

---

# Gradient-Based Feature Learning under Structured Data

---

Alireza Mousavi-Hosseini<sup>1</sup>, Denny Wu<sup>2</sup>, Taiji Suzuki<sup>3</sup>, Murat A. Erdogdu<sup>1</sup>

<sup>1</sup>University of Toronto and Vector Institute,

<sup>2</sup>New York University and Flatiron Institute,

<sup>3</sup>University of Tokyo and RIKEN AIP

{mousavi, erdogdu}@cs.toronto.edu, dennywu@nyu.edu,  
taiji@mist.i.u-tokyo.ac.jp

## Abstract

Recent works have demonstrated that the sample complexity of gradient-based learning of single index models, i.e. functions that depend on a 1-dimensional projection of the input data, is governed by their *information exponent*. However, these results are only concerned with isotropic data, while in practice the input often contains additional structure which can implicitly guide the algorithm. In this work, we investigate the effect of a *spiked covariance* structure and reveal several interesting phenomena. First, we show that in the anisotropic setting, the commonly used spherical gradient dynamics may fail to recover the true direction, even when the spike is perfectly aligned with the target direction. Next, we show that appropriate weight normalization that is reminiscent of *batch normalization* can alleviate this issue. Further, by exploiting the alignment between the (spiked) input covariance and the target, we obtain improved sample complexity compared to the isotropic case. In particular, under the spiked model with a suitably large spike, the sample complexity of gradient-based training can be made independent of the information exponent while also outperforming lower bounds for rotationally invariant kernel methods.

## 1 Introduction

A fundamental feature of neural networks is their *adaptivity* to learn unknown statistical models. For instance, when the learning problem exhibits certain low-dimensional structure or sparsity, it is expected that neural networks optimized by gradient-based algorithms can efficiently adapt to such structure via feature/representation learning. A considerable amount of research has been dedicated to understanding this phenomenon under various assumptions and to demonstrate the superiority of neural networks over non-adaptive methods such as kernel models [GMMM19, WLLM19, BES<sup>+</sup>19, LMZ20, AAM22, BES<sup>+</sup>22, DLS22, Tei23, MHPG<sup>+</sup>23].

A particular relevant problem setting for feature learning is the estimation of single index models, where the response  $y \in \mathbb{R}$  depends on the input  $\mathbf{x} \in \mathbb{R}^d$  via  $y = g(\langle \mathbf{u}, \mathbf{x} \rangle) + \epsilon$ , where  $g : \mathbb{R} \rightarrow \mathbb{R}$  is the nonlinear link function and  $\mathbf{u}$  is the unit target direction. Here, learning corresponds to recovering the unknowns  $\mathbf{u}$  and  $g$ , which requires the model to extract and adapt to the low-dimensional target direction. Recent works have shown that the sample complexity is determined by certain properties of the link function  $g$ . In particular, the complexity of gradient-based optimization is captured by the *information exponent* of  $g$  introduced by [BAGJ21]. Intuitively, a larger information exponent  $s$  corresponds to a more complex  $g$  (for gradient-based learning), and it has been proven that when the input is isotropic  $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$ , gradient flow can learn the single index model with  $\tilde{O}(d^s)$  sample complexity [BBSS22].

In practice, however, real data always exhibits certain structures such as low intrinsic dimensionality, and isotropic data assumptions fail to capture this fact. In statistics methodology, it is known that the directions along which the input  $\mathbf{x}$  has high variance are often good predictors of the target  $y$  [HTFF09]; indeed, this is the main reason principal component analysis is used in pre-training [JWHT13]. A fundamental model that captures such a structure is the *spiked matrix model* in which  $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d + \kappa \boldsymbol{\theta} \boldsymbol{\theta}^\top)$  for some unit direction  $\boldsymbol{\theta} \in \mathbb{R}^d$  and  $\kappa > 0$  [Joh01]. Along the direction  $\boldsymbol{\theta}$ , data has higher variability and predictive power. In single index models, such predictive power translates to a non-trivial alignment between the vectors  $\mathbf{u}$  and  $\boldsymbol{\theta}$  — our focus is to investigate the effect of such alignment on the sample complexity of gradient-based training.

### 1.1 Contributions: learning single index models under spiked covariance

In this paper, we study the sample complexity of learning a single index model using a two-layer neural network and show that it is determined by an interplay between

- **spike-target alignment:**  $\langle \mathbf{u}, \boldsymbol{\theta} \rangle \asymp d^{-r_1}$ ,  $r_1 \in [0, 1/2]$ ,
- **spike magnitude:**  $\kappa \asymp d^{r_2}$ , for  $r_2 \in [0, 1]$ .

Our contributions can be summarized as follows.

1. We show that even in the case of perfect spike-target alignment ( $r_1 = 0$ ), the spherical gradient flow commonly employed in recent literature (see e.g. [BAGJ21, BBSS22]) cannot recover the target direction for moderate spike magnitudes in the population limit. The failure of this covariance-agnostic procedure under anisotropic structure insinuates the necessity of an appropriate covariance-aware normalization to effectively learn the single index model.
2. We show that a covariance-aware normalization that resembles *batch normalization* resolves this issue. Indeed, the resulting gradient flow can successfully recover the target direction  $\mathbf{u}$  in this case, and depending on the amount of spike-target alignment, the sample complexity can significantly improve compared to the isotropic case.
3. Under the spiked covariance model, we prove a three-stage phase transition for the sample complexity depending on the quantities  $r_1$  and  $r_2$ . For a suitable direction and magnitude of the spike, the sample complexity can be made  $\tilde{\mathcal{O}}(d^{3+\nu})$  for any  $\nu > 0$  which is independent of the information exponent  $s$ . This should be compared against the known complexity of  $\tilde{\mathcal{O}}(d^s)$  under isotropic data.
4. We finally show that preconditioning the training dynamics with the inverse covariance improves the sample complexity. This is particularly significant for the spiked covariance model where  $\tilde{\mathcal{O}}(d^{3+\nu})$  samples can be reduced to  $\tilde{\mathcal{O}}(d^{1+\nu})$  for any  $\nu > 0$ , i.e. almost linear in  $d$ . The three-stage phase transition also emerges, as illustrated in Figure 1: in the “hard” regime, the complexity remains  $\tilde{\mathcal{O}}(d^s)$  regardless of the magnitude and direction of the spike, while in the “easy” regime the complexity only depends on the spike magnitude and not its direction. The “intermediate” regime interpolates between these two; smaller  $r_1$  and larger  $r_2$  improve the sample complexity.

The rest of the paper is organized as follows. We discuss the notation and the related work in the remainder of this section. We provide preliminaries on the statistical model and the training procedure in Section 2, and provide a negative result on the covariance-agnostic gradient flow in Section 2.1. Our main sample complexity result on a single neuron is presented in Section 3.2. We provide our results on multi-neuron neural networks in Section 4 and also discuss extensions such as preconditioning and its implications. We provide a technical summary in Section 5 and conclude in Section 6.

**Notation.** We use  $\langle \cdot, \cdot \rangle$  and  $\|\cdot\|$  to denote Euclidean inner product and norm. For matrices,  $\|\cdot\|$  denotes the usual operator norm, and  $\lambda_{\max}(\cdot)$  and  $\lambda_{\min}(\cdot)$  denote the largest and smallest eigenvalues respectively. We reserve  $\gamma$  for the standard Gaussian distribution on  $\mathbb{R}$ , and let  $\|\cdot\|_\gamma$  denote the  $L^2(\gamma)$  norm.  $\mathbb{S}^{d-1}$  is the unit  $d$ -dimensional sphere. For quantities  $a$  and  $b$ , we will use  $a \lesssim b$  to convey there exists a constant  $C$  (a universal constant unless stated otherwise, in which case may depend on polylogarithmic factors of  $d$ ) such that  $a \leq Cb$ , and  $a \asymp b$  signifies that  $a \lesssim b$  and  $b \lesssim a$ .

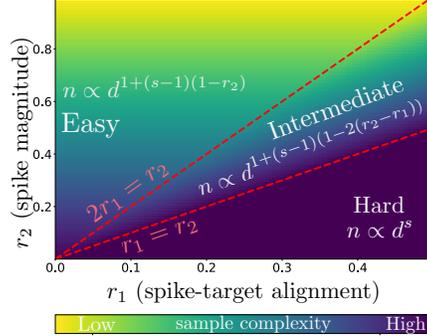


Figure 1: Sample complexity to learn  $\mathbf{u}$  and  $g$  under the spiked model. Smaller  $r_1$  denotes a better spike-target alignment, while larger  $r_2$  denotes a larger spike magnitude. The sample complexities are based on Corollary 8.

## 1.2 Further related work

**Non-Linear Feature Learning with Neural Networks.** Recently, two popular scaling regimes of neural networks have emerged for theoretical studies. A large initialization variance leads to the *lazy training* regime, where the weights do not move significantly, and the training dynamics is captured by the neural tangent kernel (NTK) [JGH18, COB19]. However, there are many instances of function classes that are efficiently learnable by neural networks and not efficiently learnable by the NTK [YS19, GMMM19]. Under a smaller initialization scaling, gradient descent on infinite-width neural networks becomes equivalent to Wasserstein gradient flow on the space of measures, known as the mean-field limit [CB18, RVE18, MMN18, MMM19, NWS22, Chi22], which can learn certain low-dimensional target functions efficiently [WLLM19, AAM22, HC22, ASKL23].

As for neural networks with smaller width, recent works showed that a two-stage feature learning procedure can outperform the NTK when the data is sampled uniformly from the hypercube [BEG<sup>+</sup>22] or isotropic Gaussian [DLS22, BES<sup>+</sup>22, BBSS22, MHPG<sup>+</sup>23, ABAM23]. However, these results do not take into account the additional structure that might be present in the covariance matrix of the input data. Two notable exceptions are [GMMM20, RGKZ21], where the authors analyzed a spiked covariance and Gaussian mixture data, respectively. Our setting is closer to [GMMM20], however, they do not provide optimization guarantees through gradient-based training. Furthermore, in a companion work [BES<sup>+</sup>23], we zoom into the setting where the spike and target are perfectly aligned ( $r_1 = 0$ ), and prove learnability in the  $n \asymp d$  regime for both kernel regression and two-layer neural network. Finally, we go over some results concurrent to our work in Appendix A.

**Learning Single Index Models.** The problem of estimating the relevant direction in a single index model is classical in statistics [LD89], with efficient dedicated algorithms ([KSK11, CM20] among others). However, these algorithms are non-standard and instead, we are concerned with standard iterative algorithms like training neural networks with gradient descent. Recently, [DH18] considered an iterative optimization procedure for learning such models with a polynomial sample complexity that is controlled by the smoothness of the link function. [BES<sup>+</sup>22] considered the effect of taking a single gradient step on the ability of a two-layer neural network to learn a single index model, and [BBSS22, MHPG<sup>+</sup>23] considered training a special two-layer neural network architecture where all neurons share the same weight with gradient flow or online SGD. However, these works only consider the isotropic Gaussian input, and the effect of anisotropy in the covariance matrix when training a neural network to learn a single index model has remained unclear.

**Training a Single Neuron with Gradient Descent.** When training the first layer, we consider a setting where there is only one effective neuron. A large body of works exists on training a single neuron using variants of gradient descent. In the realizable setting (i.e. identical link and activation), the typical assumptions on the activation correspond to information exponent 1 as the activations are required to be monotone or have similar properties, see e.g. [Sol17, YO20, DKTZ22]. In the agnostic setting, [FCG20] considered initializing from the origin which is a saddle point for information exponent larger than 1. [ATV22] also considered the agnostic learning of a ReLU activation, albeit their sample complexity is not explicit other than being polynomial in dimension.

## 2 Preliminaries: Statistical Model and Training Procedure

For a  $d$ -dimensional input  $\mathbf{x}$  and a link function  $g \in L^2(\gamma)$ , consider the single index model

$$y = g\left(\frac{\langle \mathbf{u}, \mathbf{x} \rangle}{\|\Sigma^{1/2} \mathbf{u}\|}\right) + \epsilon \quad \text{with } \mathbf{x} \sim \mathcal{N}(0, \Sigma), \quad (2.1)$$

where  $\epsilon$  is a zero-mean noise with  $\mathcal{O}(1)$  sub-Gaussian norm and  $\mathbf{u} \in \mathbb{S}^{d-1}$ . Learning the model (2.1) corresponds to approximately recovering the unknown link  $g$  and the unknown direction  $\mathbf{u}$ . Note that a normalization is needed to make this problem well-defined; without loss of generality, we write  $\langle \mathbf{u}, \mathbf{x} \rangle / \|\Sigma^{1/2} \mathbf{u}\|$  to ensure that the input variance and the scaling of  $g$  both remain independent of the conditioning of  $\Sigma$ . For this learning task, we will use a two-layer neural network of the form

$$\hat{y}(\mathbf{x}; \mathbf{W}, \mathbf{a}, \mathbf{b}) := \sum_{i=1}^m a_i \phi(\langle \mathbf{w}_i, \mathbf{x} \rangle + b_i), \quad (2.2)$$

where  $\mathbf{W} = \{\mathbf{w}_i\}_{i=1}^m$  is the  $m \times d$  matrix whose rows corresponds to first-layer weights  $\mathbf{w}_i$ ,  $\mathbf{a} = \{a_i\}_{i=1}^m$  denote the second-layer weights,  $\mathbf{b} = \{b_i\}_{i=1}^m$  denote the biases, and  $\phi$  is the non-linear

activation function. We assume  $g$  and  $\phi$  are weakly differentiable with weak derivatives  $g'$  and  $\phi'$  respectively, and  $g, g', \phi, \phi' \in L^2(\gamma)$ . We are interested in the high-dimensional regime; thus,  $d$  is assumed to be sufficiently large throughout the paper. Our ultimate goal is to learn both unknowns  $g$  and  $\mathbf{u}$  by minimizing the population risk

$$R(\mathbf{W}, \mathbf{a}, \mathbf{b}) := \frac{1}{2} \mathbb{E}[(\hat{y}(\mathbf{x}; \mathbf{W}, \mathbf{a}, \mathbf{b}) - y)^2], \quad (2.3)$$

using a gradient-based training method such as gradient flow.

We follow the two-step training procedure employed in recent works [BES<sup>+</sup>22, MHPG<sup>+</sup>23, BBSS22, DLS22]: First, we train the first-layer weights  $\mathbf{W}$  to learn the unknown direction  $\mathbf{u}$ ; at the end of this stage, the neurons  $w_i$  align with  $\mathbf{u}$ . Here, the goal is to recover only the direction. Next, using random biases and training the second-layer weights, we obtain a good approximation for the unknown link function  $g$ . In the majority of this work, we focus on the first part of this two-stage procedure as the alignment between  $w_i$ 's and  $\mathbf{u}$  essentially determines the sample complexity of the overall procedure. This problem is somewhat equivalent to the simplified problem of minimizing (2.3) with  $m = 1$ ,  $a_1 = 1$ ,  $b_1 = 0$ , i.e.,  $\hat{y}(\mathbf{x}; \mathbf{W}, \mathbf{a}, \mathbf{b})$  is replaced with  $\hat{y}(\mathbf{x}; \mathbf{w}) := \phi(\langle \mathbf{w}, \mathbf{x} \rangle)$  and we write  $R(\mathbf{w}) := R(\mathbf{W}, \mathbf{a}, \mathbf{b})$  for simplicity. We emphasize that unless  $\phi = g$  (i.e. the link function is known), the first stage of training only recovers the relevant direction  $\mathbf{u}$  and is not able to approximate  $g$ . Indeed,  $m > 1$  is often needed to learn the non-linear link function; this is the focus of Section 4.2 where we derive a complete learnability result for a two-layer neural network with  $m > 1$ .

Characteristics of the link function play an important role in the complexity of learning the model. As such, a central part of our analysis will rely on a particular property based on the Hermite expansion of functions in a basis defined by the normalized Hermite polynomials  $\{h_j\}_{j \geq 0}$  given as

$$h_j(z) = \frac{(-1)^j e^{z^2/2}}{\sqrt{j!}} \frac{d^j}{dz^j} e^{-z^2/2}. \quad (2.4)$$

These polynomials form an orthonormal basis in the space  $L^2(\gamma)$ , and the resulting expansion yields the following measure of complexity for  $g$ , which is termed as *the information exponent*.

**Definition 1** (Information exponent). *Let  $g = \sum_{j \geq 0} \alpha_j h_j$  be the Hermite expansion of  $g$ . The information exponent of  $g$  is defined to be  $s := \inf\{j > 0 : \alpha_j \neq 0\}$ .*

This concept was introduced in [BAGJ21] in a more general framework, and our definition is more in line with the setting in [BBSS22]. We remark that the definition of [BAGJ21] can be modified to handle anisotropy in which case one arrives at Definition 1. We provide a detailed discussion on this concept together with some properties of the Hermite expansion in Appendix B. Throughout the paper, we assume that the information exponent does not grow with dimension.

In the case where the  $d$ -dimensional input data is isotropic, [BBSS22] showed that learning a single index target with full-batch gradient flow requires a sample complexity of  $\tilde{O}(d^s)$  for  $s \geq 3$  where  $s$  is the information exponent of  $g$ . We will show that this sample complexity can be improved under anisotropy. More specifically, if the input covariance  $\Sigma$  has non-trivial alignment with the unknown direction  $\mathbf{u}$ , we prove in Section 3 that the resulting sample complexity can be even made independent of the information exponent if we use a certain normalization in the training. In what follows, we prove that such a normalization in training procedure is indeed necessary.

