i, \mathbf{x}') \left( \frac{1}{m_2} \sum_{j=1}^{m_2} \sigma_2(\mathbf{v}_j^\top \mathbf{x}_i) \mathbf{p}(\mathbf{v}_j) - \frac{c_2}{B(d, 2)} \mathbf{p}(\mathbf{x}_i) \right)}_{\mathbf{D}_{1,2}} \\ 1930 \\ 1931 \quad + \underbrace{\frac{c_2}{n B(d, 2)} \left( \sum_{i=1}^n f^*(\mathbf{x}_i) K^{(0)}(\mathbf{x}_i, \mathbf{x}') \mathbf{p}(\mathbf{x}_i) - \mathbb{E}_{\mathbf{x}} \left[ f^*(\mathbf{x}) K^{(0)}(\mathbf{x}, \mathbf{x}') \mathbf{p}(\mathbf{x}) \right] \right)}_{\mathbf{D}_2} \\ 1932 \\ 1933 \quad + \underbrace{\frac{c_2}{B(d, 2)} \mathbb{E}_{\mathbf{x}} \left[ f^*(\mathbf{x}) K^{(0)}(\mathbf{x}, \mathbf{x}') \mathbf{p}(\mathbf{x}) \right]}_{\mathbf{D}_3}. \\ 1934 \\ 1935 \\ 1936 \\ 1937 \\ 1938 \\ 1939 \\ 1940 \\ 1941 \\ 1942 \\ 1943$$

We will derive an upper bound on the concentration error terms  $\mathbf{D}_{1,1}$ ,  $\mathbf{D}_{1,2}$  and  $\mathbf{D}_2$ , respectively. Moreover, leveraging the asymptotic analysis in Appendix C.1, we show that  $\mathbf{D}_3 \approx d^{-6}\mathbf{H}\mathbf{p}(\mathbf{x}')$  with high probability,

**Lemma 21** (Bound  $\mathbf{D}_{1,1}$  and  $\mathbf{D}_{1,2}$ ). *Under the same assumptions in Proposition 5, with high probability on  $\mathbf{V}$ ,  $\mathcal{D}_1$  and  $\mathcal{D}_2$ , we have*

$$\|\mathbf{D}_{1,1}\|_\infty \leq \frac{9C_4^{1/4}\iota^{p+2}}{m_2d} \quad \text{and} \quad \|\mathbf{D}_{1,2}\|_\infty \leq \frac{9\iota C_4^{1/4}C_2}{\sqrt{m_2}d^3}.$$

**Lemma 22** (Bound  $\mathbf{D}_2$ ). *Under the same assumptions in Proposition 5, with high probability on  $\mathcal{D}_1$  and  $\mathcal{D}_2$ , we have*

$$\|\mathbf{D}_2\|_\infty \lesssim \frac{\iota^{p+3/2}}{\sqrt{nd^5}}.$$

**Lemma 23** (Compute  $\mathbf{D}_3$ ). *Under the same assumptions in Proposition 5, with high probability on  $\mathcal{D}_2$ , for any  $\mathbf{x}' \in \mathcal{D}_2$ , we have*

$$\left\| \mathbf{D}_3 - \frac{c_2^2}{B(d, 2)^2 d(d-1)} \cdot \mathbf{H}\mathbf{p}(\mathbf{x}') \right\|_\infty \lesssim \frac{\iota L r^2 \kappa_1 \log^2 d}{d^{6+1/6}}.$$

We defer the detailed proof of the three lemmas to Appendix C.3.3. Combining all the results above and choosing

$$\mathbf{B}^* = \frac{B(d, 2)^2 d(d-1)}{c_2^2} \cdot \frac{1}{m_2} \mathbf{H}^{-1} \mathbf{P}^\top,$$

we have with high probability on  $\mathbf{V}$ ,  $\mathcal{D}_1$  and  $\mathcal{D}_2$ ,

$$\begin{aligned} & \left\| \mathbf{B}^* \mathbf{h}^{(1)}(\mathbf{x}') - \mathbf{p}(\mathbf{x}') \right\|_2 \\ & \leq \frac{B(d, 2)^2 d(d-1)}{c_2^2} \cdot \left\| \mathbf{H}^{-1} \left( \mathbf{D}_{1,1} + \mathbf{D}_{1,2} + \mathbf{D}_2 + \mathbf{D}_3 - \frac{c_2^2}{B(d, 2)d(d-1)} \cdot \mathbf{H}\mathbf{p}(\mathbf{x}') \right) \right\|_2 \\ & \lesssim \frac{d^6 \sqrt{r}}{\lambda_{\min}(\mathbf{H})} \cdot \left( \|\mathbf{D}_{1,1}\|_\infty + \|\mathbf{D}_{1,2}\|_\infty + \|\mathbf{D}_2\|_\infty \right. \\ & \quad \left. + \left\| \mathbf{D}_3 - \frac{c_2^2}{B(d, 2)d(d-1)} \cdot \mathbf{H}\mathbf{p}(\mathbf{x}') \right\|_\infty \right) \\ & \lesssim \frac{\sqrt{r}}{\lambda_{\min}(\mathbf{H})} \cdot \left( \frac{\iota^{p+2} d^5}{m_2} + \frac{\iota d^3}{\sqrt{m_2}} + \frac{\iota^{p+3/2} d}{\sqrt{n}} + \frac{\iota L r^2 \kappa_1 \log^2 d}{d^{1/6}} \right), \end{aligned}$$

To bound  $\|\mathbf{B}^*\|_{\text{op}}$ , note that

$$\begin{aligned} \|\mathbf{B}^*\|_{\text{op}}^2 &= \|\mathbf{B}^* \mathbf{B}^{*\top}\|_{\text{op}} \\ &\lesssim \frac{d^{12}}{m_2^2 \lambda_{\min}^2(\mathbf{H})} \|\mathbf{P} \mathbf{P}^\top\|_{\text{op}} \\ &= \frac{d^{12}}{m_2 \lambda_{\min}^2(\mathbf{H})} \left\| \frac{1}{m_2} \sum_{j=1}^{m_2} \mathbf{p}(\mathbf{v}_j) \mathbf{p}(\mathbf{v}_j)^\top \right\|_{\text{op}}. \end{aligned}$$

Moreover, for any  $j \in [m_2]$ , we have

$$\|\mathbf{p}(\mathbf{v}_j) \mathbf{p}(\mathbf{v}_j)^\top\|_{\text{op}} = \|\mathbf{p}(\mathbf{v}_j)\|_2^2 = \sum_{k=1}^r (\mathbf{v}_j^\top \mathbf{A}_k \mathbf{v}_j)^2 \lesssim rd^2,$$

1998 and we have  
 1999

$$\begin{aligned}
 2000 \quad & \left\| \mathbb{E}_{\mathbf{v}} \left[ (\mathbf{p}(\mathbf{v})\mathbf{p}(\mathbf{v})^\top)^2 \right] \right\|_{\text{op}} = \left\| \mathbb{E}_{\mathbf{v}} \left[ \sum_{k=1}^r (\mathbf{v}^\top \mathbf{A}_k \mathbf{v})^2 \mathbf{p}(\mathbf{v})\mathbf{p}(\mathbf{v})^\top \right] \right\|_{\text{op}} \\
 2001 \quad & \leq \sum_{k=1}^r \left\| \mathbb{E}_{\mathbf{v}} \left[ (\mathbf{v}^\top \mathbf{A}_k \mathbf{v})^2 \mathbf{p}(\mathbf{v})\mathbf{p}(\mathbf{v})^\top \right] \right\|_{\text{op}} \\
 2002 \quad & \leq \sum_{k=1}^r d^2 \left\| \mathbb{E}_{\mathbf{v}} \left[ \mathbf{p}(\mathbf{v})\mathbf{p}(\mathbf{v})^\top \right] \right\|_{\text{op}} \\
 2003 \quad & \leq \sum_{k=1}^r d^2 \\
 2004 \quad & = rd^2.
 \end{aligned}$$

2005 The second inequality holds because  $\mathbf{p}(\mathbf{v})\mathbf{p}(\mathbf{v})^\top$  is positive semi-definite. By Matrix Bernstein  
 2006 Inequality, we have  
 2007

$$\begin{aligned}
 2008 \quad & \Pr \left[ \left\| \frac{1}{m_2} \sum_{j=1}^{m_2} \mathbf{p}(\mathbf{v}_j)\mathbf{p}(\mathbf{v}_j)^\top - \mathbf{I} \right\|_{\text{op}} \geq 1 + \frac{\sqrt{rd}\iota}{\sqrt{m_2}} \right] \leq \exp \left( -\frac{\frac{rd^2\iota^2}{2m_2}}{\frac{rd^2}{m_2} + \frac{rd^2}{3m_2} \cdot \frac{\sqrt{rd}\iota}{\sqrt{m_2}}} \right) \\
 2009 \quad & = \exp \left( -\frac{\frac{\iota^2}{2}}{1 + \frac{\sqrt{rd}\iota}{3\sqrt{m_2}}} \right).
 \end{aligned}$$

2010 Thus, when  $m_2 \geq d^4$ , we know with high probability on  $\mathbf{V}$ ,  
 2011

$$\left\| \frac{1}{m_2} \sum_{j=1}^{m_2} \mathbf{p}(\mathbf{v}_j)\mathbf{p}(\mathbf{v}_j)^\top \right\|_{\text{op}} \leq 1 + \frac{\sqrt{rd}\iota}{\sqrt{m_2}} \lesssim 1.$$

2012 Thus, we have  $\|\mathbf{B}^*\|_{\text{op}} \lesssim \frac{d^6}{\lambda_{\min}(\mathbf{H})} \sqrt{\frac{1}{m_2}}$ . The proof is complete.  $\square$   
 2013

### 2014 C.3.3 OMITTED PROOFS IN APPENDICES C.3.1 AND C.3.2

2015 *Proof of Lemma 21.* Let's first bound  $\mathbf{D}_{1,1}$ . We can rewrite  $\mathbf{D}_{1,1}$  as  
 2016

$$\mathbf{D}_{1,1} = \frac{1}{n} \sum_{i=1}^n f^*(\mathbf{x}_i) \left( K_{m_2}^{(0)}(\mathbf{x}_i, \mathbf{x}') - K^{(0)}(\mathbf{x}_i, \mathbf{x}') \right) \frac{1}{m_2} \sum_{j=1}^{m_2} \sigma_2(\mathbf{v}_j^\top \mathbf{x}_i) \mathbf{p}(\mathbf{v}_j)$$

2017 By Lemma 25, for any  $k \in [r]$ , we have with high probability on  $\mathbf{V}$

$$\left| \frac{1}{m_2} \sum_{j=1}^{m_2} (\mathbf{v}_j^\top \mathbf{A}_k \mathbf{v}_j) \sigma_2(\mathbf{v}_j^\top \mathbf{x}_i) - \frac{c_2}{B(d, 2)} \mathbf{x}_i^\top \mathbf{A}_k \mathbf{x}_i \right| \leq \frac{9\iota d^{-1} C_4^{1/4}}{\sqrt{m_2}}.$$

2018 Thus, by enumerating  $\mathbf{A}_k$  over  $\{\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_r\}$ , we have with high probability on  $\mathbf{V}$ ,  
 2019

$$\left\| \frac{1}{m_2} \sum_{j=1}^{m_2} \mathbf{p}(\mathbf{v}_j) \sigma_2(\mathbf{v}_j^\top \mathbf{x}_i) - \frac{c_2}{B(d, 2)} \mathbf{p}(\mathbf{x}_i) \right\|_\infty \leq \frac{9\iota d^{-1} C_4^{1/4}}{\sqrt{m_2}}.$$

2020 On the other hand, by Lemma 24, with high probability on  $\mathbf{V}$ , we have for any  $\mathbf{x}_i \in \mathcal{D}_1, \mathbf{x}' \in \mathcal{D}_2$ ,

$$\left| K_{m_2}^{(0)}(\mathbf{x}_i, \mathbf{x}') - K^{(0)}(\mathbf{x}_i, \mathbf{x}') \right| \leq \frac{\iota}{\sqrt{m_2} d^2}. \quad (32)$$

2021 Moreover, under the event  $E_3$  (defined in (29)), with high probability on the dataset  $\mathcal{D}_1$ , we have  
 2022  $|f(\mathbf{x})| \leq \iota^p$  for any  $\mathbf{x} \in \mathcal{D}_1$ . Thus, altogether we have  
 2023

$$\|\mathbf{D}_{1,1}\|_\infty \leq \iota^p \cdot \frac{\iota}{\sqrt{m_2} d^2} \cdot \frac{9\iota d^{-1} C_4^{1/4}}{\sqrt{m_2}} = \frac{9C_4^{1/4} \iota^{p+2}}{m_2 d}.$$

with high probability. To bound  $\mathbf{D}_{1,2}$ , from the proof above, we know with high probability,

$$\left\| \frac{1}{m_2} \sum_{j=1}^{m_2} \mathbf{p}(\mathbf{v}_j) \sigma_2(\mathbf{v}_j^\top \mathbf{x}_i) - \frac{c_2}{B(d, 2)} \mathbf{p}(\mathbf{x}_i) \right\|_\infty \leq \frac{9\iota d^{-1} C_4^{1/4}}{\sqrt{m_2}}.$$

Moreover, for any  $\mathbf{x} \in \mathcal{D}_1$  and  $\mathbf{x}' \in \mathcal{D}_2$ ,

$$K^{(0)}(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\mathbf{v}} [\sigma_2(\mathbf{v}^\top \mathbf{x}) \sigma_2(\mathbf{v}^\top \mathbf{x}')] \leq \sqrt{\mathbb{E}_{\mathbf{v}} [\sigma_2(\mathbf{v}^\top \mathbf{x})^2] \mathbb{E}_{\mathbf{v}} [\sigma_2(\mathbf{v}^\top \mathbf{x}')^2]} \leq \frac{C_2}{d^2}. \quad (33)$$