## 2.1 The spiked model and limitations of covariance-agnostic training

In practice, data often exhibit a certain structure which may have a profound impact on the statistical procedure. A well-known model that captures such a structure is the *spiked model* [Joh01] for which one or several large eigenvalues of the input covariance matrix  $\Sigma$  are separated from the bulk of the spectrum (see also [BBAP05, BS06]). Although our results hold for generic covariance matrices, they reveal interesting phenomena under the following spiked model assumption.

**Assumption 1.** *The covariance  $\Sigma$  follows the  $(\kappa, \boldsymbol{\theta})$ -spiked model if  $\Sigma = \frac{\mathbf{I}_d + \kappa \boldsymbol{\theta} \boldsymbol{\theta}^\top}{1 + \kappa}$  where  $\|\boldsymbol{\theta}\| = 1$ .*

In pursuit of the target (unit) direction  $\mathbf{u}$ , the magnitude of the neuron  $\mathbf{w}$  is immaterial; thus, recent works take advantage of this and simplify the optimization trajectory by projecting  $\mathbf{w}$  onto unit sphere  $\mathbb{S}^{d-1}$  throughout the training process [BAGJ21, BBSS22]. In the sequel, we study the same dynamics

which is agnostic to the input covariance in order to motivate our investigation of normalized gradient flow in Section 3. More specifically, we consider the spherical population gradient flow

$$\frac{d\mathbf{w}^t}{dt} = -\nabla^S R(\mathbf{w}^t) \quad \text{where} \quad \nabla^S R(\mathbf{w}) = \nabla R(\mathbf{w}) - \langle \nabla R(\mathbf{w}), \mathbf{w} \rangle \mathbf{w}. \quad (2.5)$$

where  $\nabla^S$  is the spherical gradient at the current iterate. It is straightforward to see that when the initialization  $\mathbf{w}^0$  is on the unit sphere, the entire flow will remain on the unit sphere, i.e.  $\mathbf{w}^t \in \mathbb{S}^{d-1}$  for all  $t \geq 0$ . The flow (2.5) has been proven useful for learning the direction  $\mathbf{u}$  [BBSS22] in the isotropic case  $\Sigma = \mathbf{I}_d$  when the activation  $\phi$  is ReLU. In contrast, when  $\Sigma$  follows a spiked model, we show that it can get stuck at stationary points that are almost orthogonal to  $\mathbf{u}$ . Indeed, when the input covariance  $\Sigma$  has a spike in the target direction  $\mathbf{u}$ , i.e.  $\boldsymbol{\theta} = \mathbf{u}$ , one expects that the training procedure benefits from this as the input  $\mathbf{x}$  contains information about the sought unknown  $\mathbf{u}$  without even querying the response  $y$ . The following result proves the contrary; for moderate spike magnitudes, the alignment between the first-layer weights and target  $\langle \mathbf{w}^t, \mathbf{u} \rangle$  will be insignificant for all  $t$ .

**Theorem 2.** *Let  $s > 2$  be the information exponent of  $g$  with  $\mathbb{E}[g] = 0$ , and assume  $\Sigma$  follows the  $(\kappa, \mathbf{u})$ -spiked model with  $\Omega(1) \leq \kappa \leq \mathcal{O}(d^{\frac{s-2}{s-1}})$ . For ReLU activation, let  $\mathbf{w}^t$  denote the solution to (2.5) initialized uniformly at random over  $\mathbb{S}^{d-1}$ , then with probability at least 0.99,*

$$\sup_{t \geq 0} |\langle \mathbf{w}^t, \mathbf{u} \rangle| \lesssim 1/\sqrt{d}, \quad (2.6)$$

A non-trivial alignment between the first-layer weights  $\mathbf{w}^t$  and the target direction  $\mathbf{u}$  is required to learn the single index model (2.1). However, the above result implies that in high dimensions when  $d \gg 1$ , the alignment is negligible in the population limit (when the number of samples goes to infinity). We remark that when the spike magnitude is large, i.e.  $\kappa \geq \Omega(d)$ , the flow (2.5) can achieve alignment as the problem essentially becomes one-dimensional, as we demonstrate in Appendix C.

To see why the flow (2.5) gets stuck at saddle points and fails to recover the true direction, notice that

$$R(\mathbf{w}) = \frac{1}{2} \mathbb{E} \left[ (\phi(\langle \mathbf{w}, \mathbf{x} \rangle) - y)^2 \right] = \frac{1}{2} \mathbb{E} [\phi(\langle \mathbf{w}, \mathbf{x} \rangle)^2] - \mathbb{E} [\phi(\langle \mathbf{w}, \mathbf{x} \rangle)y] + \frac{1}{2} \mathbb{E} [y^2]. \quad (2.7)$$

If the input was isotropic, i.e.  $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$ , the first term in (2.7) would be equal to  $\|\phi\|_\gamma^2$ , which is independent of  $\mathbf{w}$ . Thus, minimizing  $R(\mathbf{w})$  in this case is equivalent to maximizing the ‘‘correlation’’ term  $\mathbb{E}[\phi(\langle \mathbf{w}, \mathbf{x} \rangle)y]$ . However, under the spiked model, the alignment between  $\mathbf{w}$  and  $\mathbf{u}$  breaks the symmetry; consequently, the first term in the decomposition grows with  $\langle \mathbf{w}, \mathbf{u} \rangle$ , creating a repulsive force that traps the dynamics around the equator where  $\mathbf{w}$  is almost orthogonal to  $\mathbf{u}$ .

### 3 Main Results: Alignment via Normalized Dynamics

Having established that the covariance-agnostic training dynamics (2.5) is likely to fail, we consider a covariance-aware normalized flow in this section and show that it can achieve alignment with the unknown target and enjoy better sample complexity compared to the existing results [BAGJ21, BBSS22] in the isotropic case. We start with the population dynamics.

#### 3.1 Warm-up: Population dynamics

To simplify the exposition, we define  $\mathbf{z} := \Sigma^{-1/2} \mathbf{x}$ ,  $\bar{\mathbf{w}} := \Sigma^{1/2} \mathbf{w} / \|\Sigma^{1/2} \mathbf{w}\|$  and similarly define  $\bar{\mathbf{u}}$ , and consider the prediction function  $\hat{y}(\mathbf{x}; \bar{\mathbf{w}}) := \phi(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle)$ . Due to symmetry, the second moment of the prediction is  $\mathbb{E}[\hat{y}(\mathbf{x}; \bar{\mathbf{w}})^2] = \|\phi\|_\gamma^2$  which is independent of  $\mathbf{w}$ ; thus, the population risk reads

$$\mathcal{R}(\mathbf{w}) := \frac{1}{2} \mathbb{E} \left[ (\hat{y}(\mathbf{x}; \bar{\mathbf{w}}) - y)^2 \right] = \frac{1}{2} \|\phi\|_\gamma^2 + \frac{1}{2} \mathbb{E} [y^2] - \mathbb{E} [\phi(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle)y]. \quad (3.1)$$

In (3.1), the only term that depends on the weights  $\mathbf{w}$  is the correlation term and the source of the repulsive force in (2.7) is eliminated; we have  $\nabla_{\mathbf{w}} \mathcal{R}(\mathbf{w}) = -\nabla_{\mathbf{w}} \mathbb{E} [\phi(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle)y]$ . Based on this, we use the following normalized gradient flow for training

$$\frac{d\mathbf{w}^t}{dt} = -\eta(\mathbf{w}^t) \nabla_{\mathbf{w}} \mathcal{R}(\mathbf{w}^t) \quad \text{where} \quad \eta(\mathbf{w}) = \|\Sigma^{1/2} \mathbf{w}\|^2. \quad (3.2)$$

We remark that, though not identical, this normalization is closely related to batch normalization which is commonly employed in practice [IS15]. Under the invariance provided by the current normalization, minimizing  $\mathcal{R}(\mathbf{w})$  corresponds to maximizing  $\mathbb{E}[\phi(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle)y]$ . Thus, instead of  $\mathbf{w}$ , it will be more useful to track the dynamics of its normalized counterpart  $\bar{\mathbf{w}}$ , which is made possible by the following intermediary result that follows from Stein's lemma; also see e.g. [EDB16, MHPG<sup>+</sup>23].

**Lemma 3.** *Suppose we train  $\mathbf{w}^t$  using the gradient flow (3.2). Then  $\bar{\mathbf{w}}^t$  solves the following ODE*

$$\frac{d\bar{\mathbf{w}}^t}{dt} = -\zeta_{\phi,g}(\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle)(\mathbf{I}_d - \bar{\mathbf{w}}^t \bar{\mathbf{w}}^{t\top})\Sigma(\mathbf{I}_d - \bar{\mathbf{w}}^t \bar{\mathbf{w}}^{t\top})\bar{\mathbf{u}}, \quad (3.3)$$

where  $\zeta_{\phi,g}(\langle \bar{\mathbf{w}}, \bar{\mathbf{u}} \rangle) := -\mathbb{E}[\phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle)g'(\langle \bar{\mathbf{u}}, \mathbf{z} \rangle)]$ .

We will investigate if the modified flow (3.3) achieves alignment; in this context, alignment corresponds to  $\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle \approx 1$ . Towards that end, we make the following assumption.

**Assumption 2.** *Let  $g = \sum_{j \geq 0} \alpha_j h_j$  and  $\phi = \sum_{j \geq 0} \beta_j h_j$  be the Hermite decomposition of  $g$  and  $\phi$  respectively. Let  $s$  be the information exponent of  $g$ . For some universal constant  $c > 0$ , we assume*

$$\zeta_{\phi,g}(\omega) = -\sum_{j > 0} j \alpha_j \beta_j \omega^{j-1} \leq -c \omega^{s-1}, \quad \forall \omega \in (0, 1).$$

There are several important examples that readily satisfy Assumption 2. The obvious example is when the link function is known as in [BAGJ21], i.e.  $\phi = g$ . A more interesting example is when  $\phi$  is an activation with degree  $s$  non-zero Hermite coefficient (e.g. ReLU when  $s$  is even, see [GKK19, Claim 1]) and  $g$  is a degree  $s$  Hermite polynomial, which for  $s = 2$  corresponds to the phase retrieval problem. In this case, the assumption is satisfied if  $\alpha_s$  and  $\beta_s$  have the same sign, which occurs with probability 0.5 if we randomly choose the sign of the second layer.

Under this condition, the following result shows that the population flow (3.3) can achieve alignment.

**Proposition 4.** *Suppose Assumption 2 holds and consider the gradient flow given by (3.3) with initialization satisfying  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle > 0$ . Then, we have  $\langle \bar{\mathbf{w}}^T, \bar{\mathbf{u}} \rangle \geq 1 - \varepsilon$  as soon as*

$$T \asymp \frac{\tau_s(\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle) + \ln(1/\varepsilon)}{\lambda_{\min}(\Sigma)} \quad \text{where } \tau_s(z) := \begin{cases} 1 & s = 1 \\ \ln(1/z) & s = 2. \\ (1/z)^{s-2} & s > 2 \end{cases} \quad (3.4)$$

We remark that the information exponent enters the rate in (3.4) through the function  $\tau_s$ , and time needed to achieve  $\varepsilon$  alignment gets worse with larger information exponent. Indeed, it is understood that this quantity serves as a measure of complexity for the target function being learned.

### 3.2 Empirical dynamics and sample complexity

Given  $n$  i.i.d. samples  $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$  from the single index model (2.1), we consider the flow

$$\frac{d\mathbf{w}^t}{dt} = -\eta(\mathbf{w}^t)\nabla\hat{\mathcal{R}}(\mathbf{w}^t) \quad \text{with } \nabla\hat{\mathcal{R}}(\mathbf{w}) := -\nabla_{\mathbf{w}} \left\{ \frac{1}{n} \sum_{i=1}^n \phi \left( \frac{\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle}{\|\hat{\Sigma}^{1/2}\mathbf{w}\|} \right) y^{(i)} \right\}, \quad (3.5)$$

where we estimate the covariance matrix  $\Sigma$  using the sample mean  $\hat{\Sigma} := \frac{1}{n'} \sum_{i=1}^{n'} \mathbf{x}^{(i)} \mathbf{x}^{(i)\top}$  over  $n'$  i.i.d. samples; the above dynamics defines an empirical gradient flow. Notice that we ignored the gradient associated with the term  $\phi^2$  since the population dynamics ensures that its gradient will concentrate around zero; thus, it is redundant to estimate this term. Below, we will use  $n' = n$  for smooth activations, i.e. the same dataset can be used for covariance estimation; For ReLU, we require a more accurate covariance estimator, thus, we use  $n' \gtrsim n^2$  by assuming access to an additional  $n' - n$  unlabeled data points. Similar to the previous section, we track the dynamics of normalized  $\mathbf{w}$  by defining  $\bar{\mathbf{w}} := \hat{\Sigma}^{1/2}\mathbf{w}/\|\hat{\Sigma}^{1/2}\mathbf{w}\|$  (and leave  $\bar{\mathbf{u}}$  unchanged from Section 3.1). The same arguments as in Lemma 3 allow us to track the evolution of  $\bar{\mathbf{w}}$ , which ultimately yields the following alignment result under general covariance structure.

**Theorem 5.** *Let  $s$  be the information exponent of  $g$ , and assume it satisfies  $|g(\cdot)| \lesssim 1 + |\cdot|^p$  for some  $p > 0$ . For  $\phi$  denoting either the ReLU activation or a smooth activation satisfying  $|\phi'| \vee |\phi''| \lesssim 1$ , suppose Assumption 2 holds. For any  $\varepsilon > 0$ , suppose we run the finite sample gradient flow (3.5) with  $\eta(\mathbf{w}) = \|\hat{\Sigma}^{1/2}\mathbf{w}\|^2$ , initialized such that  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle > 0$ , and with number of samples*

$$n \gtrsim d\kappa(\Sigma)^2 \left\{ \langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle^{2(1-s)} \vee \varepsilon^{-2} \right\},$$

where  $\kappa(\Sigma)$  is the condition number of  $\Sigma$ . Then, for  $T \asymp \frac{\tau_s(\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle) + \ln(1/\varepsilon)}{\lambda_{\min}(\Sigma)}$ , we have

$$\langle \bar{\mathbf{w}}^T, \bar{\mathbf{u}} \rangle \geq 1 - \varepsilon, \quad (3.6)$$

with probability at least  $1 - c_1 d^{-c_2}$  for some universal constants  $c_1, c_2 > 0$  over the randomness of the dataset. Here,  $\tau_s$  is defined in (3.4) and  $\gtrsim$  hides poly-logarithmic factors.

**Remark.** We make the following remarks on the above theorem.

- The initial condition  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle > 0$  is required when we have odd information exponent. When  $\mathbf{w}^0$  is initialized uniformly over  $\mathbb{S}^{d-1}$ , the condition holds with probability 0.5 over the initialization. See [BAGJ21, Remark 1.8] for further discussion on this condition.
- Although  $\bar{\mathbf{w}}$  is defined using the empirical covariance unlike  $\bar{\mathbf{u}}$  which is defined by population covariance, this definition is the suitable choice to approximate the target function  $g$  (c.f. Theorem 9), since it ensures the arguments of  $\phi$  and  $g$  are sufficiently close when  $\bar{\mathbf{w}}$  recovers  $\bar{\mathbf{u}}$ .

The intuition behind the proof of Theorem 5 is presented in Section 5 with the complete proof in the appendix. We highlight that the improvement in the sample complexity compared to the isotropic setting occurs whenever the covariance structure induces a stronger initial alignment and consequently stronger signal. The following corollary demonstrates a concrete example of such improvement by specializing Theorem 5 for a spiked covariance model.

**Corollary 6.** Consider the setting of Theorem 5 with  $\Sigma$  following the  $(\kappa, \theta)$ -spiked model, where  $\langle \mathbf{u}, \theta \rangle \asymp d^{-r_1}$  and  $\kappa \asymp d^{r_2}$  with  $r_1 \in [0, 1/2]$  and  $r_2 \in [0, 1]$ . Suppose  $\mathbf{w}^0$  is sampled uniformly from  $\mathbb{S}^{d-1}$ . Then, when conditioned on  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle > 0$ , the sample complexity in Theorem 5 reads

$$n \gtrsim \begin{cases} d^{1+2r_2} (d^{s-1} \vee \varepsilon^{-2}) & 0 < r_2 < r_1 \\ d^{1+2r_2} (d^{(s-1)(1-2(r_2-r_1))} \vee \varepsilon^{-2}) & r_1 < r_2 < 2r_1, \\ d^{1+2r_2} (d^{(s-1)(1-r_2)} \vee \varepsilon^{-2}) & 2r_1 < r_2 < 1 \end{cases} \quad (3.7)$$

where  $\gtrsim$  hides poly-logarithmic factors of  $d$ .

**Remark.** We have the following observations on the above sample complexity.

- Corollary 6 demonstrates that structured data can lead to better sample complexity when the right normalization is used during training. This complements Theorem 2 where we recall that spherical training dynamics ignores the structure in data and the target direction cannot be recovered.
- When  $g$  is a polynomial of degree  $p$ , the lower bound for rotationally invariant kernels (including the neural tangent kernel at initialization) implies a complexity of at least  $d^{\Omega((1-r_2)p)}$  [DWY21]. Thus the sample complexity of Corollary 6 can always outperform the kernel lower bound when  $p$  is sufficiently large and  $s$  remains constant.