Thus, we can bound  $\mathbf{D}_{1,2}$  with high probability by

$$\|\mathbf{D}_{1,2}\|_\infty \leq \frac{C_2}{d^2} \cdot \frac{9\iota d^{-1} C_4^{1/4}}{\sqrt{m_2}} = \frac{9\iota C_4^{1/4} C_2}{\sqrt{m_2} d^3}.$$

The proof is complete.  $\square$

*Proof of Lemma 22.* Thus, let's focus on the concentration of a single element

$$Y_k(\mathbf{x}) = f^*(\mathbf{x}) K^{(0)}(\mathbf{x}, \mathbf{x}') \mathbf{x}^\top \mathbf{A}_k \mathbf{x}, \quad k = 1, 2, \dots, r.$$

Similar to the proof of Lemma 17, we denote a truncated version of  $Y_k$  by

$$\tilde{Y}_k(\mathbf{x}) = f^*(\mathbf{x}) \mathbf{1}\{f^*(\mathbf{x}) \leq R\} K^{(0)}(\mathbf{x}, \mathbf{x}') \mathbf{x}^\top \mathbf{A}_k \mathbf{x}, \quad k = 1, 2, \dots, r.$$

Here,  $R = (C\eta)^p$  for some large constant  $C$ . Now, we decompose the concentration error as

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n Y_k(\mathbf{x}_i) - \mathbb{E}_{\mathbf{x}} [Y_k(\mathbf{x})] &= \underbrace{\frac{1}{n} \sum_{i=1}^n (Y_k(\mathbf{x}_i) - \tilde{Y}_k(\mathbf{x}_i))}_{\mathcal{L}_0} + \underbrace{\frac{1}{n} \sum_{i=1}^n (\tilde{Y}_k(\mathbf{x}_i) - \mathbb{E}_{x_i} [\tilde{Y}_k(\mathbf{x}_i)])}_{\mathcal{L}_1} \\ &\quad + \underbrace{\frac{1}{n} \sum_{i=1}^n (\mathbb{E}_{x_i} [\tilde{Y}_k(\mathbf{x}_i)] - \mathbb{E}_{x_i} [Y_k(\mathbf{x}_i)])}_{\mathcal{L}_2}. \end{aligned}$$

By (29), we know with probability at least  $1 - 2ne^{-2\eta}$ ,  $\mathcal{L}_0 = 0$ .

**Bounding  $\mathcal{L}_1$ .** First we derive a uniform upper bound of  $\tilde{Y}_k(\mathbf{x})$ , which is

$$\begin{aligned} |Y_k(\mathbf{x})| &\leq R |K^{(0)}(\mathbf{x}, \mathbf{x}') \mathbf{x}^\top \mathbf{A}_k \mathbf{x}| \\ &\leq R |K^{(0)}(\mathbf{x}, \mathbf{x}')| |\mathbf{x}^\top \mathbf{A}_k \mathbf{x}| \\ &\leq R \cdot \frac{C_2}{d^2} d \|\mathbf{A}_k\|_{op} \\ &= \frac{RC_2 \|\mathbf{A}_k\|_{op}}{d}. \end{aligned}$$

Then, we bound the second moments of  $\tilde{Y}_k(\mathbf{x}) - \mathbb{E}_{\mathbf{x}} [\tilde{Y}_k(\mathbf{x})]$ . Again by Lemma 8, we know that there exists a sufficient large constant  $C > 0$  s.t.  $\Pr[|\mathbf{x}^\top \mathbf{A}_k \mathbf{x}| \geq Ct] \leq 2\exp(-\iota)$ . By taking  $\iota \geq 2\log d$ , we have

$$\begin{aligned} \text{Var} [\tilde{Y}_k(\mathbf{x})] &\leq \mathbb{E}_{\mathbf{x}} [\tilde{Y}_k^2(\mathbf{x})] \\ &= \mathbb{E}_{\mathbf{x}} [\tilde{Y}_k^2(\mathbf{x}) \mathbf{1}\{|\mathbf{x}^\top \mathbf{A}_k \mathbf{x}| \leq Ct\}] + \mathbb{E}_{\mathbf{x}} [\tilde{Y}_k^2(\mathbf{x}) \mathbf{1}\{|\mathbf{x}^\top \mathbf{A}_k \mathbf{x}| > Ct\}] \\ &\leq C^2 r^2 \kappa_1 \iota^2 \mathbb{E}_{\mathbf{x}} \left[ \left( K^{(0)}(\mathbf{x}, \mathbf{x}') \right)^2 \right] + \frac{r^2 \kappa_1 C_2^2}{d^4} \cdot \mathbb{E}_{\mathbf{x}} [\mathbf{1}\{|\mathbf{x}^\top \mathbf{A}_k \mathbf{x}| > Ct\}] \\ &\leq C^2 r^2 \kappa_1 \iota^2 \cdot \sum_{i=2}^{\infty} \frac{c_i^4}{B(d, i)^3} + \frac{2r^2 \kappa_1 C_2^2 \exp(-\iota)}{d^4} \\ &\lesssim \frac{C' r^2 \kappa_1 \iota^2}{d^6}. \end{aligned}$$

Here  $C'$  is a sufficiently large constant independent of  $d$ . We invoke (33) in the second inequality. Thus, by Bernstein's inequality, we have

$$\begin{aligned} \Pr \left[ |\mathcal{L}_1| \geq \frac{R\iota}{d^3} \sqrt{\frac{C'\iota}{n}} \right] &\leq 2 \exp \left( \frac{-\frac{\iota^3 C' r^2 \kappa_1}{2nd^6}}{\frac{C' r^2 \kappa_1 \iota^2}{nd^6} + \frac{RC_2 \|\mathbf{A}_k\|_{\text{op}}}{3nd} \sqrt{\frac{r^2 \kappa_1 \iota}{nd^6}}} \right) \\ &= 2 \exp \left( \frac{-\frac{\iota}{2}}{1 + \frac{C_2 \|\mathbf{A}_k\|_{\text{op}}}{3C'} \sqrt{\frac{d^4}{n\iota^3}}} \right). \end{aligned}$$

Thus, when  $n \geq C_2^2 \|\mathbf{A}_k\|_{\text{op}}^2 d^4$ , we have with high probability on the training dataset  $\mathcal{D}_1$ ,

$$|\mathcal{L}_1| \leq \frac{R\iota}{d^3} \sqrt{\frac{C'\iota}{n}} \lesssim \frac{R\iota^{3/2}}{\sqrt{nd^3}}.$$

**Bounding  $\mathcal{L}_2$ .** It suffices to bound

$$\begin{aligned} \left| \mathbb{E}_{\mathbf{x}} [\tilde{Y}_k(\mathbf{x})] - \mathbb{E}_{\mathbf{x}} [Y_k(\mathbf{x})] \right| &\leq \mathbb{E}_{\mathbf{x}} \left[ |f^*(\mathbf{x})| \mathbf{1}\{f^*(\mathbf{x}) > R\} \left| K^{(0)}(\mathbf{x}, \mathbf{x}') \mathbf{x}^\top \mathbf{A}_k \mathbf{x} \right| \right] \\ &\leq \mathbb{E}_{\mathbf{x}} [(f^*(\mathbf{x}))^2]^{\frac{1}{2}} \Pr [f^*(\mathbf{x}) > R]^{\frac{1}{4}} \mathbb{E}_{\mathbf{x}} \left[ K^{(0)}(\mathbf{x}, \mathbf{x}')^8 \right]^{\frac{1}{8}} \mathbb{E}_{\mathbf{x}} [(\mathbf{x}^\top \mathbf{A}_k \mathbf{x})^8]^{\frac{1}{8}} \\ &\leq 1 \cdot \exp(-\eta/2) \cdot \frac{C_2}{d^2} \cdot (8-1) \\ &= \frac{7C_2 \exp(-\eta/2)}{d^2}. \end{aligned}$$

Here we invoke (33) and Lemma 8 in the last inequality. By taking  $\eta = \iota \geq 2 \log n + 8 \log d$ , we can ensure that with high probability, we have

$$\left| \frac{1}{n} \sum_{i=1}^n Y_k(\mathbf{x}_i) - \mathbb{E}_{\mathbf{x}} [Y_k(\mathbf{x})] \right| \leq |\mathcal{L}_1| + |\mathcal{L}_2| \lesssim \frac{R\iota^{3/2}}{\sqrt{nd^3}}.$$

Thus, by taking  $k$  over  $[r]$ , we have with high probability over the training set  $\mathcal{D}_1$ , we have

$$\|\mathbf{D}_2\|_\infty \leq \frac{|c_2|}{B(d, 2)} \cdot \frac{R\iota^{3/2}}{\sqrt{nd^3}} \lesssim \frac{R\iota^{3/2}}{\sqrt{nd^5}} \lesssim \frac{\iota^{p+3/2}}{\sqrt{nd^5}}.$$

The proof is complete. □

*Proof of Lemma 23.* Note that for any  $k \in [r]$ , we have

$$\begin{aligned} &\mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x}) K^{(0)}(\mathbf{x}, \mathbf{x}') \mathbf{x}^\top \mathbf{A}_k \mathbf{x}] \\ &= \mathbb{E}_{\mathbf{x}} \left[ f^*(\mathbf{x}) \sum_{i=2}^{\infty} \frac{c_i^2 Q_i(\mathbf{x}^\top \mathbf{x}')}{B(d, i)} \cdot \mathbf{x}^\top \mathbf{A}_k \mathbf{x} \right] \\ &= \sum_{i=2}^{\infty} \frac{c_i^2}{B(d, i)} \cdot \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x}) Q_i(\mathbf{x}^\top \mathbf{x}') \mathbf{x}^\top \mathbf{A}_k \mathbf{x}] \\ &= \frac{c_2^2}{B(d, 2)d(d-1)} \cdot \langle \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x})(\mathbf{x}^\top \mathbf{A}_k \mathbf{x})(\mathbf{x}\mathbf{x}'^\top - \mathbf{I})], \mathbf{x}'\mathbf{x}'^\top - \mathbf{I} \rangle \\ &\quad + \sum_{i=3}^{\infty} \frac{c_i^2}{B(d, i)} \cdot \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x}) Q_i(\mathbf{x}^\top \mathbf{x}') \mathbf{x}^\top \mathbf{A}_k \mathbf{x}] \\ &= \frac{c_2^2}{B(d, 2)d(d-1)} \cdot \langle T(\mathbf{A}_k), \mathbf{x}'\mathbf{x}'^\top - \mathbf{I} \rangle + \sum_{i=3}^{\infty} \frac{c_i^2}{B(d, i)} \cdot \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x}) Q_i(\mathbf{x}^\top \mathbf{x}') \mathbf{x}^\top \mathbf{A}_k \mathbf{x}]. \end{aligned}$$

2160 Here  $T$  is the linear operator defined in (14). Recall by Proposition 3, we have  
 2161

$$2162 \quad \left\| T(\mathbf{A}_k) - \sum_{j=1}^r \mathbf{H}_{k,j} \mathbf{A}_j \right\|_F \lesssim d^{-1/6} L r^2 \kappa_1 \log^2 d.$$

$$2163$$

$$2164$$

2165 Let's denote  $\mathbf{R}_k = T(\mathbf{A}_k) - \sum_{j=1}^r \mathbf{H}_{k,j} \mathbf{A}_j$  so that  $\|\mathbf{R}_k\|_F \lesssim d^{-1/6} L r^2 \kappa_1 \log^2 d$ . Since  
 2166  $\langle \mathbf{R}_k, \mathbf{x}' \mathbf{x}'^\top - \mathbf{I} \rangle$  is a quadratic function of  $\mathbf{x}'$ , and  $\mathbb{E}_{\mathbf{x}'} [\langle \mathbf{R}_k, \mathbf{x}' \mathbf{x}'^\top - \mathbf{I} \rangle^2] = \frac{2d}{d+2} \|\mathbf{R}_k\|_F^2$ . By  
 2167 Lemma 8, there exists a constant  $C > 0$  such that  
 2168

$$2169 \quad \Pr \left[ |\langle \mathbf{R}_k, \mathbf{x}' \mathbf{x}'^\top - \mathbf{I} \rangle| \geq C \iota \sqrt{\mathbb{E}_{\mathbf{x}'} [\langle \mathbf{R}_k, \mathbf{x}' \mathbf{x}'^\top - \mathbf{I} \rangle^2]} \right] \leq 2 \exp(-\iota).$$

$$2170$$

2171 Thus, by enumerating  $k \in [r]$  and  $\mathbf{x}' \in \mathcal{D}_2$ , we obtain that with high probability  $(1 - nr \exp(-\iota))$   
 2172 on  $\mathcal{D}_2$ , for any  $k \in [r]$ , we have  
 2173

$$2174 \quad \left| \langle T(\mathbf{A}_k), \mathbf{x}' \mathbf{x}'^\top - \mathbf{I} \rangle - \sum_{j=1}^r H_{k,j} \mathbf{x}'^\top \mathbf{A}_j \mathbf{x}' \right| \lesssim \frac{\iota L r^2 \kappa_1 \log^2 d}{d^{1/6}}.$$

$$2175$$

$$2176$$

2177 Moreover, we have for any  $\mathbf{x}'$   
 2178

$$2179 \quad \left| \sum_{i=3}^{\infty} \frac{c_i^2}{B(d, i)} \cdot \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x}) Q_i(\mathbf{x}^\top \mathbf{x}') \mathbf{x}^\top \mathbf{A}_k \mathbf{x}] \right|$$

$$2180$$

$$2181 \leq \sum_{i=3}^{\infty} \frac{c_i^2}{B(d, i)} \cdot |\mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x}) Q_i(\mathbf{x}^\top \mathbf{x}') \mathbf{x}^\top \mathbf{A}_k \mathbf{x}]|$$

$$2182$$

$$2183 \leq \sum_{i=3}^{\infty} \frac{c_i^2}{B(d, i)} \cdot \sqrt{\mathbb{E}_{\mathbf{x}} [Q_i(\mathbf{x}^\top \mathbf{x}')^2] \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x})^2 (\mathbf{x}^\top \mathbf{A}_k \mathbf{x})^2]}$$

$$2184$$

$$2185 \leq \sum_{i=3}^{\infty} \frac{c_i^2}{B(d, i)^{3/2}} \cdot \sqrt{\mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x})^2 (\mathbf{x}^\top \mathbf{A}_k \mathbf{x})^2]}.$$

$$2186$$

$$2187$$

$$2188$$

$$2189$$

2190 Again by Lemma 8, we know that there exists a sufficient large constant  $C > 0$  s.t.  
 2191  $\Pr [|\mathbf{x}^\top \mathbf{A}_k \mathbf{x}| \geq C \iota] \leq 2 \exp(-\iota)$ . By taking  $\iota \geq (2p + 2) \log d$ , we have  
 2192

$$2193 \quad \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x})^2 (\mathbf{x}^\top \mathbf{A}_k \mathbf{x})^2] = \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x})^2 (\mathbf{x}^\top \mathbf{A}_k \mathbf{x})^2 \mathbf{1}\{|\mathbf{x}^\top \mathbf{A}_k \mathbf{x}| \leq C \iota\}]$$

$$2194 \quad + \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x})^2 (\mathbf{x}^\top \mathbf{A}_k \mathbf{x})^2 \mathbf{1}\{|\mathbf{x}^\top \mathbf{A}_k \mathbf{x}| > C \iota\}]$$

$$2195 \lesssim C^2 \iota^2 + 2d^{2p+2} \exp(-\iota)$$

$$2196 \lesssim C^2 \iota^2.$$

$$2197$$

$$2198$$

$$2199$$

2200 Altogether, with high probability on  $\mathcal{D}_2$ , for any  $k \in [r]$ , we have  
 2201

$$2202 \quad \left| \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x}) K^{(0)}(\mathbf{x}, \mathbf{x}') \mathbf{x}^\top \mathbf{A}_k \mathbf{x}] - \frac{c_2^2}{B(d, 2)d(d-1)} \cdot \sum_{j=1}^r H_{k,j} \mathbf{x}'^\top \mathbf{A}_j \mathbf{x}' \right|$$

$$2203$$

$$2204 \leq \frac{\iota L r^2 \kappa_1 \log^2 d}{B(d, 2)d^{7/6}(d-1)} + \sum_{i=3}^{\infty} \frac{C \iota c_i^2}{B(d, i)^{3/2}}.$$

$$2205$$

$$2206$$

2207 Thus, by paralleling the  $r$  entries together, we have with high probability on  $\mathcal{D}_2$   
 2208

$$2209 \quad \left\| \mathbf{D}_3 - \frac{c_2^2}{B(d, 2)^2 d(d-1)} \cdot \mathbf{H} \mathbf{p}(\mathbf{x}') \right\|_\infty \leq \frac{\iota L r^2 \kappa_1 \log^2 d}{B(d, 2)^2 d^{7/6}(d-1)} + \sum_{i=3}^{\infty} \frac{C \iota c_i^2}{B(d, 2)B(d, i)^{3/2}}$$

$$2210$$

$$2211 \lesssim \frac{\iota L r^2 \kappa_1 \log^2 d}{d^{6+1/6}}.$$

$$2212$$

$$2213$$

2214 The proof is complete. □

2214 *Proof of Lemma 3.* By the mean value theorem, we have  
 2215

$$2216 \quad |g(\mathbf{B}^* \mathbf{h}^{(1)}(\mathbf{x})) - g(\mathbf{p}(\mathbf{x}))| \lesssim \sup_{\lambda \in [0,1]} \left\| \nabla g(\lambda \mathbf{B}^* \mathbf{h}^{(1)}(\mathbf{x}) + (1-\lambda) \mathbf{p}(\mathbf{x})) \right\|_2 \left\| \mathbf{B}^* \mathbf{h}^{(1)}(\mathbf{x}) - \mathbf{p}(\mathbf{x}) \right\|_2.$$

2218 Recall by (8), we have  $\|\nabla g(\mathbf{z})\|_2 \lesssim \|g\|_{L^2} \sum_{k=1}^p r^{\frac{p-k}{4}} \|\mathbf{z}\|_2^{k-1}$ . Note that with high probability,  
 2219  $\sup_{\mathbf{x} \in \mathcal{D}_2} \|\mathbf{p}(\mathbf{x})\| \leq \tilde{O}(\sqrt{r})$ . Therefore  
 2220

$$2221 \quad \sup_{\mathbf{x} \in \mathcal{D}_2} \sup_{\lambda \in [0,1]} \left\| \nabla g(\lambda \mathbf{B}^* \mathbf{h}^{(1)}(\mathbf{x}) + (1-\lambda) \mathbf{p}(\mathbf{x})) \right\| \lesssim \|g\|_{L^2} r^{\frac{p-1}{2}}.$$

2223 Altogether, by Proposition 5,  
 2224

$$\begin{aligned} 2225 \quad & \sup_{\mathbf{x} \in \mathcal{D}_2} |g(\mathbf{B}^* \mathbf{h}^{(1)}(\mathbf{x})) - g(\mathbf{p}(\mathbf{x}))| \\ 2226 \quad & \lesssim \|g\|_{L^2} r^{\frac{p-1}{2}} \left\| \mathbf{B}^* \mathbf{h}^{(1)}(\mathbf{x}) - \mathbf{p}(\mathbf{x}) \right\|_2 \\ 2227 \quad & \leq \|g\|_{L^2} \cdot \frac{r^{p/2}}{\lambda_{\min}(\mathbf{H})} \cdot \left( \frac{\iota^{p+2} d^5}{m_2} + \frac{\iota d^3}{\sqrt{m_2}} + \frac{\iota^{p+3/2} d}{\sqrt{n}} + \frac{\iota L r^2 \kappa_1 \log^2 d}{d^{1/6}} \right) \end{aligned}$$

2228 The proof is complete.  $\square$   
 2229

#### 2230 C.4 PROOF OF OTHER SUPPORTING LEMMAS

2231 We first present the concentration of the initial kernel  $K_{m_2}^{(0)}(\mathbf{x}, \mathbf{x}')$ .  
 2232

2233 **Lemma 24.** Let  $K_{m_2}^{(0)}(\mathbf{x}, \mathbf{x}') = \frac{1}{m_2} \langle \sigma_2(\mathbf{V}\mathbf{x}), \sigma_2(\mathbf{V}\mathbf{x}') \rangle$  be the initial kernel with inner width being  
 2234  $m_2$ , and  $K^{(0)}(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\mathbf{v} \sim \text{Unif-S}^{d-1}(\sqrt{d})} [\sigma_2(\mathbf{v}^\top \mathbf{x}) \sigma_2(\mathbf{v}^\top \mathbf{x}')]^2$  be the infinite-width kernel. Then  
 2235 there exists a constant  $C$  s.t. when  $m_2 \geq Cd^4$ , with high probability probability on  $\mathbf{w}$ ,  $\mathbf{V}$  and the  
 2236 training dataset  $\mathcal{D}$ , for any  $\mathbf{x} \in \mathcal{D}_1$  and  $\mathbf{x}' \in \mathcal{D}_2$ , we have  
 2237

$$2238 \quad \left| K_{m_2}^{(0)}(\mathbf{x}, \mathbf{x}') - K^{(0)}(\mathbf{x}, \mathbf{x}') \right| \leq \frac{\iota}{\sqrt{m_2} d^2}.$$

2239 *Proof of Lemma 24.* By Assumption 4, for any  $\mathbf{x}, \mathbf{x}' \in \mathcal{D}$  and  $\mathbf{v} \in \mathbb{S}^{d-1}(\sqrt{d})$ , we have  
 2240

$$2241 \quad |\sigma_2(\mathbf{v}^\top \mathbf{x}) \sigma_2(\mathbf{v}^\top \mathbf{x}')| \leq C_\sigma^2$$

2242 and  
 2243

$$2244 \quad \mathbb{E}_{\mathbf{v}} [\sigma_2(\mathbf{v}^\top \mathbf{x})^2 \sigma_2(\mathbf{v}^\top \mathbf{x}')^2] \leq \sqrt{\mathbb{E}_{\mathbf{v}} [\sigma_2(\mathbf{v}^\top \mathbf{x})^4] \mathbb{E}_{\mathbf{v}} [\sigma_2(\mathbf{v}^\top \mathbf{x}')^4]} \leq \frac{C_4}{d^4}.$$

2245 Thus, by Bernstein inequality, we have  
 2246

$$\begin{aligned} 2247 \quad \Pr \left[ \left| K_{m_2}^{(0)}(\mathbf{x}, \mathbf{x}') - K^{(0)}(\mathbf{x}, \mathbf{x}') \right| \geq \sqrt{\frac{t}{m_2}} \right] & \leq 2 \exp \left( \frac{-\frac{t}{2m_2}}{\frac{C_4}{m_2 d^4} + \frac{C_\sigma^2}{3m_2} \sqrt{\frac{t}{m_2}}} \right) \\ 2248 \quad & = \exp \left( \frac{-t/2}{\frac{C_4}{d^4} + \frac{C_\sigma^2}{3} \sqrt{\frac{t}{m_2}}} \right). \end{aligned} \tag{34}$$

2249 By enumerating  $\mathbf{x}, \mathbf{x}'$  over  $\mathcal{D}$ , we have  
 2250

$$2251 \quad \Pr \left[ \max_{\mathbf{x}, \mathbf{x}' \in \mathcal{D}} \left| K_{m_2}^{(0)}(\mathbf{x}, \mathbf{x}') - K^{(0)}(\mathbf{x}, \mathbf{x}') \right| \geq \sqrt{\frac{t}{m_2}} \right] \leq n^2 \exp \left( \frac{-t/2}{\frac{C_4}{d^4} + \frac{C_\sigma^2}{3} \sqrt{\frac{t}{m_2}}} \right).$$

2252 Thus, when  $m_2 \geq d^4$ , we can take  $t = \iota^2/d^4$  to bound the probability by  $\text{poly}(d, n, m_2) e^{-\iota}$ , which  
 2253 concludes our proof.  $\square$   
 2254

2255 Then we present the concentration of the reconstructed features.  
 2256

2268   **Lemma 25.** Suppose  $m_2 \geq C_\sigma^2 C_4^{-1/2} d^4 \|\mathbf{A}\|_{\text{op}}^2$ . Given any  $\mathbf{A}$  such that  $\mathbf{v}^\top \mathbf{A} \mathbf{v}$  is a quadratic  
 2269 spherical harmonic, with high probability on  $\mathbf{V}$ , for any  $\mathbf{x} \in \mathcal{D}$ , we have  
 2270

$$2271 \quad \left| \frac{1}{m_2} \sum_{i=1}^{m_2} (\mathbf{v}_i^\top \mathbf{A} \mathbf{v}_i) \sigma_2(\mathbf{v}_i^\top \mathbf{x}) - \frac{c_2}{B(d, 2)} \mathbf{x}^\top \mathbf{A} \mathbf{x} \right| \leq \frac{9\iota d^{-1} C_4^{1/4}}{\sqrt{m_2}}.$$

2274   Proof of Lemma 25. Given any fixed  $\mathbf{x} \in \mathcal{D}$  and  $\mathbf{A}$  such that  $\mathbf{v}^\top \mathbf{A} \mathbf{v}$  is a quadratic spherical har-  
 2275 monic, we have  
 2276

$$2277 \quad \mathbb{E}_{\mathbf{v}} [(\mathbf{v}^\top \mathbf{A} \mathbf{v})^2 \sigma_2^2(\mathbf{v}^\top \mathbf{x})] \leq \sqrt{\mathbb{E}_{\mathbf{v}} [(\mathbf{v}^\top \mathbf{A} \mathbf{v})^4] \mathbb{E}_{\mathbf{v}} [\sigma_2^4(\mathbf{v}^\top \mathbf{x})]} \\ 2278 \quad \leq (4-1)^{2*2} \mathbb{E} [(\mathbf{v}_i^\top \mathbf{A} \mathbf{v}_i)^2] d^{-2} C_4^{1/2} \\ 2279 \quad = 81d^{-2} C_4^{1/2}$$

2280   and  
 2281

$$2282 \quad |(\mathbf{v}^\top \mathbf{A} \mathbf{v}) \sigma_2(\mathbf{v}^\top \mathbf{x})| \leq d \|\mathbf{A}\|_{\text{op}} \cdot C_\sigma = dC_\sigma \|\mathbf{A}\|_{\text{op}}.$$

2283   Since  $\mathbb{E}_{\mathbf{v}} [(\mathbf{v}^\top \mathbf{A} \mathbf{v}) \sigma_2(\mathbf{v}^\top \mathbf{x})] = \frac{c_2}{B(d, 2)} \mathbf{x}^\top \mathbf{A} \mathbf{x}$ , by Bernstein Inequality, we have  
 2284

$$2285 \quad \Pr \left[ \left| \frac{1}{m_2} \sum_{i=1}^{m_2} (\mathbf{v}_i^\top \mathbf{A} \mathbf{v}_i) \sigma_2(\mathbf{v}_i^\top \mathbf{x}) - \frac{c_2}{B(d, 2)} \mathbf{x}^\top \mathbf{A} \mathbf{x} \right| \geq \frac{9\iota d^{-1} C_4^{1/4}}{\sqrt{m_2}} \right] \\ 2286 \quad \leq 2 \exp \left( - \frac{\frac{81d^{-2} C_4^{1/2} \iota^2}{2m_2}}{\frac{81d^{-2} C_4^{1/2}}{m_2} + \frac{1}{3m_2} \cdot d(C_\sigma \|\mathbf{A}\|_{\text{op}} \cdot \frac{9\iota d^{-1} C_4^{1/4}}{\sqrt{m_2}})} \right) \\ 2287 \quad = 2 \exp \left( - \frac{\iota^2 / 2}{1 + \frac{C_\sigma d^2 \|\mathbf{A}\|_{\text{op}}}{27C_4^{1/4} \sqrt{m_2}} \cdot \iota} \right)$$