**Three-step phase transition.** Recall that in the isotropic setting  $\Sigma = \mathbf{I}_d$ , the sample complexity of learning  $g$  with information exponent  $s$  using full-batch gradient flow is  $\tilde{O}(d^s)$  for  $s \geq 3$  [BBSS22]. The sample complexity in Corollary 6 is strictly smaller than  $\tilde{O}(d^s)$  as soon as  $(s-1)r_1/(s-2) < r_2$ . Furthermore, for any  $\nu > 0$  it is at most  $\tilde{O}(d^{3+\nu})$  as soon as  $r_2 \geq 1 - \nu/(s-3)$  and  $2r_1 < r_2$ , in which case the sample complexity becomes independent of the information exponent. Interestingly, the complexity becomes independent of  $r_1$  when  $r_2 > 2r_1$  or  $r_2 < r_1$ , i.e. the direction of the spike becomes irrelevant when the spike magnitude is sufficiently large or small.

The three-stage phase transition of Corollary 6 is due to the different behaviour of the inner product  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle$  in different regimes of  $r_1$  and  $r_2$ . When  $r_2 < r_1$ , we have  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle \asymp \langle \mathbf{w}^0, \mathbf{u} \rangle$ , thus the initial alignment is just as uninformative as the isotropic case providing no improvement. Moreover, a potentially large condition number may hurt the sample complexity in this case. On the other hand, when  $r_1 < r_2 < 2r_1$  we have  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle \asymp \kappa \langle \mathbf{u}, \theta \rangle \langle \mathbf{w}^0, \theta \rangle$ , and  $r_2 > 2r_1$  leads to  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle \asymp \sqrt{\kappa} \langle \mathbf{w}^0, \theta \rangle$ , thus large  $\kappa$  or  $\langle \mathbf{u}, \theta \rangle$  in this regime may improve the sample complexity.

## 4 Implications to Neural Networks and Further Improvements

### 4.1 Improving Sample Complexity via Preconditioning

We now demonstrate that preconditioning the training dynamics with  $\hat{\Sigma}^{-1}$  can remove the dependency on  $\varkappa(\Sigma)$ , ultimately improving the sample complexity. Consider the preconditioned gradient flow

$$\frac{d\mathbf{w}^t}{dt} = -\eta(\mathbf{w}^t)\hat{\Sigma}^{-1}\nabla\hat{\mathcal{R}}(\mathbf{w}^t) \quad \text{with} \quad \eta(\mathbf{w}) = \|\hat{\Sigma}^{1/2}\mathbf{w}\|^2. \quad (4.1)$$

We have the following alignment result.

**Theorem 7.** *Consider the same setting as Theorem 5, and assume we run the preconditioned empirical gradient flow (4.1) with number of samples*

$$n \gtrsim d \left\{ \langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle^{2(1-s)} \vee \varepsilon^{-2} \right\},$$

where  $\gtrsim$  hides poly-logarithmic factors of  $d$ . Then, for  $T \asymp \tau_s(\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle) + \ln(1/\varepsilon)$ , we have

$$\langle \bar{\mathbf{w}}^T, \bar{\mathbf{u}} \rangle \geq 1 - \varepsilon,$$

with probability at least  $1 - c_1 d^{-c_2}$  for some universal constants  $c_1, c_2 > 0$ .

Preconditioning removes the condition number dependence, which is particularly important in the spiked model case where this quantity can be large.

**Corollary 8.** *Consider the setting of Theorem 7, and assume we run the preconditioned empirical gradient flow (4.1) for the  $(\kappa, \theta)$ -spiked model where  $\langle \mathbf{u}, \theta \rangle \asymp d^{-r_1}$  and  $\kappa \asymp d^{r_2}$  with  $r_1 \in [0, 1/2]$  and  $r_2 \in [0, 1]$ . Suppose  $\mathbf{w}^0$  is sampled uniformly from  $\mathbb{S}^{d-1}$ . Then, when conditioned on  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle > 0$ , the sample complexity of Theorem 7 reads*

$$n \gtrsim \begin{cases} d(d^{s-1} \vee \varepsilon^{-2}) & 0 < r_2 < r_1 \\ d(d^{(s-1)(1-2(r_2-r_1))} \vee \varepsilon^{-2}) & r_1 < r_2 < 2r_1, \\ d(d^{(s-1)(1-r_2)} \vee \varepsilon^{-2}) & 2r_1 < r_2 < 1 \end{cases} \quad (4.2)$$

where  $\gtrsim$  hides poly-logarithmic factors of  $d$ .

The above result improves upon Corollary 6; thus, making a case for preconditioning in practice. The complexity results also strictly improve upon the  $\tilde{\mathcal{O}}(d^s)$  complexity in the isotropic case [BSS22] when  $r_2 > r_1$ . Further, for any  $\nu > 0$ , we can obtain the complexity of  $\tilde{\mathcal{O}}(d^{1+\nu})$  (nearly linear in dimension) when  $r_2 > 1 - \nu/(s-1)$  and  $r_2 > 2r_1$  or  $r_1 + 1/2(1 - \nu/(s-1)) < r_2 < 2r_1$ . In addition to the remarks of Corollary 6, we note that the complexity is independent of both  $r_1$  and  $r_2$  when  $r_2 < r_1$  (cf. Figure 1 hard regime), i.e. the spike magnitude and the spike-target alignment have no effect on the complexity unless  $r_2 \geq r_1$ .

Under the spiked covariance model, one could improve the above results by instead using spectral initialization, i.e. initializing at  $\theta$ , which can be estimated from unlabeled data. Assuming perfect access to  $\theta$ , using the statement of Theorems 5 and 7, this initialization would imply a sample complexity of  $\tilde{\mathcal{O}}(d^{1+2r_2+((s-1)(2r_1-r_2)\vee 0)})$  without and  $\tilde{\mathcal{O}}(d^{1+((s-1)(2r_1-r_2)\vee 0)})$  with preconditioning.

### 4.2 Two-layer neural networks and learning the link function

Our main focus so far was learning the target direction  $\mathbf{u}$ . Next, we consider learning the unknown link function with a neural network, providing a complete learnability result for single index models.

We use Algorithm 1 and train the first-layer of the neural network with either the empirical gradient flow (3.5) or the preconditioned version (4.1). Then, we randomly choose the bias units and minimize the second layer weights using another gradient flow. Our goal is to track the sample complexity  $n$  needed to learn the single index target which we compare against the results of [BSS22]. We highlight that layer-wise training in Algorithm 1 is frequently employed in the literature [BES<sup>+</sup>22, BSS22, DLS22, MHPG<sup>+</sup>23] and in particular [BSS22] also used gradient flow for training.

---

**Algorithm 1** Layer-wise training of a two-layer ReLU network with gradient flow (GF).

---

**Input:**  $\mathbf{w}^0 \in \mathbb{R}^d, T, T', \Delta, \lambda \in \mathbb{R}_+$  and data  $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$ .

- 1: Train the first layer weights  $\mathbf{W}_j^T$  using the GF (3.5) or the preconditioned GF (4.1).
- 2: Normalize the weights  $\mathbf{W}_j^T := \mathbf{W}_j^T / \|\hat{\Sigma}^{1/2} \mathbf{W}_j^T\|$  for every  $1 \leq j \leq m$ .
- 3: Let  $b_j \stackrel{i.i.d.}{\sim} \text{Unif}(-\Delta, \Delta)$  and  $a_j^0 = 1/m$  for  $1 \leq j \leq m$ .
- 4: Train the second layer weights  $\mathbf{a}^{T'}$  via the gradient flow

$$\frac{d\mathbf{a}^t}{dt} = -\nabla_{\mathbf{a}} \left\{ \frac{1}{2n} \sum_{i=1}^n (\hat{y}(\mathbf{x}^{(i)}; \mathbf{W}^T, \mathbf{a}^t, \mathbf{b}) - y^{(i)})^2 + \frac{\lambda \|\mathbf{a}^t\|^2}{2} \right\}.$$

- 5: **return**  $(\mathbf{W}^T, \mathbf{a}^{T'}, \mathbf{b})$ .
- 

**Theorem 9.** *Let  $g$  be twice weakly differentiable with information exponent  $s$  and assume  $g''$  has at most polynomial growth. Suppose  $\phi$  is the ReLU activation, Assumption 2 holds and we run Algorithm 1 with  $\mathbf{w}^0$  initialized uniformly over  $\mathbb{S}^{d-1}$ . For any  $\varepsilon > 0$ , let  $n$  and  $T$  be chosen according to Theorem 5 when we run the gradient flow (3.5) and Theorem 7 when we run the preconditioned gradient flow (4.1). Then, for  $\Delta \asymp \sqrt{\ln(nd)}$ , some regime of  $\lambda$  given by (E.3) and sufficiently large  $T'$  given by (E.4), we have*

$$\mathbb{E}_{(\mathbf{x}, y)} \left[ \left( \hat{y}(\mathbf{x}; \mathbf{W}^T, \mathbf{a}^{T'}, \mathbf{b}) - y \right)^2 \right] \leq C_1 \mathbb{E}[\varepsilon^2] + C_2(\varepsilon + 1/m), \quad (4.3)$$

conditioned on  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle > 0$  with probability at least 0.99 over the randomness of the dataset, biases, and initialization, where  $C_1$  is a universal constant and  $C_2$  hides polylog( $m, n, d$ ) factors.

The next result immediately follows from the previous theorem together with Corollaries 6 & 8.

**Corollary 10.** *In the setting of Theorem 9, if  $\Sigma$  follows the  $(\kappa, \theta)$ -spiked model, the sample complexity  $n$  is given by (3.7) if we use the empirical gradient flow and (4.2) if we use the preconditioned version.*

We remark that for fixed  $\varepsilon$ , the sample complexity to learn  $g$  in the isotropic case is  $\tilde{O}(d^s)$  [BBSS22]. Under the spiked model, if we assume that  $r_2$  is sufficiently large and  $r_1$  is sufficiently small as discussed in the previous section, Corollary 10 improves this rate to either (3.7) when the empirical gradient flow is used without preconditioning or to (4.2) with preconditioning.

## 5 Technical Overview

In this section, we briefly discuss the key intuitions that lead to the proof of our main results. We first review the case  $\Sigma = \mathbf{I}_d$ , where we have the following decomposition for population loss

$$R(\mathbf{w}) := \frac{1}{2} \mathbb{E}[(\phi(\langle \mathbf{w}, \mathbf{x} \rangle) - y)^2] = \frac{1}{2} \|\phi\|_{\gamma}^2 + \frac{1}{2} \mathbb{E}[y^2] - \mathbb{E}[\phi(\langle \mathbf{w}, \mathbf{x} \rangle)g(\langle \mathbf{u}, \mathbf{x} \rangle)]. \quad (5.1)$$

Notice that the only term contributing to the population gradient is the last term which measures the correlation between  $\phi$  and  $g$ . Following the gradient flow and applying Stein's lemma yields

$$\frac{d\langle \mathbf{w}^t, \mathbf{u} \rangle}{dt} = \mathbb{E}[\phi'(\langle \mathbf{w}^t, \mathbf{x} \rangle)g'(\langle \mathbf{u}, \mathbf{x} \rangle)](1 - \langle \mathbf{w}^t, \mathbf{u} \rangle^2) = (1 - \langle \mathbf{w}^t, \mathbf{u} \rangle^2) \sum_{j \geq s} j \alpha_j \beta_j \langle \mathbf{w}^t, \mathbf{u} \rangle^{j-1},$$

where the second identity follows from the Hermite expansion; see also [EDB16, EBD19]. Assume  $\alpha_s \beta_s > 0$  to ensure that the population dynamics will move towards  $\mathbf{u}$  at least near initialization. When replacing the population gradient with a full-batch gradient, we need the estimation noise to be smaller than the signal existing in the gradient. When  $\langle \mathbf{w}^0, \mathbf{u} \rangle \ll 1$ , this signal is roughly of the order  $\langle \mathbf{w}^0, \mathbf{u} \rangle^{s-1}$ . As the uniform concentration error over  $\mathbb{S}^{d-1}$  scales with  $\sqrt{d/n}$ , we need  $n \asymp d \langle \mathbf{w}^0, \mathbf{u} \rangle^{2(s-1)}$  to ensure the signal remains dominant and  $\mathbf{w}^t$  moves towards  $\mathbf{u}$ . When  $\mathbf{w}^0$  is initialized uniformly over  $\mathbb{S}^{d-1}$  this translates to a sample complexity of  $n \asymp d^s$ , which is indeed obtained by [BBSS22] via similar arguments.

However, the behavior of the spherical dynamics entirely changes when we move to the anisotropic case. Suppose  $\Sigma$  follows a  $(\kappa, \mathbf{u})$ -spiked model and  $\phi$  is ReLU. Using Lemma 12, it is easy to show that with the spherical gradient flow, the alignment obeys the following ODE

$$\frac{d\langle \mathbf{w}^t, \mathbf{u} \rangle}{dt} = \left\{ \mathbb{E}[\phi'(\langle \bar{\mathbf{w}}^t, \mathbf{z} \rangle)g'(\langle \bar{\mathbf{u}}, \mathbf{z} \rangle)] - \frac{\kappa \tilde{\psi}_{\phi, g}(\mathbf{w}^t)}{1 + \kappa} \langle \mathbf{w}^t, \mathbf{u} \rangle \right\} (1 - \langle \mathbf{w}^t, \mathbf{u} \rangle^2),$$

where  $\tilde{\psi}_{\phi, g}(\mathbf{w}^t)$  is introduced in Lemma 12. The additional  $\tilde{\psi}_{\phi, g}(\mathbf{w}^t)$  term creates a repulsive force towards the equator  $\langle \mathbf{w}^t, \mathbf{u} \rangle = 0$ . The presence of this term is due to the fact that unlike (5.1), the term  $\mathbb{E}[\phi(\langle \mathbf{w}^t, \mathbf{u} \rangle)^2]$  is no longer independent of  $\mathbf{w}$  and cannot be replaced by  $\|\phi\|_\gamma^2$ . When  $\mathbf{w}^0$  is initialized uniformly over  $\mathbb{S}^{d-1}$  and  $\Omega(1) \leq \kappa \leq \mathcal{O}(d)$ , we have  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle \asymp \sqrt{\kappa} \langle \mathbf{w}^0, \mathbf{u} \rangle$ . Furthermore, at this initialization  $\tilde{\psi}_{\phi, g}(\mathbf{w}^0) \approx 1/2$ . Therefore,

$$\frac{d\langle \mathbf{w}^t, \mathbf{u} \rangle}{dt} \approx \left\{ s\alpha_s \beta_s (\sqrt{\kappa} \langle \mathbf{w}^0, \mathbf{u} \rangle)^{s-1} - \frac{\langle \mathbf{w}^0, \mathbf{u} \rangle}{2} \right\}.$$

Hence the dynamics is trapped at  $|\langle \mathbf{w}^t, \mathbf{u} \rangle| = \mathcal{O}(1/\sqrt{d})$  for all  $t > 0$  as long as  $\kappa = \mathcal{O}(d^{1-1/(s-1)})$ .

To remove the repulsive force in the spherical dynamics, we can directly normalize the input of  $\phi$ . As demonstrated by (3.1), once again the only term that varies with  $\mathbf{w}$  would be the correlation loss. Specifically, using the result of Lemma 3, in the population limit we can track  $\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle$  via

$$\frac{d\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle}{dt} = \mathbb{E}[\phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle)g'(\langle \bar{\mathbf{u}}, \mathbf{z} \rangle)] \langle \bar{\mathbf{u}}_\perp^t, \Sigma \bar{\mathbf{u}}_\perp^t \rangle, \quad (5.2)$$

where  $\bar{\mathbf{u}}_\perp^t := \bar{\mathbf{u}} - \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle \bar{\mathbf{w}}^t$ . Thus, the strength of the signal at initialization is of order  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle^{s-1} / \varkappa(\Sigma)$ , which after controlling the error in the estimate of  $\hat{\Sigma}$  and in the estimate of population gradient using finitely many samples, leads to the sample complexity  $n \asymp d \varkappa(\Sigma)^2 \langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle^{2(1-s)}$ . Importantly,  $\Sigma$  can include a much stronger initial alignment  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle$  than the isotropic case  $\langle \mathbf{w}^0, \mathbf{u} \rangle$ , which is emphasized in Corollary 6. Using preconditioning will further remove the dependency on the condition number of  $\Sigma$ .

## 6 Conclusion

We studied the dynamics of gradient flow to learn single index models when the input data covariance may contain additional structure. Under a spiked model for the covariance matrix, we showed that using spherical gradient flow, as an example of a covariance-agnostic training mechanism employed in the recent literature, is unable to learn the target direction of the single index model even when the spike and the target directions are identical. In contrast, we showed that an appropriate weight normalization removes this problem and successfully recovers the target direction. Moreover, depending on the alignment between the covariance structure and the target direction, the sample complexity can improve upon the isotropic setting, while also outperforming lower bounds for rotationally-invariant kernels. This phenomenon is due to the additional information about the target direction contained in the covariance matrix which improves the effective alignment at initialization. Additionally, we showed that a simple preconditioning of the gradient flow using the inverse empirical covariance can improve the sample complexity, achieving almost linear rate in certain settings.

We outline a few limitations of our current work and discuss directions for future research.

- While studying single index models provides a pathway to a general understanding of feature learning with structured covariance, considering multi-index models can provide a more complete picture [PSE22], e.g. by establishing incremental learning dynamics [ABAM23]. We leave the problem of learning multi-index models under structured input as an interesting future direction.
- Gradient flow under squared loss can be seen as an example of a Correlational Statistical Query (CSQ) algorithm [BF02, Rey20], i.e. an algorithm that only accesses noisy estimates of expected correlation queries from the model. Understanding the limitations of learning single index models under a structured input through a CSQ lower bound perspective is another important direction that would complement our results in this paper.
- When training the first layer, we considered a somewhat unconventional initialization and relied on the symmetry it induces. It is interesting to consider cases where we train a network with multiple neurons starting from a more standard initialization which can help relax Assumption 2.