2288   Thus, when  $m_2 \geq C_\sigma^2 C_4^{-1/2} d^4 \|\mathbf{A}\|_{\text{op}}^2$ , by enumerating  $\mathbf{x} \in \mathcal{D}$ , we obtain that with high probability  
 2289 on  $\mathbf{V}$ , for any  $\mathbf{x} \in \mathcal{D}$ , we have  
 2290

$$2291 \quad \left| \frac{1}{m_2} \sum_{i=1}^{m_2} (\mathbf{v}_i^\top \mathbf{A} \mathbf{v}_i) \sigma_2(\mathbf{v}_i^\top \mathbf{x}) - \frac{c_2}{B(d, 2)} \mathbf{x}^\top \mathbf{A} \mathbf{x} \right| \leq \frac{9\iota d^{-1} C_4^{1/4}}{\sqrt{m_2}}.$$

2292   The proof is complete. □  
 2293

## D APPROXIMATION THEORY OF THE OUTER LAYER

### D.1 PROOF OF PROPOSITION 2

2304   Since we mainly focus on the first training stage throughout this section, we may sometimes denote  
 2305  $n = n_1$  for notation simplicity, and let the training set be  $\mathcal{D}_1 = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ . Let's consider a  
 2306 formal version of Proposition 2.  
 2307

2308   **Proposition 6.** Suppose  $g$  is a degree  $p$  polynomial. By setting  $\eta = C\iota^{-5} \kappa_2^{-1} m_2^{-1/2} d^6$  for some  
 2309 constant  $C > 0$ , with high probability over  $\mathcal{D}_1, \mathcal{D}_2, \{\mathbf{w}_i\}_{i=1}^{m_1}$  and  $\mathbf{V}$ , there exists  $\mathbf{a}^* \in \mathbb{R}^{m_1}$  such  
 2310 that the parameter  $\theta^* = (\mathbf{a}^*, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{V})$  gives rise to  
 2311

$$2312 \quad \begin{aligned} \mathcal{L}_2(\theta^*) &:= \frac{1}{n_2} \sum_{\mathbf{x} \in \mathcal{D}_2} (f(\mathbf{x}; \theta^*) - g(\mathbf{p}(\mathbf{x})))^2 \\ 2313 &\lesssim \|g\|_{L^2}^2 \cdot \frac{r^p}{\lambda_{\min}(\mathbf{H})} \cdot \left( \frac{\iota^{p+2} d^5}{m_2} + \frac{\iota d^3}{\sqrt{m_2}} + \frac{\iota^{p+3/2} d}{\sqrt{n}} + \frac{\iota L r^2 \log^2 d}{d^{1/6}} \right)^2 \\ 2314 &\quad + \frac{\iota^{p+1} \|g\|_{L^2}^2}{m_1} \cdot \left( \sum_{k=0}^p \eta^{-k} \|\mathbf{B}^*\|_{\text{op}}^k r^{\frac{p-k}{4}} \right)^2. \end{aligned}$$

2322 Here  $\mathbf{a}^*$  satisfies  
 2323

$$2324 \frac{\|\mathbf{a}^*\|_2^2}{m_1} \lesssim \iota^p \|g\|_{L^2}^2 \cdot \left( \sum_{k=0}^p \eta^{-k} \|\mathbf{B}^*\|_{\text{op}}^k r^{\frac{p-k}{4}} \right)^2 = \tilde{\Omega}(\kappa_2^{2p} r^p).$$

2327 To prove the proposition, let's introduce the infinite-outer-width model as a transition term between  
 2328 the finite-outer-width model and the target function. We define the infinite-outer-width model as  
 2329

$$2330 f_{\infty, m_2}(\mathbf{x}; v) = \mathbb{E}_{a, b, \mathbf{w}} \left[ v(a, b, \mathbf{w}) \sigma_1 \left( a \eta \langle \mathbf{w}, \mathbf{h}^{(1)}(\mathbf{x}) \rangle + b \right) \right],$$

2332 where  $\mathbf{h}^{(1)}(\mathbf{x}') = \frac{1}{n} \sum_{i=1}^n f^*(\mathbf{x}_i) \cdot K_{m_2}^{(0)}(\mathbf{x}_i, \mathbf{x}') \cdot \sigma_2(\mathbf{V}^\top \mathbf{x}_i)$ .  
 2333

2334 We can decompose the  $L^2$  loss of the truth model  $f(\mathbf{x}; \theta)$  as  
 2335

$$\begin{aligned} 2336 \hat{\mathcal{L}}(\theta^*) &= \frac{1}{n} \sum_{\mathbf{x} \in \mathcal{D}_2} (f(\mathbf{x}; \theta) - f^*(\mathbf{x}))^2 \\ 2337 &= \frac{1}{n} \sum_{\mathbf{x} \in \mathcal{D}_2} \left( f(\mathbf{x}; \theta) - f_{\infty, m_2}(\mathbf{x}') + f_{\infty, m_2}(\mathbf{x}') - g(\mathbf{B}^* \mathbf{h}^{(1)}(\mathbf{x})) + g(\mathbf{B}^* \mathbf{h}^{(1)}(\mathbf{x})) - g(\mathbf{p}(\mathbf{x})) \right)^2 \\ 2338 &\lesssim \underbrace{\frac{1}{n} \sum_{\mathbf{x} \in \mathcal{D}_2} (f(\mathbf{x}; \theta) - f_{\infty, m_2}(\mathbf{x}'))^2}_{L_1} \\ 2339 &\quad + \underbrace{\frac{1}{n} \sum_{\mathbf{x} \in \mathcal{D}_2} \left( f_{\infty, m_2}(\mathbf{x}') - g(\mathbf{B}^* \mathbf{h}^{(1)}(\mathbf{x})) \right)^2}_{L_2} \\ 2340 &\quad + \underbrace{\frac{1}{n} \sum_{\mathbf{x} \in \mathcal{D}_2} \left( g(\mathbf{B}^* \mathbf{h}^{(1)}(\mathbf{x})) - g(\mathbf{p}(\mathbf{x})) \right)^2}_{L_3}. \\ 2341 & \\ 2342 & \\ 2343 & \\ 2344 & \\ 2345 & \\ 2346 & \\ 2347 & \\ 2348 & \\ 2349 & \\ 2350 & \\ 2351 & \\ 2352 & \end{aligned}$$

2353 We have bounded  $L_3$  in Corollary 3. We state Lemmas 26 and 27 as follows to bound  $L_1$  and  $L_2$ ,  
 2354 respectively.

2355 **Lemma 26** (Bound  $L_2$ ). *Given  $\mathbf{B}^* \in \mathbb{R}^{r \times m_2}$  and setting the learning rate  $\eta = C \iota^{-5} \kappa_2^{-1} m_2^{-1/2} d^6$   
 2356 for a constant  $C > 0$ , there exists  $v : \{\pm 1\} \times \mathbb{R} \times \mathbb{R}^{m_2} \rightarrow \mathbb{R}$  such that*

$$2358 \|\mathbf{v}\|_{L^2} \lesssim \|g\|_{L^2} \sum_{k=0}^p \eta^{-k} \|\mathbf{B}^*\|_{\text{op}}^k r^{\frac{p-k}{4}},$$

2361 and, with high probability over  $\mathcal{D}_1$ ,  $\mathcal{D}_2$  and  $\mathbf{V}$ , the infinite-width network satisfies

$$2363 \frac{1}{n} \sum_{\mathbf{x} \in \mathcal{D}_2} (f_{\infty, m_2}(\mathbf{x}; v) - g(\mathbf{B}^* \mathbf{h}^{(1)}(\mathbf{x})))^2 \lesssim o\left(\frac{1}{d^2 n_1^2 n_2^2 m_1^2 m_2^2}\right).$$

2365 **Lemma 27** (Bound  $L_1$ ). *Given the function  $v : \{\pm 1\} \times \mathbb{R} \times \mathbb{R}^{m_2} \rightarrow \mathbb{R}$  in Lemma 26. With high  
 2366 probability over  $\mathcal{D}_1$ ,  $\mathcal{D}_2$ ,  $\{\mathbf{w}_i\}_{i=1}^{m_1}$  and  $\mathbf{V}$ , it holds that for any  $\mathbf{x} \in \mathcal{D}_2$ ,*

$$\begin{aligned} 2368 &\left| \frac{1}{m_1} \sum_{i=1}^{m_1} v(a_i, b_i, \mathbf{w}_i) \sigma_1(\eta a_i \langle \mathbf{w}_i, \mathbf{h}^{(1)}(\mathbf{x}) \rangle + b_i) - f_{\infty, m_2}(\mathbf{x}; v) \right| \lesssim \sqrt{\frac{\iota^{p+1} \|v\|_{L^2}^2}{m_1}}, \text{ with} \\ 2369 & \\ 2370 & \\ 2371 & \\ 2372 & \\ 2373 & \\ 2374 & \\ 2375 & \end{aligned}$$

The proof of Lemmas 26 and 27 is provided in Appendix D.2. Now we begin our proof of Proposition 6.

Proof of Proposition 6. By Corollary 3, Lemma 27 and Lemma 26, by defining the vector  $\mathbf{a}^* \in \mathbb{R}^{m_1}$  by  $a_i^* = v(a_i^{(0)}, b_i^{(1)}, \mathbf{w}_i^{(1)})$  and letting  $\theta^* = (\mathbf{a}^*, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{V})$ , we have with high probability that

$$\begin{aligned} \hat{\mathcal{L}}_2(\theta^*) &\lesssim L_1 + L_2 + L_3 \\ &\lesssim \|g\|_{L^2}^2 \cdot \frac{r^p}{\lambda_{\min}^2(\mathbf{H})} \cdot \left( \frac{\iota^{p+2} d^5}{m_2} + \frac{\iota d^3}{\sqrt{m_2}} + \frac{\iota^{p+3/2} d}{\sqrt{n}} + \frac{\iota L r^2 \log^2 d}{d^{1/6}} \right)^2 \\ &\quad + \frac{\iota^{p+1} \|g\|_{L^2}^2}{m_1} \cdot \left( \sum_{k=0}^p \eta^{-k} \|\mathbf{B}^*\|_{\text{op}}^k r^{\frac{p-k}{4}} \right)^2 \\ &\quad + o\left(\frac{1}{d^2 n_1^2 n_2^2 m_1^2 m_2^2}\right) \\ &\lesssim \|g\|_{L^2}^2 \cdot \frac{r^p}{\lambda_{\min}^2(\mathbf{H})} \cdot \left( \frac{\iota^{p+2} d^5}{m_2} + \frac{\iota d^3}{\sqrt{m_2}} + \frac{\iota^{p+3/2} d}{\sqrt{n}} + \frac{\iota L r^2 \log^2 d}{d^{1/6}} \right)^2 \\ &\quad + \frac{\iota^{p+1} \|g\|_{L^2}^2}{m_1} \cdot \left( \sum_{k=0}^p \eta^{-k} \|\mathbf{B}^*\|_{\text{op}}^k r^{\frac{p-k}{4}} \right)^2. \end{aligned}$$

Here  $\mathbf{a}^*$  satisfies

$$\begin{aligned} \|\mathbf{a}^*\|_2^2 &\leq \sum_{i=1}^{m_1} v(a_i, b_i, \mathbf{w}_i)^2 \\ &\lesssim m_1 \iota^p \|v\|_{L^2}^2 \\ &\lesssim m_1 \iota^p \|g\|_{L^2}^2 \left( \sum_{k=0}^p \eta^{-k} \|\mathbf{B}^*\|_{\text{op}}^k r^{\frac{p-k}{4}} \right)^2. \end{aligned}$$

The proof is complete.  $\square$

## D.2 OMITTED PROOFS IN APPENDIX D.1

### D.2.1 RANDOM FEATURE CONSTRUCTION OF UNIVARIATE POLYNOMIALS

In this section, before proving Lemmas 26 and 27, we first construct univariate polynomials using the outer activation function  $\sigma_1$  and the random features  $a$  and  $b$  progressively.

**Lemma 28.** *There exists  $v_0(a, b)$ , supported on  $\{\pm 1\} \times [2, 3]$ , such that for any  $|z| \leq 1$*

$$\mathbb{E}_{a,b}[v_0(a, b)\sigma(az + b)] = 1, \quad \sup_{a,b} |v(a, b)| \lesssim 1.$$

*Proof.* Let  $v_0(a, b) = 12 \cdot \mathbf{1}_{a=1}(b - \frac{5}{2}) \cdot \frac{\mathbf{1}_{b \in [2, 3]}}{\mu(b)}$ . Then, since  $z + b \geq 1$ ,

$$\begin{aligned} \mathbb{E}_{a,b}[v_0(a, b)\sigma(az + b)] &= 6 \int_2^3 (b - \frac{5}{2})\sigma(z + b)db \\ &= 6 \int_2^3 (b - \frac{5}{2})(2z + 2b - 1)db \\ &= z \cdot 6 \int_2^3 (b - \frac{5}{2})db + 6 \int_2^3 (b - \frac{5}{2})(2b - 1)db \\ &= 1. \end{aligned}$$

The proof is complete.  $\square$

**Lemma 29.** *There exists  $v_1(a, b)$ , supported on  $\{\pm 1\} \times [2, 3]$ , such that for any  $|z| \leq 1$*

$$\mathbb{E}_{a,b}[v_1(a, b)\sigma(az + b)] = z, \quad \sup_{a,b} |v(a, b)| \lesssim 1.$$

2430 *Proof.* Let  $v_0(a, b) = \mathbf{1}_{a=1}(-24b + 61) \cdot \frac{\mathbf{1}_{b \in [2, 3]}}{\mu(b)}$ . Then, since  $z + b \geq 1$ ,

$$\begin{aligned} 2432 \quad \mathbb{E}_{a,b}[v_1(a, b)\sigma(az + b)] &= \frac{1}{2} \int_2^3 (-24b + 61)\sigma(z + b)db \\ 2433 \quad &= \frac{1}{2} \int_2^3 (-24b + 61)(2z + 2b - 1)db \\ 2434 \quad &= z \int_2^3 (-24b + 61)db + \frac{1}{2} \int_2^3 (-24b + 61)(2b - 1)db \\ 2435 \quad &= z. \end{aligned}$$

2440 The proof is complete.  $\square$

2442 **Lemma 30.** *There exists  $v_2(a, b)$ , supported on  $\{\pm 1\} \times [-2, 3]$ , such that for any  $|z| \leq 1$*

$$2443 \quad \mathbb{E}_{a,b}[v_2(a, b)\sigma(az + b)] = z^2, \quad \sup_{a,b} |v(a, b)| \lesssim 1.$$

2446 *Proof.* First, see that

$$\begin{aligned} 2447 \quad \int_{-2}^2 \sigma(z + b)db &= \int_{-2+z}^{2+z} \sigma(b)db \\ 2448 \quad &= \int_{-2+z}^{-1} (-2b - 1)db + \int_{-1}^1 b^2 db + \int_1^{2+z} (2b - 1)db \\ 2449 \quad &= [-b^2 - b]_{-2+z}^{-1} + \frac{2}{3} + [b^2 - b]_1^{2+z} \\ 2450 \quad &= (z - 2)^2 + (z - 2) + \frac{2}{3} + (z + 2)^2 - (z + 2) \\ 2451 \quad &= 2z^2 + \frac{14}{3}. \end{aligned}$$

2452 Let  $v_2(a, b) = \mathbf{1}_{a=1} \frac{\mathbf{1}_{b \in [-2, 2]}}{\mu(b)} - \frac{7}{3}v_0(a, b)$  Then

$$\begin{aligned} 2453 \quad \mathbb{E}_{a,b}[v_2(a, b)\sigma(az + b)] &= \frac{1}{2} \int_{-2}^2 \sigma(z + b)db - \frac{7}{3} \\ 2454 \quad &= z^2 + \frac{7}{3} - \frac{7}{3} \\ 2455 \quad &= z^2. \end{aligned}$$

2456 The proof is complete.  $\square$

2457 **Lemma 31.** *Let  $v(b) = -\frac{1}{2}k(k-1)(k-2)(1-b)^{k-3} \cdot \frac{\mathbf{1}_{b \in [0, 1]}}{\mu(b)}$ . Then*

$$2458 \quad \mathbb{E}_b[v_k(b)\sigma(z + b)] = z^k \cdot \mathbf{1}_{z>0} - \frac{k(k-1)}{2}z^2 - kz - 1.$$

2459 *Proof.* Plugging in  $v_k(b)$  and applying integration by parts yields

$$\begin{aligned} 2460 \quad \mathbb{E}_b[v_k(b)\sigma(z + b)] &= \int_0^1 -\frac{1}{2}k(k-1)(k-2)(1-b)^{k-3}\sigma(z + b)db \\ 2461 \quad &= [\frac{1}{2}k(k-1)(1-b)^{k-2}\sigma(z + b)]_0^1 - \int_0^1 \frac{1}{2}k(k-1)(1-b)^{k-2}\sigma'(z + b)db \\ 2462 \quad &= -\frac{1}{2}k(k-1)\sigma(z) + [\frac{1}{2}k(1-b)^{k-1}\sigma'(z + b)]_0^1 - \int_0^1 \frac{1}{2}k(1-b)^{k-1}\sigma''(z + b)db \\ 2463 \quad &= -\frac{1}{2}k(k-1)\sigma(z) - \frac{1}{2}k\sigma'(z) - \int_0^1 k(1-b)^{k-1}\mathbf{1}_{|z+b| \leq 1}db \end{aligned}$$

2484 When  $1 \geq z > 0$ , we have  
 2485

$$2486 - \int_0^1 k(1-b)^{k-1} \mathbf{1}_{|z+b| \leq 1} db = - \int_0^{1-z} k(1-b)^{k-1} db = [(1-b)^k]_0^{1-z} = z^k - 1.$$

2488 When  $-1 \leq z \leq 0$ , we have  
 2489

$$2490 - \int_0^1 k(1-b)^{k-1} \mathbf{1}_{|z+b| \leq 1} db = - \int_0^1 k(1-b)^{k-1} db = -1.$$

2492 Since  $z \in [-1, 1]$ , we have that  $\sigma(z) = z^2$  and  $\sigma'(z) = 2z$ . Therefore for  $z \in [-1, 1]$   
 2493

$$2494 \mathbb{E}_b[v_k(b)\sigma(z+b)] = z^k \cdot \mathbf{1}_{z>0} - \frac{k(k-1)}{2} z^2 - kz - 1.$$

2496 The proof is complete.  $\square$

2497 **Lemma 32.** *There exists  $v_k(a, b)$ , supported on  $\{\pm 1\} \times [-2, 3]$ , such that for any  $|z| \leq 1$*

$$2499 \mathbb{E}_{a,b}[v_k(a, b)\sigma(az+b)] = z^k, \quad \sup_{a,b} |v_k(a, b)| \lesssim \text{poly}(k).$$

2501 *Proof.* We focus on  $k \geq 3$ . We have that

$$2503 \mathbb{E}_b[v_k(b)\sigma(z+b)] = z^k \cdot \mathbf{1}_{z>0} - \frac{k(k-1)}{2} z^2 - kz - 1.$$

$$2506 \mathbb{E}_b[v_k(b)\sigma(-z+b)] = (-z)^k \cdot \mathbf{1}_{z<0} - \frac{k(k-1)}{2} z^2 + kz - 1.$$

2508 Therefore if  $k$  is even

$$2509 \mathbb{E}_b[v(b)\sigma(z+b) + v(b)\sigma(-z+b)] = z^k - k(k-1)z^2 - 2.$$

2511 Let  $v_k(a, b) = 2v_k(b) + k(k-1)v_2(a, b) + 2$ . Then

$$2512 \mathbb{E}_{a,b}[v_k(a, b)\sigma(az+b)] = \mathbb{E}_b[v_k(b)\sigma(z+b) + v_k(b)\sigma(z-b)] + k(k-1)z^2 + 2 = z^k.$$

2513 If  $k$  is odd,

$$2515 \mathbb{E}_b[v(b)\sigma(z+b) - v(b)\sigma(-z+b)] = z^k - 2kz.$$

2516 Let  $v_k(a, b) = 2av_k(b) + 2kv_1(a, b)$ . Then

$$2518 \mathbb{E}_{a,b}[v_k(a, b)\sigma(az+b)] = \mathbb{E}_b[v_k(b)\sigma(z+b) - v_k(b)\sigma(z-b)] + 2kz = z^k.$$

2519 The proof is complete.  $\square$

## 2521 D.2.2 PROOF OF SUPPORTING LEMMAS IN APPENDIX D.1

2523 *Proof of Lemma 26.* Let's consider a general version of Lemma 26.

2524 **Lemma 33.** *Let  $g : \mathbb{R}^r \rightarrow \mathbb{R}$  be a degree  $p$  polynomial, and let  $\mathbf{B}^\star \in \mathbb{R}^{r \times m_2}$ . Given a set of  
 2525 vectors  $\mathcal{D} = \{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n\} \subseteq \mathbb{R}^{m_2}$  that satisfies  $\eta \langle \mathbf{w}, \mathbf{z} \rangle \leq 1$  for any  $\mathbf{z} \in \mathcal{D}$  with probability at  
 2526 least  $1 - 2(n_1 + n_2) \exp(-\iota^2/2)$  over  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}_{m_2}, \mathbf{I}_{m_2})$  (uniformly over  $\mathcal{D}$ ). Then, there exists  
 2527  $v : \{\pm 1\} \times \mathbb{R} \times \mathbb{R}^m \rightarrow \mathbb{R}$  so that for all  $\mathbf{z} \in \mathcal{D}$ ,*

$$2529 \mathbb{E}_{a,b,\mathbf{w}}[v(a, b, \mathbf{w})\sigma_1(\eta a \langle \mathbf{w}, \mathbf{z} \rangle + b)] = g(\mathbf{B}^\star \mathbf{z}) + o\left(\frac{1}{dn_1 n_2 m_1 m_2}\right), \quad \text{and}$$

$$2531 \|v\|_{L^2} \lesssim \|g\|_{L^2} \sum_{k=0}^p \eta^{-k} \|\mathbf{B}^\star\|_{\text{op}}^k r^{\frac{p-k}{4}}.$$

2534 Thus, according to Proposition 4, we could set the learning rate  $\eta = C\iota^{-5}\kappa_2^{-1}m_2^{-1/2}d^6$  for a con-  
 2535 stant  $C > 0$  to ensure  $|\eta \langle \mathbf{w}, \mathbf{h}^{(1)}(\mathbf{x}') \rangle| \leq 1$  for any  $\mathbf{x}' \in \mathcal{D}_2$  with high probability on  $\mathbf{V}$ ,  $\mathcal{D}_1$ , and  
 2536 probability at least  $1 - 2(n_1 + n_2) \exp(-\iota^2/2)$  on  $\mathbf{w}$ . Thus, taking  $\mathcal{D} = \{\mathbf{h}(\mathbf{x})\}_{\mathbf{x} \in \mathcal{D}_2}$  concludes our  
 2537 proof.  $\square$

To prove Lemma 33, we first decompose  $g$  into sum of polynomials of different degrees and construct a function  $v$  to express these polynomials accordingly.

**Lemma 34.** Given  $\mathbf{z} \in \mathbb{R}^{m_2}$ . Let  $\mathbf{B}^* \in \mathbb{R}^{r \times m_2}$  and  $\mathbf{T}_k \in (\mathbb{R}^r)^{\otimes k}$ . Then, there exists  $v_k : \mathbb{R}^{m_2} \rightarrow \mathbb{R}$  such that

$$\mathbb{E}_{\mathbf{w}} [v_k(\mathbf{w})(\eta \langle \mathbf{w}, \mathbf{z} \rangle)^k] = \mathbf{T}_k((\mathbf{B}^* \mathbf{z})^{\otimes k}).$$

Here  $v_k$  satisfies

$$\|v_k\|_{L^2} \lesssim \eta^{-k} \|\mathbf{B}^*\|_{op}^k \|\mathbf{T}_k\|_F \quad \text{and} \quad \sup_w |v_k(\mathbf{w})| \lesssim m_2^{k/2} \eta^{-k} \|\mathbf{B}^*\|_{op}^k \|\mathbf{T}_k\|_F. \quad (35)$$

*Proof of Lemma 34.* It suffices to solve

$$\mathbb{E}_{\mathbf{w}} [v(\mathbf{w}) \mathbf{w}^{\otimes k}] = \eta^{-k} \mathbf{B}^{*\otimes k}(\mathbf{T}_k),$$

where  $\mathbf{B}^{*\otimes k}(\mathbf{T}_k) \in (\mathbb{R}^{m_2})^{\otimes k}$ . This is achieved by setting

$$v(\mathbf{w}) := \eta^{-k} \text{Vec}(\mathbf{w}^{\otimes k})^T \text{Mat}(\mathbb{E}[\mathbf{w}^{\otimes 2k}])^{-1} \text{Vec}(\mathbf{B}^{*\otimes k}(\mathbf{T}_k)).$$

Then,

$$\|v\|_{L^2}^2 = \eta^{-2k} \text{Vec}(\mathbf{B}^{*\otimes k}(\mathbf{T}_k))^T \text{Mat}(\mathbb{E}[\mathbf{w}^{\otimes 2k}])^{-1} \text{Vec}(\mathbf{B}^{*\otimes k}(\mathbf{T}_k)).$$

Since

$$\text{Mat}(\mathbb{E}[\mathbf{w}^{\otimes 2k}]) \succeq k! \Pi_{\text{Sym}^k(\mathbb{R}^{m_2})},$$

we have

$$\|v\|_{L^2}^2 \lesssim \eta^{-2k} \left\| \mathbf{B}^{*\otimes k}(\mathbf{T}_k) \right\|_F^2 \leq \eta^{-2k} \|\mathbf{B}^*\|_{op}^{2k} \|\mathbf{T}_k\|_F^2.$$

Finally,

$$\begin{aligned} \sup_w |v(\mathbf{w})| &= \eta^{-k} \sup_w \left| \text{Vec}(\mathbf{w}^{\otimes k})^T \text{Mat}(\mathbb{E}[\mathbf{w}^{\otimes 2k}])^{-1} \text{Vec}(\mathbf{B}^{*\otimes k}(\mathbf{T}_k)) \right| \\ &\leq \|\mathbf{w}^{\otimes k}\|_F \left\| \mathbf{B}^{*\otimes k}(\mathbf{T}_k) \right\|_F \\ &\lesssim m_2^{k/2} \eta^{-k} \|\mathbf{B}^*\|_{op}^k \|\mathbf{T}_k\|_F. \end{aligned}$$

The proof is complete.  $\square$

Then we begin our proof of Lemma 33.

*Proof of Lemma 33.* We can write

$$g(\mathbf{z}) = \sum_{k=0}^p \langle \mathbf{T}_k, \mathbf{z}^{\otimes k} \rangle.$$

By Lemma 10, we have  $\|\mathbf{T}_k\|_F \lesssim r^{\frac{p-k}{4}} \|g\|_{L^2}$ .