## Acknowledgments

The authors thank Alberto Bietti and Zhichao Wang for discussions and feedback on the manuscript. TS was partially supported by JSPS KAKENHI (20H00576) and JST CREST (JPMJCR2015). MAE was partially supported by NSERC Grant [2019-06167], CIFAR AI Chairs program, CIFAR AI Catalyst grant.

## References

- [AAM22] Emmanuel Abbe, Enric Boix Adsera, and Theodor Misiakiewicz, *The merged-staircase property: a necessary and nearly sufficient condition for sgd learning of sparse functions on two-layer neural networks*, Conference on Learning Theory, 2022.
- [ABAM23] Emmanuel Abbe, Enric Boix-Adsera, and Theodor Misiakiewicz, *Sgd learning on neural networks: leap complexity and saddle-to-saddle dynamics*, arXiv preprint arXiv:2302.11055 (2023).
- [ASKL23] Luca Arnaboldi, Ludovic Stephan, Florent Krzakala, and Bruno Loureiro, *From high-dimensional & mean-field dynamics to dimensionless odes: A unifying approach to sgd in two-layers networks*, arXiv preprint arXiv:2302.05882 (2023).
- [ATV22] Pranjal Awasthi, Alex Tang, and Aravindan Vijayaraghavan, *Agnostic learning of general relu activation using gradient descent*, arXiv preprint arXiv:2208.02711 (2022).
- [BAGJ21] Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath, *Online stochastic gradient descent on non-convex losses from high-dimensional inference.*, J. Mach. Learn. Res. **22** (2021), 106–1.
- [BBAP05] Jinho Baik, Gérard Ben Arous, and Sandrine Péché, *Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices*.
- [BBSS22] Alberto Bietti, Joan Bruna, Clayton Sanford, and Min Jae Song, *Learning single-index models with shallow neural networks*, Advances in Neural Information Processing Systems, 2022.
- [BEG<sup>+</sup>22] Boaz Barak, Benjamin L Edelman, Surbhi Goel, Sham Kakade, Eran Malach, and Cyril Zhang, *Hidden Progress in Deep Learning: SGD Learns Parities Near the Computational Limit*, arXiv preprint arXiv:2207.08799 (2022).
- [BES<sup>+</sup>19] Jimmy Ba, Murat Erdogdu, Taiji Suzuki, Denny Wu, and Tianzong Zhang, *Generalization of two-layer neural networks: An asymptotic viewpoint*, International Conference on Learning Representations, 2019.
- [BES<sup>+</sup>22] Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg Yang, *High-dimensional Asymptotics of Feature Learning: How One Gradient Step Improves the Representation*, arXiv preprint arXiv:2205.01445 (2022).
- [BES<sup>+</sup>23] Jimmy Ba, Murat A. Erdogdu, Taiji Suzuki, Zhichao Wang, and Denny Wu, *Learning in the presence of low-dimensional structure: a spiked random matrix perspective*, Thirty-seventh Conference on Neural Information Processing Systems (NeurIPS 2023), 2023.
- [BF02] Nader H Bshouty and Vitaly Feldman, *On using extended statistical queries to avoid membership queries*, Journal of Machine Learning Research **2** (2002), no. Feb, 359–395.
- [BMZ23] Raphaël Berthier, Andrea Montanari, and Kangjie Zhou, *Learning time-scales in two-layers neural networks*, arXiv preprint arXiv:2303.00055 (2023).
- [BPVZ23] Joan Bruna, Loucas Pillaud-Vivien, and Aaron Zweig, *On single index models beyond gaussian data*, arXiv preprint arXiv:2307.15804 (2023).
- [BS06] Jinho Baik and Jack W Silverstein, *Eigenvalues of large sample covariance matrices of spiked population models*, Journal of multivariate analysis **97** (2006), no. 6, 1382–1408.

- [CB18] Lenaic Chizat and Francis Bach, *On the Global Convergence of Gradient Descent for Over-parameterized Models using Optimal Transport*, Advances in Neural Information Processing Systems, 2018.
- [Chi22] Lénaïc Chizat, *Mean-field langevin dynamics: Exponential convergence and annealing*, arXiv preprint arXiv:2202.01009 (2022).
- [CM20] Sitan Chen and Raghu Meka, *Learning polynomials in few relevant dimensions*, Conference on Learning Theory, 2020.
- [COB19] Lenaic Chizat, Edouard Oyallon, and Francis Bach, *On Lazy Training in Differentiable Programming*, Advances in Neural Information Processing Systems, 2019.
- [CWPPS23] Elizabeth Collins-Woodfin, Courtney Paquette, Elliot Paquette, and Inbar Seroussi, *Hitting the high-dimensional notes: An ode for sgd learning dynamics on glms and multi-index models*, arXiv preprint arXiv:2308.08977 (2023).
- [DH18] Rishabh Dubeja and Daniel Hsu, *Learning single-index models in gaussian space*, Conference On Learning Theory, PMLR, 2018, pp. 1887–1930.
- [DKL<sup>+</sup>23] Yatin Dandi, Florent Krzakala, Bruno Loureiro, Luca Pesce, and Ludovic Stephan, *Learning two-layer neural networks, one (giant) step at a time*, arXiv preprint arXiv:2305.18270 (2023).
- [DKTZ22] Ilias Diakonikolas, Vasilis Kontonis, Christos Tzamos, and Nikos Zarifis, *Learning a single neuron with adversarial label noise via gradient descent*, Conference on Learning Theory, PMLR, 2022, pp. 4313–4361.
- [DLS22] Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi, *Neural Networks can Learn Representations with Gradient Descent*, Conference on Learning Theory, 2022.
- [DNGL23] Alex Damian, Eshaan Nichani, Rong Ge, and Jason D Lee, *Smoothing the landscape boosts the signal for sgd: Optimal sample complexity for learning single index models*, arXiv preprint arXiv:2305.10633 (2023).
- [DWY21] Konstantin Donhauser, Mingqi Wu, and Fanny Yang, *How rotational invariance of common kernels prevents generalization in high dimensions*, International Conference on Machine Learning, 2021.
- [EBD19] Murat A. Erdogdu, Mohsen Bayati, and Lee H. Dicker, *Scalable Approximations to Generalized Linear Problems*, Journal of Machine Learning Research (2019).
- [EDB16] Murat A Erdogdu, Lee H Dicker, and Mohsen Bayati, *Scaled least squares estimator for glms in large-scale problems*, Advances in Neural Information Processing Systems **29** (2016).
- [Erd15] Murat A Erdogdu, *Newton-stein method: a second order method for glms via stein’s lemma*, Proceedings of Advances in Neural Information Processing Systems, 2015, pp. 1216–1224.
- [FCG20] Spencer Frei, Yuan Cao, and Quanquan Gu, *Agnostic learning of a single neuron with gradient descent*, Advances in Neural Information Processing Systems, vol. 33, Curran Associates, Inc., 2020, pp. 5417–5428.
- [GKK19] Surbhi Goel, Sushrut Karmalkar, and Adam Klivans, *Time/accuracy tradeoffs for learning a relu with respect to gaussian marginals*, Advances in neural information processing systems **32** (2019).
- [GMMM19] B. Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari, *Limitations of Lazy Training of Two-layers Neural Networks*, Advances in Neural Information Processing Systems, 2019.
- [GMMM20] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari, *When Do Neural Networks Outperform Kernel Methods?*, Advances in Neural Information Processing Systems, 2020.

- [HC22] Karl Hajjar and Lenaic Chizat, *On the symmetries in the dynamics of wide two-layer neural networks*, arXiv preprint arXiv:2211.08771 (2022).
- [HTFF09] Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman, *The elements of statistical learning: data mining, inference, and prediction*, vol. 2, Springer, 2009.
- [IS15] Sergey Ioffe and Christian Szegedy, *Batch normalization: Accelerating deep network training by reducing internal covariate shift*, International conference on machine learning, pmlr, 2015, pp. 448–456.
- [JGH18] Arthur Jacot, Franck Gabriel, and Clement Hongler, *Neural Tangent Kernel: Convergence and Generalization in Neural Networks*, Advances in Neural Information Processing Systems, 2018.
- [Joh01] Iain M Johnstone, *On the distribution of the largest eigenvalue in principal components analysis*, The Annals of statistics **29** (2001), no. 2, 295–327.
- [JWHT13] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani, *An introduction to statistical learning*, vol. 112, Springer, 2013.
- [KKSK11] Sham M Kakade, Varun Kanade, Ohad Shamir, and Adam Kalai, *Efficient learning of generalized linear and single index models with isotonic regression*, Advances in Neural Information Processing Systems **24** (2011).
- [LD89] Ker-Chau Li and Naihua Duan, *Regression Analysis Under Link Violation*, The Annals of Statistics (1989).
- [LMZ20] Yuanzhi Li, Tengyu Ma, and Hongyang R Zhang, *Learning over-parametrized two-layer neural networks beyond NTK*, Conference on Learning Theory, 2020.
- [MHD<sup>+</sup>23] Arvind Mahankali, Jeff Z Haochen, Kefan Dong, Margalit Glasgow, and Tengyu Ma, *Beyond ntk with vanilla gradient descent: A mean-field analysis of neural networks with polynomial width, samples, and time*, arXiv preprint arXiv:2306.16361 (2023).
- [MHPG<sup>+</sup>23] Alireza Mousavi-Hosseini, Sejun Park, Manuela Girotti, Ioannis Mitliagkas, and Murat A Erdogdu, *Neural networks efficiently learn low-dimensional representations with SGD*, The Eleventh International Conference on Learning Representations, 2023.
- [MMM19] Song Mei, Theodor Misiakiewicz, and Andrea Montanari, *Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit*, Conference on Learning Theory, 2019.
- [MMN18] Song Mei, Andrea Montanari, and Phan-Minh Nguyen, *A mean field view of the landscape of two-layer neural networks*, Proceedings of the National Academy of Sciences **115** (2018), no. 33, E7665–E7671.
- [NWS22] Atsushi Nitanda, Denny Wu, and Taiji Suzuki, *Convex analysis of the mean field langevin dynamics*, International Conference on Artificial Intelligence and Statistics, PMLR, 2022, pp. 9741–9757.
- [O’D14] Ryan O’Donnell, *Analysis of boolean functions*, Cambridge University Press, 2014.
- [PSE22] Sejun Park, Umut Simsekli, and Murat A. Erdogdu, *Generalization Bounds for Stochastic Gradient Descent via Localized  $\varepsilon$ -Covers*, arXiv preprint arXiv:2209.08951 (2022).
- [Rey20] Lev Reyzin, *Statistical queries and statistical algorithms: Foundations and applications*, arXiv preprint arXiv:2004.00557 (2020).
- [RGKZ21] Maria Refinetti, Sebastian Goldt, Florent Krzakala, and Lenka Zdeborová, *Classifying high-dimensional Gaussian mixtures: Where kernel methods fail and neural networks succeed*, International Conference on Machine Learning, 2021.

- [RVE18] Grant M Rotskoff and Eric Vanden-Eijnden, *Neural networks as Interacting Particle Systems: Asymptotic convexity of the Loss Landscape and Universal Scaling of the Approximation Error*, arXiv preprint arXiv:1805.00915 (2018).
- [Sol17] Mahdi Soltanolkotabi, *Learning relus via gradient descent*, Advances in neural information processing systems **30** (2017).
- [SWON23] Taiji Suzuki, Denny Wu, Kazusato Oko, and Atsushi Nitanda, *Feature learning via mean-field langevin dynamics: classifying sparse parities and beyond*, Thirty-seventh Conference on Neural Information Processing Systems (NeurIPS 2023), 2023.
- [Tel23] Matus Telgarsky, *Feature selection and low test error in shallow low-rotation relu networks*, The Eleventh International Conference on Learning Representations, 2023.
- [Ver18] Roman Vershynin, *High-dimensional probability: An introduction with applications in data science*, Cambridge University Press, 2018.
- [VH16] Ramon Van Handel, *Probability in high dimension*, 2016.
- [Wai19] Martin J. Wainwright, *High-dimensional statistics: A non-asymptotic viewpoint*, Cambridge University Press, 2019.
- [WLLM19] Colin Wei, Jason D Lee, Qiang Liu, and Tengyu Ma, *Regularization matters: Generalization and optimization of neural nets vs their induced kernel*, Advances in Neural Information Processing Systems **32** (2019).
- [YO20] Gilad Yehudai and Shamir Ohad, *Learning a single neuron with gradient methods*, Proceedings of Thirty Third Conference on Learning Theory, Proceedings of Machine Learning Research, vol. 125, PMLR, 2020, pp. 3756–3786.
- [YS19] Gilad Yehudai and Ohad Shamir, *On the Power and Limitations of Random Features for Understanding Neural Networks*, Advances in Neural Information Processing Systems, 2019.

## A Concurrent Works

In this paragraph we briefly summarize a few relevant results concurrent to our submission. [BMZ23] provided a precise analysis of the two-timescale dynamics in learning a single index model with information exponent  $s = 1$ . [MHD<sup>+</sup>23] considered the learning of a single index target with  $s = 2$  using a neural network in the mean-field regime. [BPVZ23] extended the information exponent-based characterization of online SGD to input data beyond Gaussian. [DNGL23] showed that a gradient-smoothed dynamics can improve the sample complexity and match the CSQ lower bound. Finally, beyond the single-index setting, [DKL<sup>+</sup>23, SWON23, CWPPS23] considered learning low-dimensional target functions supported on  $k > 1$  dimensions via gradient-based feature learning.

## B Background on Hermite Expansion

The normalized Hermite polynomials  $\{h_j\}_{j \geq 0}$  given by (2.4) provide an orthonormal basis for  $L^2(\gamma)$ , thus for every  $f \in L^2(\gamma)$  we have

$$f = \sum_{j=0}^{\infty} \langle f, h_j \rangle_{\gamma} h_j,$$

where  $\langle f, h_j \rangle_{\gamma} := \mathbb{E}_{z \sim \mathcal{N}(0,1)}[f(z)h_j(z)]$ . We will commonly invoke the following well-known properties of Hermite polynomials. If  $j \geq 1$ , then  $h'_j = \sqrt{j}h_{j-1}$ , where  $h'_j$  stands for the derivative of  $h_j$ . Furthermore, if  $z_1$  and  $z_2$  are two standard Gaussian random variables with  $\mathbb{E}[z_1 z_2] = \rho$ , then  $\mathbb{E}[h_i(z_1)h_j(z_2)] = \delta_{ij}\rho^j$  where  $\delta_{ij}$  is the Kronecker delta. We refer the interested reader to [O'D14, Chapter 11.2] for additional discussions and properties of these polynomials.

We will now discuss how our Definition 1 relates to the original definition of information exponent of [BAGJ21]. In their setting, they assume the true data distribution  $\mathbb{P}_{\mathbf{u}}$  is parameterized by some unit vector  $\mathbf{u} \in \mathbb{S}^{d-1}$ , and we know the parametric family  $\{\mathbb{P}_{\mathbf{w}}\}_{\mathbf{w} \in \mathbb{S}^{d-1}}$ ; thus the problem is to estimate the direction  $\mathbf{u}$ . Furthermore, they assume the population loss, which is the expectation of some per-sample loss, has spherical symmetry, i.e. the population loss  $R(\mathbf{w})$  can be written as  $R(\mathbf{w}) = \tilde{R}(\langle \mathbf{w}, \mathbf{u} \rangle)$ . Then, [BAGJ21, Definition 1.2] defines the information exponent to be the degree of the first non-zero coefficient of  $\tilde{R}$  in its Taylor expansion around the origin. In other words, we say  $R$  has information exponent  $s$  if

$$\begin{cases} \frac{d^k \tilde{R}}{dz^k}(0) = 0 & 1 \leq k < s \\ \frac{d^k \tilde{R}}{dz^k}(0) = -c < 0 & k = s \\ \left| \frac{d^k \tilde{R}}{dz^k}(z) \right| \leq C & k > s, \forall z \in [-1, 1] \end{cases},$$

where  $C, c > 0$  are universal constants. To specialize the above abstract definition to the Gaussian case, consider the setting where the input data is standard Gaussian  $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$  and the problem is to estimate  $\mathbf{u} \in \mathbb{S}^{d-1}$  given a response variable  $y = f(\langle \mathbf{u}, \mathbf{x} \rangle)$  with known  $f$ . Via the Hermite expansion of  $f$ , one can write

$$\tilde{R}(\langle \mathbf{w}, \mathbf{u} \rangle) := \frac{1}{2} \mathbb{E}[(f(\langle \mathbf{w}, \mathbf{x} \rangle) - f(\langle \mathbf{u}, \mathbf{x} \rangle))^2] = - \sum_{j \geq 1} \langle f, h_j \rangle_{\gamma}^2 \langle \mathbf{w}, \mathbf{u} \rangle^j + \text{const.}$$

Thus, the information exponent of  $\tilde{R}$  is indeed the degree of the first non-zero term in the Hermite expansion of  $f$ .