Define  $v_k(a, b)$  to be the function so that  $\mathbb{E}_{a,b}[v_k(a, b)\sigma_1(az + b)] = z^k$ , and let  $v_k(\mathbf{w})$  be the function where  $\mathbb{E}_{\mathbf{w}}[v(\mathbf{w})(\eta \langle \mathbf{w}, \mathbf{z} \rangle)^k] = \langle \mathbf{T}_k, (\mathbf{B}^* \mathbf{z})^{\otimes k} \rangle$ . Next, define

$$v(a, b, \mathbf{w}) = \sum_{k=0}^p v_k(a, b) v_k(\mathbf{w}).$$

Here  $v_k(a, b)$  is defined in Lemma 34. Then we have that

$$\|v\|_{L^2} \lesssim \sum_{k=0}^p (\mathbb{E}[v_k(\mathbf{w})^2])^{1/2} \leq \|g\|_{L^2} \sum_{k=0}^p \eta^{-k} \|\mathbf{B}^*\|_{op}^k r^{\frac{p-k}{4}}.$$

Note that  $\|v\|_{L^2} = \mathcal{O}(\text{poly}(m_2, d))$  and  $|\sigma_1(\eta a \langle \mathbf{w}, \mathbf{z} \rangle + b)| \leq \eta a \langle \mathbf{w}, \mathbf{z} \rangle + b$  has polynomial growth. Since we have taken  $\iota = C \log(dn_1 n_2 m_1 m_2)$  for some sufficiently large  $C > 0$ , we know by Cauchy inequality,

$$|\mathbb{E}_{a,b,\mathbf{w}}[v(\mathbf{w})\sigma_1(\eta \langle \mathbf{w}, \mathbf{z} \rangle + b)\mathbf{1}\{\eta \langle \mathbf{w}, \mathbf{z} \rangle > 1\}]| \leq o\left(\frac{1}{dn_1 n_2 m_1 m_2}\right).$$

Thus, we then have that

$$\begin{aligned} & \mathbb{E}_{a,b,\mathbf{w}}[v(a, b, \mathbf{w})\sigma_1(\eta a \langle \mathbf{w}, \mathbf{z} \rangle + b)] \\ &= \mathbb{E}_{a,b,\mathbf{w}}[v(a, b, \mathbf{w})\sigma_1(\eta a \langle \mathbf{w}, \mathbf{z} \rangle + b) \cdot \mathbf{1}_{|\eta \langle \mathbf{w}, \mathbf{z} \rangle| \leq 1}] + o\left(\frac{1}{dn_1 n_2 m_1 m_2}\right) \\ &= \sum_{k=0}^p \mathbb{E}_{a,b,\mathbf{w}}[v_k(a, b)v_k(\mathbf{w})\sigma_1(\eta a \langle \mathbf{w}, \mathbf{z} \rangle + b) \cdot \mathbf{1}_{|\eta \langle \mathbf{w}, \mathbf{z} \rangle| \leq 1}] + o\left(\frac{1}{dn_1 n_2 m_1 m_2}\right) \\ &= \sum_{k=0}^p \mathbb{E}_{\mathbf{w}}[v_k(\mathbf{w})(\eta \langle \mathbf{w}, \mathbf{z} \rangle)^k \cdot \mathbf{1}_{|\eta \langle \mathbf{w}, \mathbf{z} \rangle| \leq 1}] + o\left(\frac{1}{dn_1 n_2 m_1 m_2}\right) \\ &= \sum_{k=0}^p \mathbb{E}_{\mathbf{w}}[v_k(\mathbf{w})(\eta \langle \mathbf{w}, \mathbf{z} \rangle)^k] + o\left(\frac{1}{dn_1 n_2 m_1 m_2}\right) \\ &= \sum_{k=0}^p \langle \mathbf{T}_k, (\mathbf{B}^* \mathbf{z})^{\otimes k} \rangle + o\left(\frac{1}{dn_1 n_2 m_1 m_2}\right) \\ &= g(\mathbf{B}^* \mathbf{z}) + o\left(\frac{1}{dn_1 n_2 m_1 m_2}\right). \end{aligned}$$

The proof is complete.  $\square$

*Proof of Lemma 27.* Fix  $\mathbf{x} \in \mathcal{D}_2$ . For notation simplicity, we denote  $f_v^\infty(\mathbf{x}) = f_{\infty, m_2}(\mathbf{x}; v)$ . Consider a truncation radius  $R > 0$  to be chosen later and let  $E_x$  be the set of  $\mathbf{w}$  such that

$$\sup_{a,b} |v(a, b, \mathbf{w})| \leq R \text{ and } \eta \langle \mathbf{w}, \mathbf{h}^{(1)}(\mathbf{x}) \rangle \leq 1.$$

By the construction of  $v(a, b, \mathbf{w})$  in the proof of Lemma 33, we know it can be seen as a degree- $p$  polynomial of  $\mathbf{w}$ . Thus, by Lemma 7, by taking  $R = C\iota^{p/2} \|v\|_{L^2}$  for some sufficiently large  $C > 0$ , we can ensure that

$$\Pr_{a,b}[\sup |v(a, b, \mathbf{w})| \leq R] \geq 1 - \exp(-\iota).$$

Moreover, by Proposition 4, conditional on a high probability event on  $\mathbf{V}$ ,  $\mathcal{D}_1$  and  $\mathcal{D}_2$ , by taking  $\eta = C\iota^{-5}d^6m_2^{-1}$ , we have  $\Pr[\eta \langle \mathbf{w}, \mathbf{h}^{(1)}(\mathbf{x}) \rangle \leq 1] \geq 1 - 4\exp(-\iota^2/2)$  for a single  $\mathbf{x}$ . Now consider the random variables

$$Z_i := \mathbf{1}\{\mathbf{w}_i \in E_x\}v(a_i, b_i, \mathbf{w}_i)\sigma_1(\eta a_i \langle \mathbf{w}_i, \mathbf{h}^{(1)}(\mathbf{x}) \rangle + b_i), \quad i = 1, 2, \dots, m_1.$$

We directly have that  $|Z_i| \lesssim \iota^{p/2} \|v\|_{L^2}$ , and with high probability,

$$\frac{1}{m_1} \sum_{i=1}^{m_1} v(a_i, b_i, \mathbf{w}_i)^2 \lesssim \iota^p \|v\|_{L^2}^2.$$

Therefore by Hoeffding inequality, with probability at least  $1 - 2\exp(-\iota)$ , we have

$$\begin{aligned} & \left| \frac{1}{m_1} \sum_{i=1}^{m_1} \mathbf{1}_{\mathbf{w}_i \in E_x} v(a_i, b_i, \mathbf{w}_i)\sigma_1(\eta a_i \langle \mathbf{w}_i, \mathbf{h}^{(1)}(\mathbf{x}) \rangle + b_i) \right. \\ & \quad \left. - \mathbb{E}[\mathbf{1}_{\mathbf{w} \in E_x} v(a, b, \mathbf{w})\sigma_1(\eta a \langle \mathbf{w}, \mathbf{h}^{(1)}(\mathbf{x}) \rangle + b)] \right| \lesssim \sqrt{\frac{\iota^{p+1} \|v\|_{L^2}^2}{m_1}}. \end{aligned}$$

Similar to the proof of Lemma 33, note that both  $v(a, b, \mathbf{w})$  and  $|\sigma_1(\eta a \langle \mathbf{w}, \mathbf{z} \rangle + b)|$  has polynomial growth. Since we have taken  $\iota = C \log(dn_1 n_2 m_1 m_2)$  for some sufficiently large  $C > 0$ , we know by Cauchy inequality,

$$\begin{aligned} & \left| \mathbb{E}[\mathbf{1}_{w \in E_x} v(a, b, \mathbf{w}) \sigma_1(\eta a \langle \mathbf{w}, \mathbf{h}^{(1)}(\mathbf{x}) \rangle + b)] - f_v^\infty(\mathbf{x}) \right| \\ &= \left| \mathbb{E}[\mathbf{1}_{w \notin E_x} v(a, b, \mathbf{w}) \sigma_1(\eta a \langle \mathbf{w}, \mathbf{h}^{(1)}(\mathbf{x}) \rangle + b)] \right| \\ &\leq \mathbb{P}(\mathbf{w} \notin E_x) (\mathbb{E}[v(a, b, \mathbf{w})^2 \sigma_1(\eta a \langle \mathbf{w}, \mathbf{h}^{(1)}(\mathbf{x}) \rangle + b)^2])^{1/2} \\ &\lesssim \exp(-C \log(dm_1 m_2 n_1 n_2)) \tilde{O}(\|v\|_{L^2}) \\ &\lesssim \frac{1}{m_1}. \end{aligned}$$

Finally, union bounding over  $\mathbf{x} \in \mathcal{D}_2$ , we see that

$$\begin{aligned} \sup_{\mathbf{x} \in \mathcal{D}_2} \left| \frac{1}{m_1} \sum_{i=1}^{m_1} v(a_i, b_i, \mathbf{w}_i) \sigma_1(\eta a_i \langle \mathbf{w}_i, \mathbf{h}^{(1)}(\mathbf{x}) \rangle + b_i) - f_v^\infty(\mathbf{x}) \right| &\lesssim \sqrt{\frac{\iota^{p+1} \|v\|_{L^2}^2}{m_1}} + \frac{1}{m_1} \\ &\lesssim \sqrt{\frac{\iota^{p+1} \|v\|_{L^2}^2}{m_1}}. \end{aligned}$$

The proof is complete.  $\square$

## E GENERALIZATION THEORY

### E.1 FORMAL PROOF OF THEOREM 1

The proof is divided into two parts. The first part of proof formalizes the proof we present in Section 4. The second part presents the generalization theory after we construct  $\mathbf{a}^*$  that gives small  $L^2$  error by Proposition 2, with the formal version presented in Proposition 6.

#### E.1.1 PART1: ANALYSIS BEFORE FEATURE RECONSTRUCTION

Denote  $\mathbf{w}_j = \epsilon^{-1} \mathbf{w}_j^{(0)} \sim \mathcal{N}(0, \mathbf{I}_{m_2})$ . Note that for any  $\mathbf{x} \in \mathcal{D}_1$  and  $j \in [m_1]$ , we have

$$\langle \mathbf{w}_j^{(0)}, \mathbf{h}^{(0)}(\mathbf{x}) \rangle = \langle \epsilon \mathbf{w}_j, \mathbf{h}^{(0)}(\mathbf{x}) \rangle \sim \mathcal{N}\left(0, \epsilon^2 \left\| \mathbf{h}^{(0)}(\mathbf{x}) \right\|_2^2\right).$$

Since  $\left\| \mathbf{h}^{(0)}(\mathbf{x}) \right\|_2^2 = \sum_{k=1}^{m_2} \sigma_2^2(\mathbf{v}_k^\top \mathbf{x}) \leq m_2 C_\sigma^2$ . By setting  $\epsilon^{-1} = C_\sigma \sqrt{2\iota m_2}$ , we know  $\langle \mathbf{w}_j^{(0)}, \mathbf{h}^{(0)}(\mathbf{x}) \rangle \leq 1$  with probability at least  $1 - 2 \exp(-\iota)$ . Thus, uniformly bounding over  $\mathbf{x} \in \mathcal{D}_1$  and  $j \in [m_1]$ , we know with high probability over  $\mathbf{W}$ , we have

$$\langle \mathbf{w}_j^{(0)}, \mathbf{h}^{(0)}(\mathbf{x}) \rangle \leq 1 \text{ for any } \mathbf{x} \in \mathcal{D}_1, \quad j \in [m_1].$$

Then, according Algorithm 1, after one-step gradient descent on  $\mathbf{W}$ , we know with high probability, for each  $j \in [m_1]$ ,

$$\begin{aligned} \eta_1 \nabla_{\mathbf{w}_j^{(0)}} \mathcal{L}(\theta^{(0)}) &= -\eta_1 \frac{a_j^{(0)}}{m_1} \cdot \frac{1}{n_1} \sum_{\mathbf{x} \in \mathcal{D}_1} f^*(\mathbf{x}) \mathbf{h}^{(0)}(\mathbf{x}) \sigma'_1(\langle \epsilon \mathbf{w}_j, \mathbf{h}^{(0)}(\mathbf{x}) \rangle) \\ &= -\frac{2\epsilon\eta_1}{m_1} a_j^{(0)} \cdot \frac{1}{n_1} \sum_{i=1}^n f^*(\mathbf{x}) \mathbf{h}^{(0)}(\mathbf{x}) \mathbf{h}^{(0)}(\mathbf{x})^\top \mathbf{w}_j, \end{aligned}$$

which is a linear transformation on  $\mathbf{w}_j$ . By taking  $\eta_1 = \frac{m_1}{2\epsilon m_2} \cdot \eta$  for some  $\eta > 0$  to be chosen later and  $\lambda_1 = \eta_1^{-1}$ , we have

$$\begin{aligned} \mathbf{w}_j^{(1)} &= \mathbf{w}_j^{(0)} - \eta_1 \left[ \nabla_{\mathbf{w}_j^{(0)}} \mathcal{L}(\theta^{(0)}) + \lambda_1 \mathbf{w}_j^{(0)} \right] \\ &= -\eta_1 \nabla_{\mathbf{w}_j^{(0)}} \mathcal{L}(\theta^{(0)}) \\ &= \frac{\eta a_j^{(0)}}{m_2} \cdot \frac{1}{n_1} \sum_{i=1}^n f^*(\mathbf{x}) \mathbf{h}^{(0)}(\mathbf{x}) \mathbf{h}^{(0)}(\mathbf{x})^\top \mathbf{w}_j. \end{aligned}$$

Then for any second-stage training sample  $\mathbf{x}' \in \mathcal{D}_2$ , the inner-layer neuron becomes

$$\begin{aligned} \left\langle \mathbf{w}_j^{(1)}, \sigma_2(\mathbf{V}\mathbf{x}') \right\rangle &= \frac{\eta a_j^{(0)}}{m_2} \left\langle \frac{1}{n_1} \sum_{i=1}^n f^*(\mathbf{x}) \mathbf{h}^{(0)}(\mathbf{x}) \mathbf{h}^{(0)}(\mathbf{x})^\top \mathbf{w}_j, \mathbf{h}^{(0)}(\mathbf{x}') \right\rangle \\ &= \eta a_j^{(0)} \cdot \left\langle \mathbf{w}_j, \frac{1}{n_1} \sum_{i=1}^n K_{m_2}^{(0)}(\mathbf{x}, \mathbf{x}') \mathbf{h}^{(0)}(\mathbf{x}) \right\rangle \\ &= \eta a_j^{(0)} \cdot \langle \mathbf{w}_j, \mathbf{h}^{(1)}(\mathbf{x}') \rangle. \end{aligned}$$

Thus, after the first training stage and reinitialization on  $\mathbf{b} = \mathbf{b}^{(1)}$ , the model becomes the following random-feature model in the second stage:

$$f(\mathbf{x}'; \theta) = \frac{1}{m_1} \sum_{j=1}^{m_1} a_j \sigma_1 \left( \eta a_j^{(0)} \langle \mathbf{w}_j, \mathbf{h}^{(1)}(\mathbf{x}') \rangle + b_j^{(1)} \right).$$

By Proposition 6, we know there exists  $\mathbf{a}^* \in \mathbb{R}^{m_1}$  such that with high probability over  $\mathcal{D}_1, \mathcal{D}_2$ ,  $\{\mathbf{w}_i\}_{i=1}^{m_1}$  and  $\mathbf{V}$ , by taking the parameter  $\theta^* = (\mathbf{a}^*, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{V})$ , it holds that

$$\begin{aligned} \hat{\mathcal{L}}_2(\theta^*) &\lesssim \|g\|_{L^2}^2 \cdot \frac{r^p}{\lambda_{\min}(\mathbf{H})} \cdot \left( \frac{\iota^{p+2} d^5}{m_2} + \frac{\iota d^3}{\sqrt{m_2}} + \frac{\iota^{p+3/2} d}{\sqrt{n_1}} + \frac{\iota L r^2 \kappa_1 \log^2 d}{d^{1/6}} \right)^2 \\ &\quad + \frac{\iota^{p+1} \|g\|_{L^2}^2}{m_1} \cdot \left( \sum_{k=0}^p \eta^{-k} \|\mathbf{B}^*\|_{\text{op}}^k r^{\frac{p-k}{4}} \right)^2. \end{aligned}$$

Here  $\mathbf{a}^*$  satisfies

$$\frac{\|\mathbf{a}^*\|_2^2}{m_1} \lesssim \iota^p \|g\|_{L^2}^2 \cdot \left( \sum_{k=0}^p \eta^{-k} \|\mathbf{B}^*\|_{\text{op}}^k r^{\frac{p-k}{4}} \right)^2.$$

The first part of the proof is complete.

### E.1.2 PART2: GENERALIZATION THEORY

Denote the population absolute loss as  $\mathcal{L}_1(f, g) = \mathbb{E}_{\mathbf{x}} [|f(\mathbf{x}) - g(\mathbf{x})|]$ . Moreover, we consider a truncated loss function as

$$\ell_\tau(z) = \min(|z|, \tau) \text{ and } \mathcal{L}_{1,\tau}(f, g) = \mathbb{E}_{\mathbf{x}} [\ell_\tau(f(\mathbf{x}) - g(\mathbf{x}))],$$

where  $\tau > 0$  is the truncation radius. Moreover, we denote the empirical truncated absolute loss as

$$\hat{\mathcal{L}}_{1,\tau}(f, g) = \frac{1}{n_2} \sum_{\mathbf{x} \in \mathcal{D}_2} \ell_\tau(f(\mathbf{x}) - g(\mathbf{x})).$$

Suppose Algorithm 1 gives rise to a set of parameters  $\hat{\theta} = (\hat{\mathbf{a}}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{V})$ , and we have constructed  $\theta^* = (\mathbf{a}^*, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{V})$  that leads to small empirical loss, we decompose the population absolute loss as

$$\begin{aligned} \mathcal{L}_1(f(\cdot; \hat{\theta}), f^*) &= \underbrace{\hat{\mathcal{L}}_{1,\tau}(f(\cdot; \hat{\theta}), f^*)}_{L_1} + \underbrace{\mathcal{L}_{1,\tau}(f(\cdot; \hat{\theta}), f^*) - \hat{\mathcal{L}}_{1,\tau}(f(\cdot; \hat{\theta}), f^*)}_{L_2} \\ &\quad + \underbrace{\mathcal{L}_1(f(\cdot; \hat{\theta}), f^*) - \mathcal{L}_{1,\tau}(f(\cdot; \hat{\theta}), f^*)}_{L_3}. \end{aligned}$$

Here with a little abuse of notation, we consider  $f^* = g^*(\mathbf{p})$  for learning the original target function and denote  $f^* = g(\mathbf{p})$  with  $g$  being any degree  $p$  polynomial for the transfer learning setting. Next, we bound  $L_1, L_2$  and  $L_3$  respectively.

**Bound  $L_1$**  With a little abuse of notation, we denote  $\hat{\mathcal{L}}_2(\mathbf{a}) = \hat{\mathcal{L}}_2(\theta)$  for  $\theta = (\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{V})$  since we only optimize  $\mathbf{a}$  in the second stage. By Proposition 6, we know with high probability, the

2754 empirical  $L^2$  loss of  $\theta^*$  is bounded by  
 2755

$$\begin{aligned} \hat{\mathcal{L}}_2(\mathbf{a}^*) &= \frac{1}{n_2} \sum_{\mathbf{x} \in \mathcal{D}_2} (f(\mathbf{x}; \theta^*) - f^*(\mathbf{x}))^2 \\ &\lesssim \|g\|_{L^2}^2 \cdot \frac{r^p}{\lambda_{\min}^2(\mathbf{H})} \cdot \left( \frac{\iota^{p+2} d^5}{m_2} + \frac{\iota d^3}{\sqrt{m_2}} + \frac{\iota^{p+3/2} d}{\sqrt{n_1}} + \frac{\iota L r^2 \kappa_1 \log^2 d}{d^{1/6}} \right)^2 \\ &\quad + \frac{\iota^{p+1} \|g\|_{L^2}^2}{m_1} \cdot \left( \sum_{k=0}^p \eta^{-k} \|\mathbf{B}^*\|_{\text{op}}^k r^{\frac{p-k}{4}} \right)^2. \end{aligned}$$

2761 Here  $\mathbf{a}^*$  satisfies  
 2762

$$\frac{\|\mathbf{a}^*\|_2^2}{m_1} \lesssim \iota^p \|g\|_{L^2}^2 \cdot \left( \sum_{k=0}^p \eta^{-k} \|\mathbf{B}^*\|_{\text{op}}^k r^{\frac{p-k}{4}} \right)^2.$$

2763 In the second training stage, let's set the weight decay in the second training stage as  
 2764

$$\begin{aligned} \lambda_2 = \lambda &= \|\mathbf{a}^*\|_2^{-2} \|g\|_{L^2}^2 \cdot \left( \frac{r^p}{\lambda_{\min}^2(\mathbf{H})} \cdot \left( \frac{\iota^{p+2} d^5}{m_2} + \frac{\iota d^3}{\sqrt{m_2}} + \frac{\iota^{p+3/2} d}{\sqrt{n_1}} + \frac{\iota L r^2 \kappa_1 \log^2 d}{d^{1/6}} \right)^2 \right. \\ &\quad \left. + \frac{\iota^{p+1}}{m_1} \cdot \left( \sum_{k=0}^p \eta^{-k} \|\mathbf{B}^*\|_{\text{op}}^k r^{\frac{p-k}{4}} \right)^2 \right) \end{aligned}$$

2775 so that the empirical  $L^2$  loss is directly bounded by  
 2776

$$\hat{\mathcal{L}}_2(\mathbf{a}^*) := \frac{1}{n_2} \sum_{\mathbf{x} \in \mathcal{D}_2} (f(\mathbf{x}; \theta^*) - f^*(\mathbf{x}))^2 \lesssim \lambda \|\mathbf{a}^*\|_2^2.$$

2777 We further consider the regularized second-stage training loss to be  
 2778

$$\hat{\mathcal{L}}_{2,\lambda}(\mathbf{a}) = \frac{1}{n_2} \sum_{\mathbf{x} \in \mathcal{D}_2} (f(\mathbf{x}; (\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{V})) - f^*(\mathbf{x}))^2 + \frac{\lambda}{2} \|\mathbf{a}\|_2^2.$$

2782 Note that this loss is strongly convex, so it has a global minimum  $\mathbf{a}^{(\infty)} = \operatorname{argmin} \hat{\mathcal{L}}_{2,\lambda}(\mathbf{a})$ . Thus,  
 2783 we have

$$\mathcal{L}_{2,\lambda}(\mathbf{a}^{(\infty)}) \leq \mathcal{L}_{2,\lambda}(\mathbf{a}^*) \lesssim \lambda \|\mathbf{a}^*\|_2^2.$$

2785 Since  $\mathcal{L}_{2,\lambda}(\mathbf{a})$  is  $\lambda$ -strongly convex, and we can write  $f(\mathbf{x}; (\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{V})) = \mathbf{a}^\top \Psi(\mathbf{x})$ , where  
 2786  $\Psi(\mathbf{x}) = \operatorname{Vec}(m_1^{-1} \sigma_1(\eta a_i^{(0)} \langle \mathbf{w}_i, \mathbf{h}^{(1)}(\mathbf{x}) \rangle + b_i^{(1)}))$ . Therefore, by Lemma 4 and our choice of  $\eta$   
 2787 to ensure  $\eta \langle \mathbf{w}_i, \mathbf{h}^{(1)}(\mathbf{x}) \rangle \leq 1$  with high probability, we know with high probability,  
 2788

$$\lambda_{\max}(\nabla_{\mathbf{a}}^2 \hat{\mathcal{L}}_{2,\lambda}) \leq \frac{2}{n_2} \sum_{\mathbf{x} \in \mathcal{D}_2} \|\Psi(\mathbf{x})\|_2 \lesssim \frac{1}{m_1}.$$

2789 Thus,  $\mathcal{L}_{2,\lambda}(\mathbf{a})$  is  $\lambda + \mathcal{O}(\frac{1}{m_1})$ -smooth. By choosing the second-stage learning rate  $\eta_2 = \Omega(m_1)$ ,  
 2790 after  $T = \tilde{\mathcal{O}}(\lambda^{-1}) = \operatorname{poly}(d, n, m_1, m_2, \|g\|_{L^2})$  steps, we can reach an iterate  $\hat{\mathbf{a}} = \mathbf{a}^{(T)}$  so that  
 2791

$$\hat{\mathcal{L}}_2(\hat{\mathbf{a}}) \lesssim \hat{\mathcal{L}}_2(\mathbf{a}^*) \text{ and } \|\hat{\mathbf{a}}\|_2 \lesssim \|\mathbf{a}^*\|_2.$$

2792 Denoting  $\hat{\theta} = (\hat{\mathbf{a}}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{V})$ , it holds that  
 2793

$$\hat{\mathcal{L}}_{1,\tau}(f(\cdot; \hat{\theta}), f^*) \leq \frac{1}{n_2} \sum_{\mathbf{x} \in \mathcal{D}_2} |f(\mathbf{x}; \hat{\theta}) - f^*(\mathbf{x})| \leq \sqrt{\hat{\mathcal{L}}_2(\hat{\mathbf{a}})} \leq \sqrt{\hat{\mathcal{L}}_2(\mathbf{a}^*)}.$$

2794 Thus, we have  
 2795

$$\begin{aligned} L_1 &= \hat{\mathcal{L}}_{1,\tau}(f(\cdot; \hat{\theta}), f^*) \leq \sqrt{\frac{1}{n_2} \sum_{\mathbf{x} \in \mathcal{D}_2} (f(\mathbf{x}; \theta^*) - f^*(\mathbf{x}))^2} \\ &\lesssim \|g\|_{L^2} \cdot \frac{r^{p/2}}{\lambda_{\min}(\mathbf{H})} \cdot \left( \frac{\iota^{p+2} d^5}{m_2} + \frac{\iota d^3}{\sqrt{m_2}} + \frac{\iota^{p+3/2} d}{\sqrt{n}} + \frac{\iota L r^2 \kappa_1 \log^2 d}{d^{1/6}} \right) \\ &\quad + \sqrt{\frac{\iota^{p+1} \|g\|_{L^2}^2}{m_1}} \cdot \left( \sum_{k=0}^p \eta^{-k} \|\mathbf{B}^*\|_{\text{op}}^k r^{\frac{p-k}{4}} \right). \end{aligned}$$

2808 Here  $\hat{\mathbf{a}}$  satisfies  
 2809

2810

$$2811 \frac{\|\hat{\mathbf{a}}\|_2^2}{m_1} \lesssim \frac{\|\mathbf{a}^*\|_2^2}{m_1} \lesssim \iota^p \|g\|_{L^2}^2 \cdot \left( \sum_{k=0}^p \eta^{-k} \|\mathbf{B}^*\|_{\text{op}}^k r^{\frac{p-k}{4}} \right)^2.$$

2813

2814

2815 We assume  $\|\mathbf{a}\|_2^2 \leq m_1 B_a^2$ , where  $B_a$  satisfies  
 2816

2817

2818

$$2819 B_a^2 \lesssim \iota^p \|g\|_{L^2}^2 \cdot \left( \sum_{k=0}^p \eta^{-k} \|\mathbf{B}^*\|_{\text{op}}^k r^{\frac{p-k}{4}} \right)^2.$$

2820

2821

2822

2823 **Bound  $L_2$**  To bound  $L_2$ , we rely on standard Rademacher complexity analysis. The following  
 2824 lemma provides an upper bound on the Rademacher complexity of the random feature model.