Now consider the general case where  $\mathbf{x} \sim \mathcal{N}(0, \Sigma)$ . The spherical symmetry assumed in [BAGJ21] no longer holds. However, after proper normalization of weights, if we consider the population loss

$$R(\mathbf{w}) := \frac{1}{2} \mathbb{E} \left[ \left( f \left( \frac{\langle \mathbf{w}, \mathbf{x} \rangle}{\|\Sigma^{1/2} \mathbf{w}\|} \right) - f \left( \frac{\langle \mathbf{u}, \mathbf{x} \rangle}{\|\Sigma^{1/2} \mathbf{u}\|} \right) \right)^2 \right],$$

then  $R(\mathbf{w}) = \tilde{R} \left( \frac{\langle \mathbf{w}, \Sigma \mathbf{u} \rangle}{\|\Sigma^{1/2} \mathbf{w}\| \|\Sigma^{1/2} \mathbf{u}\|} \right)$ . Indeed, a close examination of the arguments of [BAGJ21] reveals that for their results to hold, the proper symmetry to consider is the ellipsoidal symmetry, and

the proper definition of information exponent is the degree of the first non-zero term in the Hermite expansion of  $\tilde{R}$ , which reads

$$\tilde{R}(z) = - \sum_{j \geq 1} \langle f, h_j \rangle^2 z^j + \text{const.}$$

Once again, we can consistently define the information exponent to be the degree of the first non-zero term in the Hermite expansion of  $f$ , as long as the input is Gaussian (potentially anisotropic).

## C Proofs of Section 2.1

Before beginning our main discussions, we state the following lemma which is a generalization of Stein's lemma (Gaussian integration by parts), and will help obtain a closed-form expression for the population gradient. We refer to [Erd15, MHPG<sup>+</sup>23] for similar statements.

**Lemma 11.** *Let  $f, g : \mathbb{R} \rightarrow \mathbb{R}$  with  $g$  weakly differentiable. Suppose  $z \sim \mathcal{N}(0, \mathbf{I}_d)$ . Then, for any  $\mathbf{w}, \mathbf{u} \in \mathbb{S}^{d-1}$ , we have*

$$\mathbb{E}[f(\langle \mathbf{w}, z \rangle)g(\langle \mathbf{u}, z \rangle)z] = \mathbb{E}[f(\langle \mathbf{w}, z \rangle)g'(\langle \mathbf{u}, z \rangle)]\mathbf{u} + \mathbb{E}[f(\langle \mathbf{w}, z \rangle)\{g(\langle \mathbf{u}, z \rangle)\langle \mathbf{w}, z \rangle - g'(\langle \mathbf{u}, z \rangle)\langle \mathbf{u}, \mathbf{w} \rangle\}]\mathbf{w}.$$

**Proof.** Consider the conditional distribution  $z | \langle \mathbf{w}, z \rangle \sim \mathcal{N}(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\Sigma}})$ , where

$$\bar{\boldsymbol{\mu}} = \langle \mathbf{w}, z \rangle \mathbf{w} \quad \text{and} \quad \bar{\boldsymbol{\Sigma}} = \mathbf{I}_d - \mathbf{w} \mathbf{w}^\top.$$

Recall that Stein's lemma (Gaussian integration by parts) states that when  $\bar{z} \sim \mathcal{N}(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\Sigma}})$ , then

$$\mathbb{E}[g(\bar{z})\bar{z}] = \mathbb{E}[g(\bar{z})]\bar{\boldsymbol{\mu}} + \bar{\boldsymbol{\Sigma}} \mathbb{E}[\nabla g(\bar{z})].$$

Hence,

$$\mathbb{E}[g(\langle \mathbf{u}, z \rangle)z | \langle \mathbf{w}, z \rangle] = \mathbb{E}[g(\langle \mathbf{u}, z \rangle) | \langle \mathbf{w}, z \rangle]\langle \mathbf{w}, z \rangle \mathbf{w} + (\mathbf{I}_d - \mathbf{w} \mathbf{w}^\top) \mathbb{E}[g'(\langle \mathbf{u}, z \rangle) | \langle \mathbf{w}, z \rangle] \mathbf{u}.$$

Applying the tower property of conditional expectation and rearranging the terms yields the desired result.  $\square$

We are now ready to state and prove the expression for the population gradient when using the ReLU activation.

**Lemma 12.** *Suppose  $\phi$  is the ReLU activation and  $g$  is weakly differentiable. Let  $z \sim \mathcal{N}(0, \mathbf{I}_d)$ . Define  $\bar{\mathbf{w}} := \frac{\boldsymbol{\Sigma}^{1/2} \mathbf{w}}{\|\boldsymbol{\Sigma}^{1/2} \mathbf{w}\|}$  and similarly define  $\bar{\mathbf{u}}$ . Then*

$$\nabla R(\mathbf{w}) = \boldsymbol{\Sigma} \left\{ \tilde{\psi}_{\phi, g}(\mathbf{w}) \mathbf{w} + \tilde{\zeta}_{\phi, g}(\mathbf{w}) \mathbf{u} \right\}$$

where

$$\tilde{\psi}_{\phi, g}(\mathbf{w}) := \frac{\mathbb{E}[-\phi(\langle \bar{\mathbf{w}}, z \rangle)g(\langle \bar{\mathbf{u}}, z \rangle) + \phi'(\langle \bar{\mathbf{w}}, z \rangle)g'(\langle \bar{\mathbf{u}}, z \rangle)\langle \bar{\mathbf{u}}, \bar{\mathbf{w}} \rangle]}{\|\boldsymbol{\Sigma}^{1/2} \mathbf{w}\|} + \frac{1}{2},$$

and

$$\tilde{\zeta}_{\phi, g}(\mathbf{w}) := - \frac{\mathbb{E}[\phi'(\langle \bar{\mathbf{w}}, z \rangle)g'(\langle \bar{\mathbf{u}}, z \rangle)]}{\|\boldsymbol{\Sigma}^{1/2} \mathbf{u}\|}.$$

**Proof.** Notice that the population risk is given by

$$R(\mathbf{w}) = \frac{1}{2} \mathbb{E} \left[ \left( \phi(\langle \mathbf{w}, \mathbf{x} \rangle) - g \left( \frac{\langle \mathbf{u}, \mathbf{x} \rangle}{\|\boldsymbol{\Sigma}^{1/2} \mathbf{u}\|} \right) \right)^2 \right] + \frac{\mathbb{E}[\epsilon^2]}{2} = \frac{1}{2} \mathbb{E}[\phi(\langle \mathbf{w}, \mathbf{x} \rangle)^2] - \mathbb{E}[\phi(\langle \mathbf{w}, \mathbf{x} \rangle)g] + \frac{\mathbb{E}[y^2]}{2},$$

Notice that  $\mathbf{x} = \boldsymbol{\Sigma}^{1/2} z$ . By the homogeneity of ReLU, we can rewrite the first term as

$$\mathbb{E}[\phi(\langle \mathbf{w}, \mathbf{x} \rangle)^2] = \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w} \|\phi\|_\gamma^2 = \frac{\mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w}}{2}.$$

Then we have,

$$\nabla R(\mathbf{w}) = \frac{\Sigma \mathbf{w}}{2} - \underbrace{\mathbb{E} \left[ \phi(\langle \mathbf{w}, \mathbf{x} \rangle)' g \left( \frac{\langle \mathbf{u}, \mathbf{x} \rangle}{\|\Sigma^{1/2} \mathbf{u}\|} \right) \mathbf{x} \right]}_{=: v(\mathbf{w})}$$

Note that  $\phi'(\langle \mathbf{w}, \mathbf{x} \rangle) = \phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle)$ . Then,

$$\begin{aligned} v(\mathbf{w}) &= \Sigma^{1/2} \mathbb{E}[\phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle)' g(\langle \bar{\mathbf{u}}, \mathbf{z} \rangle) \mathbf{z}] \\ &= \Sigma^{1/2} \{ \mathbb{E}[\phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle)' g(\langle \bar{\mathbf{u}}, \mathbf{z} \rangle)] \bar{\mathbf{u}} + \mathbb{E}[\phi(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) g(\langle \bar{\mathbf{u}}, \mathbf{z} \rangle) - \phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) g'(\langle \bar{\mathbf{u}}, \mathbf{z} \rangle) \langle \bar{\mathbf{u}}, \bar{\mathbf{w}} \rangle] \bar{\mathbf{w}} \} \\ &= \Sigma \left\{ \frac{\mathbb{E}[\phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle)' g'(\langle \bar{\mathbf{u}}, \mathbf{z} \rangle)]}{\|\Sigma^{1/2} \mathbf{u}\|} \mathbf{u} + \frac{\mathbb{E}[\phi(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) g(\langle \bar{\mathbf{u}}, \mathbf{z} \rangle) - \phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) g'(\langle \bar{\mathbf{u}}, \mathbf{z} \rangle) \langle \bar{\mathbf{u}}, \bar{\mathbf{w}} \rangle]}{\|\Sigma^{1/2} \mathbf{w}\|} \mathbf{w} \right\} \end{aligned}$$

where we used Lemma 11 and the fact that  $\phi'(z)z = \phi(z)$ . Therefore,

$$\nabla R(\mathbf{w}) = \Sigma \left\{ \tilde{\psi}_{\phi,g}(\mathbf{w}) \mathbf{w} + \tilde{\zeta}_{\phi,g}(\mathbf{w}) \mathbf{u} \right\}$$

which concludes the proof.  $\square$

Particularly, the above lemma yields the following corollary for the spherical dynamics in the population limit.

**Corollary 13.** *Suppose  $\{\mathbf{w}^t\}_{t \geq 0}$  is a solution to the population spherical gradient flow (2.5),  $\phi$  is the ReLU activation, and  $\Sigma$  follows the  $(\kappa, \boldsymbol{\theta})$ -spiked model. Then,*

$$\begin{aligned} \frac{d\langle \mathbf{w}^t, \mathbf{u} \rangle}{dt} &= - \frac{\tilde{\zeta}_{\phi,g}(\mathbf{w}^t)(1 - \langle \mathbf{w}^t, \mathbf{u} \rangle^2)}{1 + \kappa} \\ &\quad - \frac{\kappa \left\{ \tilde{\psi}_{\phi,g}(\mathbf{w}^t) \langle \mathbf{w}^t, \boldsymbol{\theta} \rangle + \tilde{\zeta}_{\phi,g}(\mathbf{w}^t) \langle \mathbf{u}, \boldsymbol{\theta} \rangle \right\}}{1 + \kappa} (\langle \boldsymbol{\theta}, \mathbf{u} \rangle - \langle \mathbf{w}^t, \boldsymbol{\theta} \rangle \langle \mathbf{w}^t, \mathbf{u} \rangle). \end{aligned} \quad (\text{C.1})$$

### C.1 Proof of Theorem 2

Plugging in  $\boldsymbol{\theta} = \mathbf{u}$  in Corollary 13, we obtain

$$\frac{d\langle \mathbf{w}^t, \mathbf{u} \rangle}{dt} = - \left\{ \tilde{\zeta}_{\phi,g}(\mathbf{w}^t) + \frac{\kappa \tilde{\psi}_{\phi,g}(\mathbf{w}^t) \langle \mathbf{w}^t, \mathbf{u} \rangle}{1 + \kappa} \right\} (1 - \langle \mathbf{u}, \mathbf{w}^t \rangle^2). \quad (\text{C.2})$$

To prove the statement of the theorem, we will show that whenever

$$|\langle \mathbf{w}^t, \mathbf{u} \rangle| \leq \frac{C}{\sqrt{d}},$$

we will have  $\frac{d\langle \mathbf{w}^t, \mathbf{u} \rangle}{dt} < 0$ , thus when initialized from  $\mathcal{O}(1/\sqrt{d})$ ,  $\langle \mathbf{w}^t, \mathbf{u} \rangle$  can never escape the saddle point near the equator  $\langle \mathbf{w}, \mathbf{u} \rangle = 0$ .

Recall from the properties of the Hermite expansion in Appendix B that

$$g' = \sum_{j \geq 1} \sqrt{j} \alpha_j h_{j-1} \quad \text{and} \quad \phi' = \sum_{j \geq 1} \sqrt{j} \beta_j h_{j-1}.$$

Since we additionally assume  $\mathbb{E}[g] = \alpha_0 = 0$ , by the definition of  $\tilde{\psi}_{\phi,g}$  in Lemma 12 and the properties of Hermite expansion discussed in Appendix B, we have

$$\tilde{\psi}_{\phi,g}(\mathbf{w}^t) = \frac{\sum_{j \geq s} (j-1) \alpha_j \beta_j \langle \bar{\mathbf{w}}^t, \mathbf{u} \rangle^j}{\|\Sigma^{1/2} \mathbf{w}\|} + \frac{1}{2},$$

and similarly

$$\tilde{\zeta}_{\phi,g}(\mathbf{w}^t) := - \sum_{j \geq s} j \alpha_j \beta_j \langle \bar{\mathbf{w}}^t, \mathbf{u} \rangle^{j-1}.$$

Thus we obtain

$$\frac{d\langle \mathbf{w}^t, \mathbf{u} \rangle^2}{dt} = 2F(\mathbf{w}^t)(1 - \langle \mathbf{w}^t, \mathbf{u} \rangle^2),$$

where

$$F(\mathbf{w}^t) := \sum_{j \geq s} j \alpha_j \beta_j \langle \bar{\mathbf{w}}^t, \mathbf{u} \rangle^{j-1} \langle \mathbf{w}^t, \mathbf{u} \rangle - \frac{\kappa}{1 + \kappa} \sum_{j \geq s} (j-1) \alpha_j \beta_j \langle \bar{\mathbf{w}}^t, \mathbf{u} \rangle^{j+1} \langle \mathbf{w}^t, \mathbf{u} \rangle - \frac{\kappa \langle \mathbf{w}^t, \mathbf{u} \rangle^2}{2(1 + \kappa)}.$$

We proceed by upper bounding  $F$ . To do so, first note that

$$\|\Sigma^{1/2} \mathbf{w}\| = \sqrt{\frac{1 + \kappa \langle \mathbf{w}, \mathbf{u} \rangle^2}{1 + \kappa}} \geq \sqrt{\frac{1}{1 + \kappa}}.$$

To bound the first term of  $F$ , we have

$$\begin{aligned} \sum_{j \geq s} j \alpha_j \beta_j \langle \bar{\mathbf{w}}^t, \mathbf{u} \rangle^{j-1} \langle \mathbf{w}^t, \mathbf{u} \rangle &\leq |\langle \mathbf{w}^t, \mathbf{u} \rangle| |\langle \bar{\mathbf{w}}^t, \mathbf{u} \rangle|^{s-1} \sum_{j \geq s} j |\alpha_j \beta_j| |\langle \bar{\mathbf{w}}^t, \mathbf{u} \rangle|^{j-s} \\ &\leq \|\phi'\|_\gamma \|g'\|_\gamma |\langle \bar{\mathbf{w}}^t, \mathbf{u} \rangle|^{s-1} |\langle \mathbf{w}^t, \mathbf{u} \rangle| \\ &\leq \|\phi'\|_\gamma \|g'\|_\gamma (1 + \kappa)^{(s-1)/2} |\langle \mathbf{w}^t, \mathbf{u} \rangle|^s. \end{aligned}$$

Similarly,

$$-\frac{\kappa}{1 + \kappa} \sum_{j \geq s} (j-1) \alpha_j \beta_j \langle \bar{\mathbf{w}}^t, \mathbf{u} \rangle^{j+1} \langle \mathbf{w}^t, \mathbf{u} \rangle \leq \|\phi'\|_\gamma \|g'\|_\gamma (1 + \kappa)^{(s+1)/2} |\langle \mathbf{w}^t, \mathbf{u} \rangle|^{s+2}.$$

Hence, for  $\kappa \geq 1$ ,

$$F(\mathbf{w}^t) \leq \|\phi'\|_\gamma \|g'\|_\gamma (1 + \kappa)^{(s-1)/2} |\langle \mathbf{w}^t, \mathbf{u} \rangle|^s \left(1 + (1 + \kappa) \langle \mathbf{w}^t, \mathbf{u} \rangle^2\right) - \frac{\langle \mathbf{w}^t, \mathbf{u} \rangle^2}{4}.$$

Suppose  $\kappa < d/C^2 - 1$ , then

$$\begin{aligned} F(\mathbf{w}^t) &\leq \langle \mathbf{w}^t, \mathbf{u} \rangle^2 \left(2\|\phi'\|_\gamma \|g'\|_\gamma (1 + \kappa)^{(s-1)/2} |\langle \mathbf{w}^t, \mathbf{u} \rangle|^{s-2} - 1/4\right) \\ &\leq \langle \mathbf{w}^t, \mathbf{u} \rangle^2 \left(2\|\phi'\|_\gamma \|g'\|_\gamma C^{s-2} \sqrt{\frac{(1 + \kappa)^{s-1}}{d^{s-2}}} - 1/4\right). \end{aligned}$$

Thus, for any  $\kappa$  such that

$$1 \leq \kappa \leq \left\{ \frac{d^{\frac{s-2}{s-1}}}{(8C^{s-2} \|\phi'\|_\gamma \|g'\|_\gamma)^{\frac{2}{s-1}}} \wedge \frac{d}{C^2} \right\} - 1$$

and any  $\mathbf{w}^t$  such that  $|\langle \mathbf{w}^t, \mathbf{u} \rangle| \leq C/\sqrt{d}$ , we have  $\frac{d\langle \mathbf{w}^t, \mathbf{u} \rangle^2}{dt} \leq 0$ , hence  $\sup_{t \geq 0} |\langle \mathbf{w}^t, \mathbf{u} \rangle| \leq C/\sqrt{d}$ , as long as the above holds true at initialization.

Finally, we will show  $|\langle \mathbf{w}^0, \mathbf{u} \rangle| \leq C/\sqrt{d}$  with probability at least 0.99 for a suitable choice of constant  $C$ . Indeed, this is an elementary concentration of measure result on the unit sphere. For simplicity, we avoid performing sharp probability of failure analysis and only remark that  $\mathbb{E}[\langle \mathbf{w}^0, \mathbf{u} \rangle^2] = 1/d$ , thus by the Markov inequality

$$\mathbb{P}(\langle \mathbf{w}^0, \mathbf{u} \rangle^2 \geq C^2/d) \leq 1/C^2,$$

hence a choice of  $C \geq 10$  suffices, and the proof is complete.  $\square$

## C.2 Extremely Large Spike

In this section, we will show that under extremely large spike, the spherical gradient flow (2.5) can potentially recover the true direction. Namely, we will prove the following proposition.

**Proposition 14.** *Suppose we initialize the spherical population gradient flow (2.5) from  $\mathbf{w}^0$ . Let  $\phi$  be the ReLU activation and assume*

$$\langle \phi, g \rangle_\gamma := \mathbb{E}_{z \sim \mathcal{N}(0,1)}[\phi(z)g(z)] = \alpha > 1/2,$$

and  $\kappa \geq \frac{C}{\langle \mathbf{w}^0, \mathbf{u} \rangle^2}$  for a sufficiently large constant  $C > 0$  depending only on  $g$ . Then, the gradient flow on the sphere satisfies

$$\begin{cases} \langle \mathbf{w}^T, \mathbf{u} \rangle \geq 1 - \varepsilon & \text{if } \langle \mathbf{w}^0, \mathbf{u} \rangle > 0 \\ \langle \mathbf{w}^T, \mathbf{u} \rangle \leq -1 + \varepsilon & \text{if } \langle \mathbf{w}^0, \mathbf{u} \rangle < 0 \end{cases} \quad (\text{C.3})$$

whenever

$$T \geq \frac{1}{\alpha - 1/2} \ln(2/\varepsilon). \quad (\text{C.4})$$

Before proceeding to the proof, we notice that if we uniformly initialize  $\mathbf{w}^0$  over  $\mathbb{S}^{d-1}$ , then the typical value for  $\langle \mathbf{w}^0, \mathbf{u} \rangle$  is of order  $d^{-1/2}$ , meaning that the above proposition asks for  $\kappa = \Omega(d)$ . This is a regime where lower bounds for the sample complexity of kernel methods are  $\Omega(1)$  [DWY21], thus no meaningful separation in terms of dimension dependency of the sample complexity between neural networks and kernel methods is possible, as the problem becomes effectively one-dimensional.

**Proof.** The cases where  $\langle \mathbf{w}^0, \mathbf{u} \rangle > 0$  and  $\langle \mathbf{w}^0, \mathbf{u} \rangle < 0$  are symmetric, thus we only present the proof for the former. Using (C.2), we can write the dynamics on the sphere more explicitly as

$$\frac{d\langle \mathbf{w}^t, \mathbf{u} \rangle}{dt} = \left\{ \underbrace{\frac{\mathbb{E}[\phi'(\langle \bar{\mathbf{w}}^t, \mathbf{z} \rangle)g'(\langle \mathbf{u}, \mathbf{z} \rangle)]}{1 + \kappa\langle \mathbf{w}^t, \mathbf{u} \rangle^2}}_{=:B_1} + \underbrace{\frac{\kappa\langle \mathbf{w}^t, \mathbf{u} \rangle \mathbb{E}[\phi(\langle \bar{\mathbf{w}}^t, \mathbf{z} \rangle)g(\langle \mathbf{u}, \mathbf{z} \rangle)]}{\sqrt{1+\kappa}\sqrt{1+\kappa\langle \mathbf{w}^t, \mathbf{u} \rangle^2}}}_{=:B_2} - \frac{\kappa\langle \mathbf{w}^t, \mathbf{u} \rangle}{2(1+\kappa)} \right\} (1 - \langle \mathbf{w}^t, \mathbf{u} \rangle^2).$$

Our goal is to study the regime of large  $\kappa$ , therefore we will bound how much  $B_1$  and  $B_2$  can deviate from their corresponding  $\kappa = \infty$  values. In particular, we have

$$B_1 \geq -\frac{\|\phi'\|_\gamma \|g'\|_\gamma}{1 + \kappa\langle \mathbf{w}^t, \mathbf{u} \rangle^2} = \frac{-\|g'\|_\gamma/2}{1 + \kappa\langle \mathbf{w}^t, \mathbf{u} \rangle^2}.$$

Furthermore, assuming  $\langle \mathbf{w}^t, \mathbf{u} \rangle > 0$  and  $\langle \phi, g \rangle_\gamma > 0$ , let  $c_\kappa(\mathbf{w}^t) := \frac{\kappa\langle \mathbf{w}^t, \mathbf{u} \rangle}{\sqrt{1+\kappa}\sqrt{1+\kappa\langle \mathbf{w}^t, \mathbf{u} \rangle^2}}$ . Then, by the Lipschitzness of  $\phi$ ,

$$\begin{aligned} B_2 &= c_\kappa(\mathbf{w}^t)\langle \phi, g \rangle_\gamma + c_\kappa(\mathbf{w}^t) \mathbb{E}[(\phi(\langle \bar{\mathbf{w}}^t, \mathbf{z} \rangle) - \phi(\langle \mathbf{u}, \mathbf{z} \rangle))g(\langle \mathbf{u}, \mathbf{z} \rangle)] \\ &\geq c_\kappa(\mathbf{w}^t)\langle \phi, g \rangle_\gamma - |c_\kappa(\mathbf{w}^t)| \|g\|_\gamma \|\bar{\mathbf{w}}^t - \mathbf{u}\| \\ &\geq c_\kappa(\mathbf{w}^t)\langle \phi, g \rangle_\gamma - |c_\kappa(\mathbf{w}^t)| \|g\|_\gamma \sqrt{2(1 - \langle \bar{\mathbf{w}}^t, \mathbf{u} \rangle^2)} \\ &\geq c_\kappa(\mathbf{w}^t)\langle \phi, g \rangle_\gamma - \frac{\sqrt{2\kappa}\|g\|_\gamma |\langle \mathbf{w}^t, \mathbf{u} \rangle|}{1 + \kappa\langle \mathbf{w}^t, \mathbf{u} \rangle^2}. \end{aligned}$$

where we used  $\bar{\mathbf{w}}^t = \frac{\Sigma^{1/2}\mathbf{w}^t}{\|\Sigma^{1/2}\mathbf{w}^t\|}$  in the last step. Suppose  $\langle \mathbf{w}^t, \mathbf{u} \rangle > 0$  (which holds at least on a neighborhood around initialization, and as we will see below holds for all  $t > 0$ ), then,

$$c_\kappa(\mathbf{w}^t) \geq \frac{\kappa\langle \mathbf{w}^t, \mathbf{u} \rangle^2}{1 + \kappa\langle \mathbf{w}^t, \mathbf{u} \rangle^2}.$$

As a result, we obtain

$$B_1 + B_2 - \frac{\kappa\langle \mathbf{w}^t, \mathbf{u} \rangle}{2(1+\kappa)} \geq \langle \phi, g \rangle_\gamma \left(1 - \frac{1}{\kappa\langle \mathbf{w}^t, \mathbf{u} \rangle^2}\right) - \frac{\|g'\|_\gamma/2 + \sqrt{2}\|g\|_\gamma}{\sqrt{\kappa\langle \mathbf{w}^t, \mathbf{u} \rangle^2}} - \frac{1}{2}.$$

Consequently, the lower bound of the time derivative of  $\langle \mathbf{w}^t, \mathbf{u} \rangle$  becomes larger as  $\langle \mathbf{w}^t, \mathbf{u} \rangle$  increases. Therefore, assuming  $\langle \mathbf{w}^0, \mathbf{u} \rangle > 0$ , we only need to control this lower bound at initialization. Assume

$$\kappa \geq \frac{4}{\langle \mathbf{w}^0, \mathbf{u} \rangle^2 (\alpha - 1/2)} \left\{ \alpha \vee \frac{(\|g'\|_\gamma + \sqrt{8}\|g\|_\gamma)^2}{\alpha - 1/2} \right\}.$$

From this, we conclude that when  $\langle \mathbf{w}^t, \mathbf{u} \rangle > 0$  and  $\langle \phi, g \rangle_\gamma = \alpha > 1/2$ , we have

$$\frac{d\langle \mathbf{w}^t, \mathbf{u} \rangle}{dt} \geq \frac{\alpha - 1/2}{2} (1 - \langle \mathbf{w}^t, \mathbf{u} \rangle^2),$$

integration yields the desired result.  $\square$

## D Proofs of Section 3

We begin by stating the closed-form expression for the population gradient, i.e. the counterpart of Lemma 12 in the normalized setting.

**Lemma 15.** *Consider the population risk  $\mathcal{R}(\mathbf{w})$  defined by (3.1), recall that*

$$\mathcal{R}(\mathbf{w}) = \mathbb{E} \left[ -\phi \left( \frac{\langle \mathbf{w}, \mathbf{x} \rangle}{\|\Sigma^{1/2} \mathbf{w}\|} \right) g(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) \right] + \frac{1}{2} \|\phi\|_\gamma^2 + \frac{1}{2} \mathbb{E}[y^2].$$

Then,

$$\nabla \mathcal{R}(\mathbf{w}) = \frac{\Sigma^{1/2} (\mathbf{I}_d - \bar{\mathbf{w}} \bar{\mathbf{w}}^\top) \zeta_{\phi, g}(\langle \bar{\mathbf{w}}, \bar{\mathbf{u}} \rangle) \bar{\mathbf{u}}}{\|\Sigma^{1/2} \mathbf{w}\|}, \quad (\text{D.1})$$

where

$$\zeta_{\phi, g}(\langle \bar{\mathbf{w}}, \bar{\mathbf{u}} \rangle) := -\mathbb{E}[\phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) g'(\langle \bar{\mathbf{u}}, \mathbf{z} \rangle)] = -\sum_{j \geq s} j \alpha_j \beta_j \langle \bar{\mathbf{w}}, \bar{\mathbf{u}} \rangle^{j-1}. \quad (\text{D.2})$$

**Proof.** Recall from (3.1) that

$$\begin{aligned} \nabla \mathcal{R}(\mathbf{w}) &= \nabla_{\mathbf{w}} \mathbb{E} \left[ -\phi \left( \frac{\langle \mathbf{w}, \mathbf{x} \rangle}{\|\Sigma^{1/2} \mathbf{w}\|} \right) g(\langle \bar{\mathbf{u}}, \mathbf{z} \rangle) \right] \\ &= -\frac{1}{\|\Sigma^{1/2} \mathbf{w}\|} \left( \mathbf{I}_d - \frac{\Sigma \mathbf{w} \mathbf{w}^\top}{\|\Sigma^{1/2} \mathbf{w}\|^2} \right) \Sigma^{1/2} \mathbb{E}[\phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) g'(\langle \bar{\mathbf{u}}, \mathbf{z} \rangle) \mathbf{z}] \\ &= \frac{\Sigma^{1/2} (\mathbf{I}_d - \bar{\mathbf{w}} \bar{\mathbf{w}}^\top)}{\|\Sigma^{1/2} \mathbf{w}\|} \{ \zeta_{\phi, g}(\langle \bar{\mathbf{w}}, \bar{\mathbf{u}} \rangle) \bar{\mathbf{u}} + \psi_{\phi, g}(\langle \bar{\mathbf{w}}, \bar{\mathbf{u}} \rangle) \bar{\mathbf{w}} \} \quad (\text{by Lemma 11}) \\ &= \frac{\Sigma^{1/2} (\mathbf{I}_d - \bar{\mathbf{w}} \bar{\mathbf{w}}^\top) \zeta_{\phi, g}(\langle \bar{\mathbf{w}}, \bar{\mathbf{u}} \rangle) \bar{\mathbf{u}}}{\|\Sigma^{1/2} \mathbf{w}\|}, \end{aligned}$$

where

$$\psi_{\phi, g}(\langle \bar{\mathbf{w}}, \bar{\mathbf{u}} \rangle) := -\mathbb{E}[\phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) g'(\langle \bar{\mathbf{u}}, \mathbf{z} \rangle) \langle \bar{\mathbf{w}}, \mathbf{z} \rangle - \phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) g'(\langle \bar{\mathbf{u}}, \mathbf{z} \rangle) \langle \bar{\mathbf{w}}, \bar{\mathbf{u}} \rangle] \quad (\text{D.3})$$

(the above is only a function of  $\langle \bar{\mathbf{w}}, \bar{\mathbf{u}} \rangle$  due to the Hermite expansion).  $\square$

Given the closed form of the population gradient, the proof of Lemma 3 is immediate by noticing that

$$\frac{d\bar{\mathbf{w}}^t}{dt} = \frac{(\mathbf{I}_d - \bar{\mathbf{w}}^t \bar{\mathbf{w}}^{t\top})}{\|\Sigma^{1/2} \mathbf{w}\|} \Sigma^{1/2} \frac{d\mathbf{w}^t}{dt}. \quad (\text{D.4})$$

Next, we move on to prove Proposition 4.

## D.1 Proof of Proposition 4

From Lemma 3, we have

$$\begin{aligned} \frac{d\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle}{dt} &= -\zeta_{\phi, g}(\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle) \|\Sigma^{1/2} \bar{\mathbf{u}}_{\perp}^t\|^2 \\ &\geq c\lambda_{\min}(\Sigma) \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle^{s-1} (1 - \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle^2), \end{aligned}$$

where  $\bar{\mathbf{u}}_{\perp}^t := \bar{\mathbf{u}} - \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle \bar{\mathbf{w}}^t$ . The above inequality and the fact that  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle > 0$  imply that  $\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle$  is non-decreasing in time. Let

$$T_1 := \sup\{t > 0 : \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle < 1/2\}.$$

Then, on  $t \in [0, T_1]$ , we have

$$\frac{d\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle}{dt} \geq \frac{3c\lambda_{\min}(\Sigma)}{4} \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle^{s-1},$$

and integration yields

$$T_1 \leq 0 \vee \frac{4}{3c\lambda_{\min}(\Sigma)} \begin{cases} 1/2 - \langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle & s = 1 \\ \ln(1/(2\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle)) & s = 2 \\ \frac{1}{s-2} ((1/\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle)^{s-2} - 2^{s-2}) & s > 2 \end{cases}$$

Therefore,  $T_1 \lesssim \tau_s(\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle) / \lambda_{\min}(\Sigma)$ .

For  $t > T_1$ , we have

$$\frac{d\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle}{dt} \geq \frac{c\lambda_{\min}(\Sigma)}{2^{s-1}} (1 - \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle^2).$$

Let

$$T_2 = \sup\{t > 0 : \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle < 1 - \varepsilon\}.$$

Once again, integration implies

$$T_2 \leq T_1 + \frac{2^{s-2}}{c\lambda_{\min}(\Sigma)} \ln(2/(3\varepsilon)),$$

which completes the proof.  $\square$

## D.2 Preliminary Lemmas for proving Theorem 5

We first introduce a number of concentration (and anti-concentration) lemmas that will be useful for proving Theorem 5.

**Lemma 16.** *Suppose  $\{\mathbf{z}^{(i)}\}_{i=1}^n \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d)$ , and  $g : \mathbb{R} \rightarrow \mathbb{R}$  satisfies  $|g(\cdot)| \leq C(1 + |\cdot|^p)$  for some  $C > 0$  and  $p \geq 1$ . Additionally, suppose  $\{\epsilon^{(i)}\}_{i=1}^n$  are i.i.d.  $\sigma$ -sub-Gaussian zero-mean noise independent of  $\{\mathbf{z}^{(i)}\}_{i=1}^n$ . Let  $y^{(i)} := g(\langle \bar{\mathbf{u}}, \mathbf{z}^{(i)} \rangle) + \epsilon^{(i)}$  for some  $\bar{\mathbf{u}} \in \mathbb{S}^{d-1}$ . Then, for any  $q > 0$ , with probability at least  $1 - 4d^{-q}$ , we have*

$$\left| y^{(i)} \right| \leq C + C(2 \ln(nd^q))^{p/2} + \sigma \sqrt{2 \ln(nd^q)} \lesssim \ln(nd^q)^{p/2},$$

for all  $i$ .