2825

2826

2827 **Lemma 35.** Let  $\mathcal{F} = \{f_\theta : \theta = (\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{V}), \|\mathbf{a}\|_2 \leq \sqrt{m_1} B_a\}$ . Recall the empirical  
 2828 Rademacher complexity of  $\mathcal{F}$  as

2829

$$2830 \mathcal{R}_n(\mathcal{F}) = \mathbb{E}_{\sigma \in \{\pm 1\}^n} \left[ \sup_{f \in \mathcal{F}} \frac{1}{n_2} \sum_{i=1}^n \sigma_i f(\mathbf{x}_i) \right],$$

2831

2832

2833

2834 Here the dataset  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n_2}\} = \mathcal{D}_2$ . Then with high probability, we have

2835

2836

2837

$$\mathcal{R}_n(\mathcal{F}) \lesssim \frac{B_a}{\sqrt{n_2}}.$$

2838

2839

2840

2841 The proof is provided in Appendix E.2. Since the  $\ell_\tau$  is 1-Lipschitz, by standard Rademacher com-  
 2842 plexity analysis, we have that with high probability that

2843

2844

2845

$$2846 L_2 = \mathbb{E}_{\mathbf{x}} \ell_\tau(f(\mathbf{x}; \hat{\theta}) - f^*(\mathbf{x})) - \frac{1}{n_2} \sum_{\mathbf{x} \in \mathcal{D}_2} \ell_\tau(f(\mathbf{x}; \hat{\theta}) - f^*(\mathbf{x})) \\ 2847 \lesssim \mathcal{R}_n(\mathcal{F}) + \tau \sqrt{\frac{\iota}{n_2}} \\ 2848 \lesssim \sqrt{\frac{B_a^2}{n_2}} + \tau \sqrt{\frac{\iota}{n_2}}.$$

2849

2850

2851

2852

2853

2854 **Bound  $L_3$**  Finally, we relate the truncated loss  $\ell_\tau$  to the  $L_1$  population loss.

2855

2856

2857

2858 **Lemma 36.** By letting  $\tau = \Omega(\max(\iota^p, B_a))$ , with high probability over  $\hat{\theta}$ , we have

2859

2860

2861

$$2862 L_3 = \mathbb{E}_{\mathbf{x}} [ |f(\mathbf{x}; \hat{\theta}) - f^*(\mathbf{x})| ] - \mathbb{E}_{\mathbf{x}} [\ell_\tau(f(\mathbf{x}; \hat{\theta}) - f^*(\mathbf{x}))] \leq o\left(\frac{1}{n_1 n_2 m_1 m_2 d}\right).$$

Here we recall that  $n = n_1 + n_2$ . The proof is provided in Appendix E.2.

**Put the loss together** By invoking the upper bound of  $L_1$ ,  $L_2$  and  $L_3$  and plugging the values of  $\tau$ ,  $\eta$ ,  $\|\mathbf{B}^*\|_{\text{op}}$ ,  $\lambda_{\min}(\mathbf{H})$ ,  $B_a$ ,  $L$  and  $\|g\|_{L^2}$ , we have

$$\begin{aligned}
 \mathbb{E}_{\mathbf{x}} [ |f(\mathbf{x}; \hat{\theta}) - f^*(\mathbf{x})| ] &= L_1 + L_2 + L_3 \\
 &\lesssim \frac{1}{n_2} \sum_{\mathbf{x} \in \mathcal{D}_2} \ell_\tau(f(\mathbf{x}; \theta^*) - f^*(\mathbf{x})) + \sqrt{\frac{B_a^2}{n_2} + \tau \sqrt{\frac{\iota}{n_2}}} + o\left(\frac{1}{n_1 n_2 m_1 m_2 d}\right) \\
 &\lesssim \|g\|_{L^2} \cdot \frac{r^{p/2}}{\lambda_{\min}(\mathbf{H})} \cdot \left( \frac{\iota^{p+2} d^5}{m_2} + \frac{\iota d^3}{\sqrt{m_2}} + \frac{\iota^{p+3/2} d}{\sqrt{n_1}} + \frac{\iota L r^2 \kappa_1 \log^2 d}{d^{1/6}} \right) \\
 &\quad + \sqrt{\frac{\iota^{p+1} \|g\|_{L^2}}{m_1}} \cdot \left( \sum_{k=0}^p \eta^{-k} \|\mathbf{B}^*\|_{\text{op}}^k r^{\frac{p-k}{4}} \right) + \sqrt{\frac{B_a^2}{n_2} + \tau \sqrt{\frac{\iota}{n_2}}} \\
 &\lesssim \frac{r^{p/2}}{\lambda_{\min}(\mathbf{H})} \cdot \left( \frac{\iota^{p+2} d^5}{m_2} + \frac{\iota d^3}{\sqrt{m_2}} + \frac{\iota^{p+3/2} d}{\sqrt{n_1}} + \frac{\iota L r^2 \kappa_1 \log^2 d}{d^{1/6}} \right) \\
 &\quad + \sqrt{\frac{\iota^{6p+1} r^{p/2} \kappa_2^{2p} (r^{1/4} \vee \lambda_{\min}^{-1}(\mathbf{H}))^p}{m_1}} \\
 &\quad + \sqrt{\frac{\iota^{6p+1} r^{p/2} \kappa_2^{2p} (r^{1/4} \vee \lambda_{\min}^{-1}(\mathbf{H}))^p}{n_2}} \\
 &= \tilde{\mathcal{O}} \left( \sqrt{\frac{r^p \kappa_2^{2p}}{\min(n_2, m_1)}} + \sqrt{\frac{d^6 r^{p+1}}{m_2}} + \sqrt{\frac{d^2 r^{p+1}}{n}} + \frac{r^{p+2} \kappa_1}{d^{1/6}} \right).
 \end{aligned}$$

The proof is complete.

## E.2 OMITTED PROOFS IN APPENDIX E.1

*Proof of Lemma 35.* Given  $\theta = (\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{V})$ , since we can write

$$f_\theta(\mathbf{x}) = \mathbf{a}^\top \Psi(\mathbf{x}), \text{ where } \Psi(\mathbf{x}) = \text{Vec} \left( m_1^{-1} \sigma_1 \left( \eta a_i^{(0)} \langle \mathbf{w}_i, \mathbf{h}^{(1)}(\mathbf{x}) \rangle + b_i^{(1)} \right) \right).$$

By Proposition 4 and our choice of  $\eta$  to ensure  $|\eta \langle \mathbf{w}_i, \mathbf{h}^{(1)}(\mathbf{x}) \rangle| \leq 1$  with high probability for any  $\mathbf{x} \in \mathcal{D}_2$ , we obtain that for any  $i \in [m_1]$  and  $\mathbf{x} \in \mathcal{D}_2$ ,

$$\left| \eta a_i^{(0)} \langle \mathbf{w}_i, \mathbf{h}^{(1)}(\mathbf{x}) \rangle + b_i^{(1)} \right| \leq a_i^{(0)} |\eta \langle \mathbf{w}_i, \mathbf{h}^{(1)}(\mathbf{x}) \rangle| + b_i^{(1)} \lesssim 1.$$

Thus,  $\|\Psi(\mathbf{x})\|_2^2 \leq m_1^{-1}$ . by the standard linear Rademacher bound, with high probability, the empirical Rademacher complexity is upper bounded by

$$\mathcal{R}_n(\mathcal{F}) \lesssim \frac{\sqrt{m_1} B_a}{n_2} \sqrt{\sum_{\mathbf{x} \in \mathcal{D}_2} \|\Psi(\mathbf{x})\|_2^2} \leq \frac{B_a}{\sqrt{n_2}}.$$

The proof is complete.  $\square$

*Proof of Lemma 36.* We can bound the difference between  $\ell_\tau$  and  $L_1$  loss by

$$\begin{aligned}
 &\mathbb{E}_{\mathbf{x}} [ |f(\mathbf{x}; \hat{\theta}) - f^*(\mathbf{x})| ] - \mathbb{E}_{\mathbf{x}} [ \ell_\tau(f(\mathbf{x}; \hat{\theta}) - f^*(\mathbf{x})) ] \\
 &\leq \mathbb{E}_{\mathbf{x}} [ |f(\mathbf{x}; \hat{\theta}) - f^*(\mathbf{x})| \mathbf{1} \{ |f(\mathbf{x}; \hat{\theta}) - f^*(\mathbf{x})| \geq \tau \} ] \\
 &\leq \sqrt{\mathbb{E}_{\mathbf{x}} [ |f(\mathbf{x}; \hat{\theta}) - f^*(\mathbf{x})|^2 ] \Pr [|f(\mathbf{x}; \hat{\theta}) - f^*(\mathbf{x})| \geq \tau]} \\
 &\lesssim \sqrt{\mathbb{E}_{\mathbf{x}} [ (f(\mathbf{x}; \hat{\theta})^2 + f^*(\mathbf{x})^2) ] [\Pr [|f(\mathbf{x}; \hat{\theta})| \geq \tau/2] + \Pr [|f^*(\mathbf{x})| \geq \tau/2]]} \tag{36}
 \end{aligned}$$

2916 Recall that we can write  
 2917

$$2918 \quad f(\mathbf{x}; \hat{\theta}) = \hat{\mathbf{a}}^\top \Psi(\mathbf{x}), \text{ where } \Psi(\mathbf{x}) = \text{Vec} \left( m_1^{-1} \sigma_1 \left( \eta a_i^{(0)} \langle \mathbf{w}_i, \mathbf{h}^{(1)}(\mathbf{x}) \rangle + b_i^{(1)} \right) \right).$$

2919 By following the proof of Lemma 35 and applying Proposition 4 for one single sample point  $\mathbf{x}$   
 2920 (instead of the whole set  $\mathcal{D}_2$ ), we know with high probability over  $\mathbf{V}, \mathbf{w}$  and  $\mathcal{D}_1$  (we denote this  
 2921 event by  $E_1$ ), we have for any  $i \in [m_1]$ ,  
 2922

$$2923 \quad \left| \eta \langle \mathbf{w}_i, \mathbf{h}^{(1)}(\mathbf{x}) \rangle \right| \leq 1$$

2925 holds with high probability on  $\mathbf{x}$ . Also, since for any  $i \in [m_1]$ ,  $\mathbf{w}_i \sim \mathcal{N}(\mathbf{0}_{m_2}, \mathbf{I}_{m_2})$ , we  
 2926 know  $\|\mathbf{w}_i\|_2 \lesssim \sqrt{m_2 \iota}$  for any  $i \in [m_1]$  with high probability. We denote this joint event on  
 2927  $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{m_1}$  by  $E_2$ . Thus, conditional on events  $E_1$  and  $E_2$ , we have

$$2928 \quad \left| f(\mathbf{x}; \hat{\theta}) \right| \lesssim \frac{\|\hat{\mathbf{a}}\|_2}{\sqrt{m_1}} \text{ with high probability on } \mathbf{x}.$$

2929 We denote this conditional event by  $E_{x,1}$ . Moreover, since  $\eta = C\iota^{-5}m_2^{-1/2}d^6$ , we have  
 2930

$$\begin{aligned} 2933 \quad \left| f(\mathbf{x}; \hat{\theta}) \right| &\leq \frac{\|\hat{\mathbf{a}}\|_2}{m_1} \sum_{j=1}^{m_1} \left( \eta \|\mathbf{w}_j\|_2 \left\| \frac{1}{n_2} \sum_{i=1}^n K_{m_2}^{(0)}(\mathbf{x}_i, \mathbf{x}') \mathbf{h}^{(0)}(\mathbf{x}_i) \right\|_2 + 3 \right) \\ 2936 \quad &\leq \frac{\sqrt{m_2 \iota} \|\hat{\mathbf{a}}\|_2}{m_1} \sum_{j=1}^{m_1} \frac{\eta}{n} \sum_{i=1}^n C_\sigma^2 \cdot \sqrt{m_2} C_\sigma + 3 \\ 2939 \quad &\lesssim \sqrt{m_2 d^6} \|\hat{\mathbf{a}}\|_2 \end{aligned}$$

2940 holds for any  $\mathbf{x}$ . Moreover, since  $f^*(\mathbf{x})$  is a degree- $2p$  polynomial of  $\mathbf{x}$ , we know by Lemma  
 2941 8, with probability at least  $1 - \exp(-\iota)$ , we have  $|f^*| \leq C_f \iota^p$  for sufficiently large  $C_f > 0$ .  
 2942 Besides, we have  $\mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x})^2] \lesssim 1$ . Altogether, conditional on  $E_1$  and  $E_2$ , by choosing  $\tau =$   
 2943  $C' \max(\iota^p, m_1^{-1/2} \|\hat{\mathbf{a}}\|_2) = \Omega(\max(\iota^p, B_a))$  for some sufficiently large  $C$ , we have

$$\begin{aligned} 2945 \quad &\mathbb{E}_{\mathbf{x}} \left[ (f(\mathbf{x}; \hat{\theta})^2 + f^*(\mathbf{x})^2) \right] \left[ \Pr \left[ |f(\mathbf{x}; \hat{\theta})| \geq \tau/2 \right] + \Pr [|f^*(\mathbf{x})| \geq \tau/2] \right] \\ 2947 \quad &\lesssim (m_2 d^{12} \|\hat{\mathbf{a}}\|_2^2 + 1) (\Pr[\mathbf{x} \notin E_{x,1}] + \Pr[\mathbf{x} \notin E_{x,2}]) \\ 2949 \quad &\lesssim o \left( \frac{1}{d^2 m_1^2 m_2^2 n_1^2 n_2^2} \right). \end{aligned}$$

2950 The last inequality holds because of the definition of high probability events and the choice of  $\iota$   
 2951 with  $\iota = C \log(dm_1 m_2 n_1 n_2)$  for sufficiently large  $C$ . Plugging the result into (36) concludes our  
 2952 proof.  $\square$

2954  
 2955  
 2956  
 2957  
 2958  
 2959  
 2960  
 2961  
 2962  
 2963  
 2964  
 2965  
 2966  
 2967  
 2968  
 2969