**Proof.** Notice that  $\langle \bar{\mathbf{u}}, \mathbf{z}^{(i)} \rangle \sim \mathcal{N}(0, 1)$ , thus

$$\mathbb{P}\left(\left|\langle \bar{\mathbf{u}}, \mathbf{z}^{(i)} \rangle\right| \geq \sqrt{2 \ln(nd^q)}\right) \leq 2n^{-1}d^{-q}.$$

Similarly, by the sub-Gaussian and zero-mean property of  $\epsilon^{(i)}$ ,

$$\mathbb{P}\left(\left|\epsilon^{(i)}\right| \geq \sigma \sqrt{2 \ln(nd^q)}\right) \leq 2n^{-1}d^{-q}.$$

Thus, by a union bound, we have

$$\left|\langle \bar{\mathbf{u}}, \mathbf{z}^{(i)} \rangle\right| \leq \sqrt{2 \ln(nd^q)} \quad \text{and} \quad \left|\epsilon^{(i)}\right| \leq \sigma \sqrt{2 \ln(nd^q)}, \quad \text{for all } 1 \leq i \leq n,$$

with probability at least  $1 - 4d^{-q}$ . Using the upper bound on  $|g|$  finishes the proof.  $\square$

**Lemma 17.** Suppose  $\{\bar{\mathbf{z}}^{(i)}\}_{i=1}^n$  are i.i.d. samples drawn uniformly from  $\mathbb{S}^{d-1}$ . Then,

$$\mathbb{P}\left(\sup_{\bar{\mathbf{w}} \in \mathbb{S}^{d-1}} \sum_{i=1}^n \mathbf{1}\left(\left|\langle \bar{\mathbf{w}}, \bar{\mathbf{z}}^{(i)} \rangle\right| \leq \frac{3\sqrt{d}}{8n}\right)\right) \geq 3d\left(2 + \ln(8n/\sqrt{d})\right) \leq e^{-d}.$$

**Proof.** Fix some  $\epsilon \in (0, 1)$ . Let  $N_\epsilon$  be a minimal  $\epsilon$ -covering of  $\mathbb{S}^{d-1}$ . Let  $\hat{\mathbf{w}}$  be the projection of  $\bar{\mathbf{w}}$  onto  $N_\epsilon$ . Notice that by the triangle inequality and the union bound

$$\begin{aligned} \mathbb{P}\left(\sup_{\bar{\mathbf{w}} \in \mathbb{S}^{d-1}} \sum_{i=1}^n \mathbf{1}\left(\left|\langle \bar{\mathbf{w}}, \bar{\mathbf{z}}^{(i)} \rangle\right| \leq \epsilon\right) \geq \alpha\right) &\leq \mathbb{P}\left(\sup_{\hat{\mathbf{w}} \in N_\epsilon} \sum_{i=1}^n \mathbf{1}\left(\left|\langle \hat{\mathbf{w}}, \bar{\mathbf{z}}^{(i)} \rangle\right| \leq 2\epsilon\right) \geq \alpha\right) \\ &\leq |N_\epsilon| \mathbb{P}\left(\sum_{i=1}^n \mathbf{1}\left(\left|\langle \hat{\mathbf{w}}, \bar{\mathbf{z}}^{(i)} \rangle\right| \leq 2\epsilon\right) \geq \alpha\right) \\ &\leq (3/\epsilon)^d \mathbb{P}\left(\sum_{i=1}^n \mathbf{1}\left(\left|\langle \hat{\mathbf{w}}, \bar{\mathbf{z}}^{(i)} \rangle\right| \leq 2\epsilon\right) \geq \alpha\right). \end{aligned}$$

Moreover, due to [BBSS22, Lemma A.7],

$$\mathbb{E}[\mathbf{1}(|\langle \hat{\mathbf{w}}, \bar{\mathbf{z}} \rangle| \leq 2\epsilon)] = \mathbb{P}(\langle \hat{\mathbf{w}}, \bar{\mathbf{z}} \rangle \leq 2\epsilon) \leq 8\sqrt{d}\epsilon.$$

Choose  $\epsilon = \frac{3\sqrt{d}}{8n}$ . By Lemma 25

$$\mathbb{P}\left(\sum_{i=1}^n \mathbf{1}\left(\left|\langle \hat{\mathbf{w}}, \bar{\mathbf{z}}^{(i)} \rangle\right| \leq 2\epsilon\right) \geq 3d\left(2 + \ln(8n/\sqrt{d})\right)\right) \leq (3/\epsilon)^d e^{-d},$$

which completes the proof.  $\square$

We summarize the above statements into a “good event”, as characterized by the following lemma.

**Lemma 18.** Let  $\{\mathbf{z}^{(i)}\}_{i=1}^n \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d)$ , and  $\bar{\mathbf{z}}^{(i)} := \mathbf{z}^{(i)} / \|\mathbf{z}^{(i)}\|$  for every  $i$ . We say event  $\mathcal{G}$  occurs whenever:

1.  $|y|^{(i)} \lesssim \ln(nd^q)^{p/2}$  for all  $1 \leq i \leq n$ .
2.  $\sup_{\bar{\mathbf{w}} \in \mathbb{S}^{d-1}} \sum_{i=1}^n \mathbf{1}\left(\left|\langle \bar{\mathbf{w}}, \bar{\mathbf{z}}^{(i)} \rangle\right| \lesssim \sqrt{d}/n\right) \lesssim d \ln(n/\sqrt{d})$ .
3.  $\lambda_{\max}\left(\frac{1}{n} \sum_{i=1}^n \mathbf{z}^{(i)} \mathbf{z}^{(i)\top}\right) - 1 \lesssim \sqrt{d/n}$ .
4.  $1 - \lambda_{\min}\left(\frac{1}{n} \sum_{i=1}^n \mathbf{z}^{(i)} \mathbf{z}^{(i)\top}\right) \lesssim \sqrt{d/n}$ .

For  $n \gtrsim d$ , event  $\mathcal{G}$  occurs with probability at least  $1 - \mathcal{O}(d^{-q})$ .

**Proof.** The first and second statements of the lemma follow from Lemmas 16 and 17 respectively. The third and fourth statements are standard Gaussian covariance concentration bounds (see e.g. [Wai19, Theorem 6.1] where both statements hold with probability at least  $1 - 2e^{-d}$ ).  $\square$

### D.3 Proof of Theorem 5

We begin by recalling the definition  $\bar{\mathbf{w}} := \frac{\hat{\Sigma}^{1/2} \mathbf{w}}{\|\hat{\Sigma}^{1/2} \mathbf{w}\|}$  and the finite-samples dynamics (3.5), which we copy here for the reader’s convenience,

$$\frac{d\mathbf{w}^t}{dt} = \eta(\mathbf{w}^t) \nabla_{\mathbf{w}^t} \left\{ \frac{1}{n} \sum_{i=1}^n \phi\left(\frac{\langle \mathbf{w}^t, \mathbf{x}^{(i)} \rangle}{\|\hat{\Sigma}^{1/2} \mathbf{w}^t\|}\right) y^{(i)} \right\},$$

where  $\eta(\mathbf{w}^t) = \|\hat{\Sigma}^{1/2} \mathbf{w}^t\|^2$ . Moreover, via chain rule, we obtain

$$\frac{d\bar{\mathbf{w}}^t}{dt} = \frac{(\mathbf{I}_d - \bar{\mathbf{w}}^t \bar{\mathbf{w}}^{t\top}) \hat{\Sigma}^{1/2} d\mathbf{w}^t}{\|\hat{\Sigma}^{1/2} \mathbf{w}^t\|}. \quad (\text{D.5})$$

Let  $\tilde{\mathbf{z}}^{(i)} := \hat{\Sigma}^{-1/2} \mathbf{x}^{(i)} = \hat{\Sigma}^{-1/2} \Sigma^{1/2} \mathbf{z}^{(i)}$ . Then,

$$\frac{d\bar{\mathbf{w}}^t}{dt} = (\mathbf{I}_d - \bar{\mathbf{w}}^t \bar{\mathbf{w}}^{t\top}) \hat{\Sigma} (\mathbf{I}_d - \bar{\mathbf{w}}^t \bar{\mathbf{w}}^{t\top}) \left\{ \frac{1}{n} \sum_{i=1}^n \phi' \left( \frac{\langle \bar{\mathbf{w}}^t, \mathbf{x}^{(i)} \rangle}{\|\hat{\Sigma}^{1/2} \bar{\mathbf{w}}^t\|} \right) y^{(i)} \tilde{\mathbf{z}}^{(i)} \right\}$$

To simplify the notation, define

$$\boldsymbol{\nu}(\bar{\mathbf{w}}) := (\mathbf{I}_d - \bar{\mathbf{w}} \bar{\mathbf{w}}^\top) \hat{\Sigma} (\mathbf{I}_d - \bar{\mathbf{w}} \bar{\mathbf{w}}^\top) \bar{\mathbf{w}}.$$

Then

$$\frac{d\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle}{dt} = \left\langle \boldsymbol{\nu}(\bar{\mathbf{w}}^t), \frac{1}{n} \sum_{i=1}^n \phi' \left( \langle \bar{\mathbf{w}}^t, \tilde{\mathbf{z}}^{(i)} \rangle \right) y^{(i)} \tilde{\mathbf{z}}^{(i)} \right\rangle.$$

We can decompose the above dynamics into a population term and three different error terms in the following manner:

$$\begin{aligned} \frac{d\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle}{dt} &= \langle \boldsymbol{\nu}(\bar{\mathbf{w}}^t), \mathbb{E}_{\mathbf{z}, y} [\phi'(\langle \bar{\mathbf{w}}^t, \mathbf{z} \rangle) y \mathbf{z}] \rangle \\ &\quad + \underbrace{\left\langle \boldsymbol{\nu}(\bar{\mathbf{w}}^t), \frac{1}{n} \sum_{i=1}^n \phi' \left( \langle \bar{\mathbf{w}}^t, \mathbf{z}^{(i)} \rangle \right) y^{(i)} \mathbf{z}^{(i)} - \mathbb{E}_{\mathbf{z}, y} [\phi'(\langle \bar{\mathbf{w}}^t, \mathbf{z} \rangle) y \mathbf{z}] \right\rangle}_{=:\mathcal{E}_1} \\ &\quad + \underbrace{\frac{1}{n} \sum_{i=1}^n \left\{ \phi' \left( \langle \bar{\mathbf{w}}^t, \tilde{\mathbf{z}}^{(i)} \rangle \right) - \phi' \left( \langle \bar{\mathbf{w}}^t, \mathbf{z}^{(i)} \rangle \right) \right\} y^{(i)} \langle \mathbf{z}^{(i)}, \boldsymbol{\nu}(\bar{\mathbf{w}}^t) \rangle}_{=:\mathcal{E}_2} \\ &\quad + \underbrace{\frac{1}{n} \sum_{i=1}^n \phi' \left( \langle \bar{\mathbf{w}}, \tilde{\mathbf{z}}^{(i)} \rangle \right) y^{(i)} \langle \tilde{\mathbf{z}}^{(i)} - \mathbf{z}^{(i)}, \boldsymbol{\nu}(\bar{\mathbf{w}}^t) \rangle}_{=:\mathcal{E}_3}. \end{aligned}$$

We will proceed in three steps. In the first, we bound  $\mathcal{E}_1$ , the concentration error. In the second, we bound  $\mathcal{E}_2$  and  $\mathcal{E}_3$ , the errors due to estimating  $\Sigma$  with  $\hat{\Sigma}$  (i.e. replacing  $\mathbf{z}^{(i)}$  with  $\tilde{\mathbf{z}}^{(i)}$ ). Finally, we will analyze the convergence time similar to that of Proposition 4. Throughout the proof, we will assume that the event  $\mathcal{G}$  of Lemma 18 occurs.

**Step 1. Controlling the concentration error  $\mathcal{E}_1$ .** Let  $K \asymp \ln(nd^q)^{p/2}$ , and notice that on event  $\mathcal{G}$  we have  $|y^{(i)}| \lesssim K$  for all  $i$ . Let  $y_K := y \mathbf{1}(|y| \leq K)$ . On the event  $\mathcal{G}$ , we have  $y_K^{(i)} = y^{(i)}$  for all  $i$ , and

$$\mathcal{E}_1 = \langle \boldsymbol{\nu}(\bar{\mathbf{w}}^t), \boldsymbol{\Delta}_n \rangle \geq -\|\boldsymbol{\Delta}_n\| \|\boldsymbol{\nu}(\bar{\mathbf{w}}^t)\|,$$

where

$$\boldsymbol{\Delta}_n := \frac{1}{n} \sum_{i=1}^n \phi' \left( \langle \bar{\mathbf{w}}, \mathbf{z}^{(i)} \rangle \right) y_K^{(i)} \mathbf{z}^{(i)} - \mathbb{E}_{\mathbf{z}, y} [\phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) y \mathbf{z}].$$

Thus, our objective is to bound  $\|\boldsymbol{\Delta}_n\|$  uniformly for all  $\bar{\mathbf{w}} \in \mathbb{S}^{d-1}$ . To that end, we first modify the expectation in the above definition so that the empirical average and expected value match in terms of their random variables. Specifically,

$$\begin{aligned} \sup_{\bar{\mathbf{w}}, \mathbf{v} \in \mathbb{S}^{d-1}} \langle \boldsymbol{\Delta}_n, \mathbf{v} \rangle &= \sup_{\bar{\mathbf{w}}, \mathbf{v} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n \phi' \left( \langle \bar{\mathbf{w}}, \mathbf{z}^{(i)} \rangle \right) y_K^{(i)} \langle \mathbf{z}^{(i)}, \mathbf{v} \rangle - \mathbb{E}_{\mathbf{z}, y} [\phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) y_K \langle \mathbf{z}, \mathbf{v} \rangle] \\ &\quad - \mathbb{E}_{\mathbf{z}, y} [\phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) y \langle \mathbf{z}, \mathbf{v} \rangle \mathbf{1}(|y| > K)]. \end{aligned}$$

By the Cauchy-Schwartz inequality,

$$\begin{aligned} |\mathbb{E}_{\mathbf{z}, y} [\phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) y \langle \mathbf{z}, \mathbf{v} \rangle \mathbf{1}(|y| > K)]| &\leq \mathbb{E}_{\mathbf{z}, y} \left[ \phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle)^2 y^2 \langle \mathbf{z}, \mathbf{v} \rangle^2 \right]^{1/2} \mathbb{E}[\mathbf{1}(|y| > K)]^{1/2} \\ &\lesssim \mathbb{E}[y^4]^{1/4} \mathbb{E}_{\mathbf{z}} \left[ \langle \mathbf{z}, \mathbf{v} \rangle^4 \right]^{1/4} \mathbb{P}(|y| > K)^{1/2} \\ &\lesssim d^{-q/2}, \end{aligned}$$

where the last inequality follows from Lemma 16. Hence,

$$\begin{aligned} \sup_{\bar{\mathbf{w}}, \mathbf{v} \in \mathbb{S}^{d-1}} \langle \Delta_n, \mathbf{v} \rangle &\leq \sup_{\bar{\mathbf{w}}, \mathbf{v} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n \phi'(\langle \bar{\mathbf{w}}, \mathbf{z}^{(i)} \rangle) y_K^{(i)} \langle \mathbf{z}^{(i)}, \mathbf{v} \rangle - \mathbb{E}[\phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) y_K \langle \mathbf{z}, \mathbf{v} \rangle] \\ &\quad + \mathcal{O}(d^{-q/2}). \end{aligned}$$

Next, we need to establish high-probability bounds via a covering argument. To simplify the exposition, define the stochastic process indexed by  $\bar{\mathbf{w}} \in \mathbb{S}^{d-1}$  and  $\mathbf{v} \in \mathbb{S}^{d-1}$  via

$$X_{\bar{\mathbf{w}}, \mathbf{v}}^{(i)} := \phi'(\langle \bar{\mathbf{w}}, \mathbf{z}^{(i)} \rangle) y_K^{(i)} \langle \mathbf{z}^{(i)}, \mathbf{v} \rangle.$$

Fix some  $\epsilon_w, \epsilon_v > 0$ . Let  $\Theta_w$  and  $\Theta_v$  be  $\epsilon_w$  and  $\epsilon_v$  coverings of  $\mathbb{S}^{d-1}$ , and let  $\hat{\mathbf{w}}$  and  $\hat{\mathbf{v}}$  denote the projection of  $\bar{\mathbf{w}}$  onto  $\Theta_w$  and of  $\mathbf{v}$  onto  $\Theta_v$  respectively, then

$$\begin{aligned} \sup_{\bar{\mathbf{w}}, \mathbf{v} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n X_{\bar{\mathbf{w}}, \mathbf{v}}^{(i)} - \mathbb{E}[X_{\bar{\mathbf{w}}, \mathbf{v}}] &= \sup_{\bar{\mathbf{w}}, \mathbf{v} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n (X_{\bar{\mathbf{w}}, \mathbf{v}}^{(i)} - X_{\bar{\mathbf{w}}, \hat{\mathbf{v}}}^{(i)}) + \frac{1}{n} \sum_{i=1}^n (X_{\bar{\mathbf{w}}, \hat{\mathbf{v}}}^{(i)} - X_{\hat{\mathbf{w}}, \hat{\mathbf{v}}}^{(i)}) \\ &\quad + \mathbb{E}_{\mathbf{z}, y}[X_{\bar{\mathbf{w}}, \hat{\mathbf{v}}} - X_{\bar{\mathbf{w}}, \mathbf{v}}] + \mathbb{E}_{\mathbf{z}, y}[X_{\hat{\mathbf{w}}, \hat{\mathbf{v}}} - X_{\bar{\mathbf{w}}, \hat{\mathbf{v}}}] \\ &\quad + \frac{1}{n} \sum_{i=1}^n X_{\hat{\mathbf{w}}, \hat{\mathbf{v}}}^{(i)} - \mathbb{E}_{\mathbf{z}, y}[X_{\hat{\mathbf{w}}, \hat{\mathbf{v}}}] \end{aligned}$$

We bound each of the terms using Cauchy-Schwartz. Specifically,

$$\frac{1}{n} \sum_{i=1}^n (X_{\bar{\mathbf{w}}, \mathbf{v}}^{(i)} - X_{\bar{\mathbf{w}}, \hat{\mathbf{v}}}^{(i)}) \leq \sqrt{\frac{1}{n} \sum_{i=1}^n \phi'(\langle \bar{\mathbf{w}}, \mathbf{z}^{(i)} \rangle)^2 y_K^{(i)2}} \sqrt{\frac{1}{n} \sum_{i=1}^n \langle \mathbf{z}^{(i)}, \mathbf{v} - \hat{\mathbf{v}} \rangle^2} \lesssim K \epsilon_v,$$

where we used the upper bound on the operator norm of  $\frac{1}{n} \sum_{i=1}^n \mathbf{z}^{(i)} \mathbf{z}^{(i)\top}$  from Lemma 18 together with the fact that  $n \gtrsim d$ . Similarly,

$$\mathbb{E}_{\mathbf{z}, y}[X_{\bar{\mathbf{w}}, \hat{\mathbf{v}}} - X_{\bar{\mathbf{w}}, \mathbf{v}}] \leq \mathbb{E}_{\mathbf{z}, y}[\phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle)^2 y_K^2]^{1/2} \mathbb{E}_{\mathbf{z}}[\langle \mathbf{z}, \mathbf{v} - \hat{\mathbf{v}} \rangle^2]^{1/2} \lesssim K \epsilon_v.$$

To bound the differences when we replace  $\bar{\mathbf{w}}$  with  $\hat{\mathbf{w}}$ , we need to make a distinction between ReLU and smooth activations as the respective arguments are to some extent different. When  $\phi'$  is Lipschitz,

$$\frac{1}{n} \sum_{i=1}^n (X_{\bar{\mathbf{w}}, \hat{\mathbf{v}}}^{(i)} - X_{\hat{\mathbf{w}}, \hat{\mathbf{v}}}^{(i)}) \leq \sqrt{\frac{1}{n} \sum_{i=1}^n (y_K^{(i)})^2 (\phi'(\langle \bar{\mathbf{w}}, \mathbf{z}^{(i)} \rangle) - \phi'(\langle \hat{\mathbf{w}}, \mathbf{z}^{(i)} \rangle))^2} \sqrt{\frac{1}{n} \sum_{i=1}^n \langle \mathbf{z}^{(i)}, \hat{\mathbf{v}} \rangle^2} \lesssim K \epsilon_w,$$

and

$$\mathbb{E}[X_{\hat{\mathbf{w}}, \hat{\mathbf{v}}} - X_{\bar{\mathbf{w}}, \hat{\mathbf{v}}}] \leq \mathbb{E}[y_K^2 (\phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) - \phi'(\langle \hat{\mathbf{w}}, \mathbf{z} \rangle))^2]^{1/2} \mathbb{E}[\langle \mathbf{z}, \hat{\mathbf{v}} \rangle^2]^{1/2} \lesssim K \epsilon_w.$$

Therefore, for a smooth activation  $\phi$  we choose  $\epsilon_v = \epsilon_w = \sqrt{d/n}$ , and obtain

$$\sup_{\bar{\mathbf{w}}, \mathbf{v} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n X_{\bar{\mathbf{w}}, \mathbf{v}}^{(i)} - \mathbb{E}_{\mathbf{z}, y}[X_{\bar{\mathbf{w}}, \mathbf{v}}] \leq \sup_{\hat{\mathbf{w}}, \hat{\mathbf{v}}} \frac{1}{n} \sum_{i=1}^n X_{\hat{\mathbf{w}}, \hat{\mathbf{v}}}^{(i)} - \mathbb{E}_{\mathbf{z}, y}[X_{\hat{\mathbf{w}}, \hat{\mathbf{v}}}] + \tilde{\mathcal{O}}(\sqrt{d/n}).$$

When  $\phi$  is the ReLU activation, we need to show that the sign of the preactivation changes only for a small number of samples when we change the weight  $\bar{\mathbf{w}}$  to  $\hat{\mathbf{w}}$ . Notice that

$$\begin{aligned} \text{sign}(\langle \bar{\mathbf{w}}, \mathbf{z}^{(i)} \rangle) \neq \text{sign}(\langle \hat{\mathbf{w}}, \mathbf{z}^{(i)} \rangle) &\implies \left| \langle \bar{\mathbf{w}}, \mathbf{z}^{(i)} \rangle \right| \leq \left| \langle \hat{\mathbf{w}} - \bar{\mathbf{w}}, \mathbf{z}^{(i)} \rangle \right| \\ &\implies \left| \langle \bar{\mathbf{w}}, \bar{\mathbf{z}}^{(i)} \rangle \right| \leq \epsilon_w. \end{aligned}$$

Recall that  $\bar{\mathbf{z}}^{(i)} := \mathbf{z}^{(i)} / \|\mathbf{z}^{(i)}\|$ . Choose  $\epsilon_w \asymp \sqrt{d/n}$ . On event  $\mathcal{G}$ , we know from Lemma 18 that at most  $\mathcal{O}(d \ln(n/\sqrt{d}))$  samples can satisfy the above condition. Therefore,

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n (X_{\bar{\mathbf{w}}, \hat{\mathbf{v}}}^{(i)} - X_{\hat{\mathbf{w}}, \hat{\mathbf{v}}}^{(i)}) &\leq \sqrt{\frac{1}{n} \sum_{i=1}^n (y_K^{(i)})^2 (\phi'(\langle \bar{\mathbf{w}}, \mathbf{z}^{(i)} \rangle) - \phi'(\langle \hat{\mathbf{w}}, \mathbf{z}^{(i)} \rangle))^2} \sqrt{\frac{1}{n} \sum_{i=1}^n \langle \mathbf{z}^{(i)}, \hat{\mathbf{v}} \rangle^2} \\ &\lesssim K \sqrt{\frac{d \ln(n/\sqrt{d})}{n}}. \end{aligned}$$

and

$$\begin{aligned}
\mathbb{E}[X_{\hat{\mathbf{w}}, \hat{\mathbf{v}}} - X_{\bar{\mathbf{w}}, \hat{\mathbf{v}}}] &\leq \mathbb{E}\left[y_K^2 (\phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) - \phi'(\langle \hat{\mathbf{w}}, \mathbf{z} \rangle))^2\right]^{1/2} \mathbb{E}[\langle \mathbf{z}, \hat{\mathbf{v}} \rangle^2]^{1/2} \\
&\leq K \mathbb{P}(\text{sign}(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) \neq \text{sign}(\langle \hat{\mathbf{w}}, \mathbf{z} \rangle))^{1/2} \\
&\leq K \mathbb{P}(|\langle \bar{\mathbf{w}}, \mathbf{z} \rangle| \leq \epsilon_w)^{1/2} \\
&\leq 2K \sqrt{d^{1/2} \epsilon_w},
\end{aligned}$$

where the last inequality follows from the anti-concentration on the sphere [BBSS22, Lemma A.7]. Thus, for ReLU we choose  $\epsilon_v \asymp \sqrt{d/n}$  and  $\epsilon_w \asymp \sqrt{d/n}$ , and once again obtain

$$\sup_{\bar{\mathbf{w}}, \mathbf{v} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n X_{\bar{\mathbf{w}}, \mathbf{v}}^{(i)} - \mathbb{E}_{\mathbf{z}, y}[X_{\bar{\mathbf{w}}, \mathbf{v}}] \leq \sup_{\hat{\mathbf{w}}, \hat{\mathbf{v}}} \frac{1}{n} \sum_{i=1}^n X_{\hat{\mathbf{w}}, \hat{\mathbf{v}}}^{(i)} - \mathbb{E}_{\mathbf{z}, y}[X_{\hat{\mathbf{w}}, \hat{\mathbf{v}}}] + \tilde{\mathcal{O}}(\sqrt{d/n}).$$

It remains to bound the term

$$\sup_{\hat{\mathbf{w}}, \hat{\mathbf{v}}} \frac{1}{n} \sum_{i=1}^n X_{\hat{\mathbf{w}}, \hat{\mathbf{v}}}^{(i)} - \mathbb{E}[X_{\hat{\mathbf{w}}, \hat{\mathbf{v}}}] .$$

Notice that for fixed  $\hat{\mathbf{w}}, \hat{\mathbf{v}}$ ,  $X_{\hat{\mathbf{w}}, \hat{\mathbf{v}}}$  is sub-Gaussian with sub-Gaussian norm  $\mathcal{O}(K)$ . Thus, via the sub-Gaussian maximal inequality [VH16, Lemma 5.2],

$$\sup_{\hat{\mathbf{w}}, \hat{\mathbf{v}}} \frac{1}{n} \sum_{i=1}^n X_{\hat{\mathbf{w}}, \hat{\mathbf{v}}}^{(i)} - \mathbb{E}[X_{\hat{\mathbf{w}}, \hat{\mathbf{v}}}] \lesssim \sqrt{K^2 d/n \ln(1/(\epsilon_w \epsilon_v))},$$

with probability at least  $1 - e^{-d}$ . Consequently, we have

$$\sup_{\bar{\mathbf{w}} \in \mathbb{S}^{d-1}} \|\Delta_n\| \leq \tilde{\mathcal{O}}(\sqrt{d/n} + d^{-q/2}),$$

with probability at least  $1 - \mathcal{O}(d^{-q})$ . Assuming that  $n$  grows at most polynomially in dimension and choosing a sufficiently large  $q$ , we have  $\sup_{\bar{\mathbf{w}} \in \mathbb{S}^{d-1}} \|\Delta_n\| \leq \tilde{\mathcal{O}}(\sqrt{d/n})$  with probability at least  $1 - \mathcal{O}(d^{-q})$ .

Finally, by Lemma 23,

$$\|\nu(\bar{\mathbf{w}}^t)\| \leq \lambda_{\max}(\hat{\Sigma}) \lesssim \lambda_{\max}(\Sigma), \tag{D.6}$$

with probability at least  $1 - e^{-n'/2}$ . Combining the above with the bound on  $\|\Delta_n\|$ , we have  $\mathcal{E}_1 \geq -\lambda_{\max}(\Sigma) \tilde{\mathcal{O}}(\sqrt{d/n})$  with probability at least  $1 - \mathcal{O}(d^{-q})$ , which concludes the first step of the proof.

**Step 2. Bounding the error due to the estimation of  $\Sigma$ , i.e.  $\mathcal{E}_2$  and  $\mathcal{E}_3$ .** Recall that we are considering the event  $\mathcal{G}$ , thus  $y^{(i)} = y_K^{(i)}$ . We can control each of the error terms separately. We begin by  $\mathcal{E}_2$ , where by Cauchy-Schwartz

$$\begin{aligned}
\mathcal{E}_2 &= \frac{1}{n} \sum_{i=1}^n \left\{ \phi'(\langle \bar{\mathbf{w}}^t, \tilde{\mathbf{z}}^{(i)} \rangle) - \phi'(\langle \bar{\mathbf{w}}^t, \mathbf{z}^{(i)} \rangle) \right\} y_K^{(i)} \langle \mathbf{z}^{(i)}, \nu(\bar{\mathbf{w}}^t) \rangle \\
&\geq -\sqrt{\frac{1}{n} \sum_{i=1}^n \left\{ \phi'(\langle \bar{\mathbf{w}}^t, \mathbf{z}^{(i)} \rangle) - \phi'(\langle \bar{\mathbf{w}}^t, \tilde{\mathbf{z}}^{(i)} \rangle) \right\}^2} \sqrt{\frac{1}{n} \sum_{i=1}^n y_K^{(i)2} \langle \mathbf{z}^{(i)}, \nu(\bar{\mathbf{w}}^t) \rangle^2} \\
&\geq -K \|\nu(\bar{\mathbf{w}}^t)\| \sqrt{\frac{1}{n} \sum_{i=1}^n \left\{ \phi'(\langle \bar{\mathbf{w}}^t, \mathbf{z}^{(i)} \rangle) - \phi'(\langle \bar{\mathbf{w}}^t, \tilde{\mathbf{z}}^{(i)} \rangle) \right\}^2},
\end{aligned}$$

where the last line follows from Lemma 18 and the fact that  $n \gtrsim d$ . When  $\phi'$  is additionally Lipschitz, we have

$$\mathcal{E}_2 \gtrsim -K \|\nu(\bar{\mathbf{w}}^t)\| \sqrt{\frac{1}{n} \sum_{i=1}^n \langle \bar{\mathbf{w}}^t, \mathbf{z}^{(i)} - \tilde{\mathbf{z}}^{(i)} \rangle^2}.$$

Moreover, for any  $\bar{\mathbf{w}} \in \mathbb{S}^{d-1}$ ,

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \langle \bar{\mathbf{w}}, \mathbf{z}^{(i)} - \tilde{\mathbf{z}}^{(i)} \rangle^2 &\leq \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{z}^{(i)} \mathbf{z}^{(i)\top} \right\|^2 \|(\mathbf{I}_d - \hat{\Sigma}^{-1/2} \Sigma^{1/2}) \bar{\mathbf{w}}\|^2 \\ &\leq \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{z}^{(i)} \mathbf{z}^{(i)\top} \right\|^2 \|\mathbf{I}_d - \hat{\Sigma}^{-1/2} \Sigma^{1/2}\|^2 \\ &\lesssim \frac{d}{n'}, \end{aligned}$$

where the last inequality holds with probability at least  $1 - 2e^{-d}$  on the event of Lemma 24. Hence for smooth activations we conclude

$$\mathcal{E}_2 \gtrsim -K \|\nu(\bar{\mathbf{w}}^t)\| \sqrt{d/n'}.$$

When  $\phi$  is the ReLU activation, we need a more involved argument to control  $\mathcal{E}_2$ . In particular, we will show that for any  $\bar{\mathbf{w}}$ , at most only  $\tilde{\mathcal{O}}(d)$  datapoints can have  $\text{sign}(\langle \bar{\mathbf{w}}, \mathbf{z}^{(i)} \rangle) \neq \text{sign}(\langle \bar{\mathbf{w}}, \tilde{\mathbf{z}}^{(i)} \rangle)$ . Notice that

$$\begin{aligned} \text{sign}(\langle \bar{\mathbf{w}}, \mathbf{z}^{(i)} \rangle) \neq \text{sign}(\langle \bar{\mathbf{w}}, \tilde{\mathbf{z}}^{(i)} \rangle) &\implies \left| \langle \bar{\mathbf{w}}, \mathbf{z}^{(i)} \rangle \right| \leq \left| \langle \bar{\mathbf{w}}, \mathbf{z}^{(i)} - \tilde{\mathbf{z}}^{(i)} \rangle \right| \\ &\implies \left| \langle \bar{\mathbf{w}}, \bar{\mathbf{z}}^{(i)} \rangle \right| \leq \|\mathbf{I}_d - \hat{\Sigma}^{-1/2} \Sigma^{1/2}\| \end{aligned} \quad (\text{D.7})$$

where  $\bar{\mathbf{z}}^{(i)} := \frac{\mathbf{z}^{(i)}}{\|\mathbf{z}^{(i)}\|}$  is distributed uniformly over  $\mathbb{S}^{d-1}$ . From Lemma 24 we have  $\|\mathbf{I}_d - \hat{\Sigma}^{-1/2} \Sigma^{1/2}\| \lesssim \sqrt{\frac{d}{n'}}$  with probability at least  $1 - 2e^{-d}$ . On the other hand, from Lemma 17 we know with probability at least  $1 - e^{-d}$ , for any  $\bar{\mathbf{w}} \in \mathbb{S}^{d-1}$  at most  $\tilde{\mathcal{O}}(d)$  of the labeled samples have  $|\langle \bar{\mathbf{w}}, \mathbf{z}^{(i)} \rangle| \lesssim \sqrt{d}/n$ . Recall that  $n' \gtrsim n^2$  when using the ReLU activation. This is precisely why we make this choice for the ReLU activation, as we need to balance the RHS of (D.7) which is of order  $\sqrt{d}/n'$  with the LHS of (D.7) which should at most be of order  $\sqrt{d}/n$  if we want to ensure only  $\tilde{\mathcal{O}}(d)$  samples satisfy the bound. When  $n' = n^2$  we can balance these two terms, thus with probability at least  $1 - 3e^{-d}$  the sign change can occur for at most  $\tilde{\mathcal{O}}(d)$  many samples, and

$$\frac{1}{n} \sum_{i=1}^n \left\{ \phi'(\langle \bar{\mathbf{w}}^t, \mathbf{z}^{(i)} \rangle) - \phi'(\langle \bar{\mathbf{w}}^t, \tilde{\mathbf{z}}^{(i)} \rangle) \right\}^2 \leq \tilde{\mathcal{O}}\left(\frac{d}{n}\right).$$

In this case, we end up with

$$\mathcal{E}_2 \geq -K \|\nu(\bar{\mathbf{w}}^t)\| \tilde{\mathcal{O}}(\sqrt{d/n}).$$

Bounding  $\mathcal{E}_3$  for ReLU and Lipschitz  $\phi'$  is identical. In both cases, by Cauchy-Schwartz,

$$\begin{aligned} \mathcal{E}_3 &\geq -\sqrt{\frac{1}{n} \sum_{i=1}^n \phi'(\langle \bar{\mathbf{w}}, \mathbf{z}^{(i)} \rangle)^2 y_K^{(i)2}} \sqrt{\frac{1}{n} \sum_{i=1}^n \langle \mathbf{z}^{(i)} - \tilde{\mathbf{z}}^{(i)}, \nu(\bar{\mathbf{w}}^t) \rangle^2} \\ &\gtrsim -K \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{z}^{(i)} \mathbf{z}^{(i)\top} \right\| \|\mathbf{I}_d - \hat{\Sigma}^{-1/2} \Sigma^{1/2}\| \|\nu(\bar{\mathbf{w}}^t)\| \\ &\gtrsim -K \|\nu(\bar{\mathbf{w}}^t)\| \sqrt{d/n'}, \end{aligned}$$

which holds on the intersection of event  $\mathcal{G}$  and of Lemma 23. At last, using the bound on  $\|\nu(\bar{\mathbf{w}}^t)\|$  from (D.6), we obtain

$$\mathcal{E}_2 \wedge \mathcal{E}_3 \geq -\lambda_{\max}(\Sigma) \tilde{\mathcal{O}}(\sqrt{d/n}),$$

with probability at least  $1 - \mathcal{O}(d^{-q})$ .

**Step 3. Analyzing the Convergence.** As a result of the previous steps, we have established

$$\frac{d\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle}{dt} \geq \langle \nu(\bar{\mathbf{w}}^t), \mathbb{E}[\phi'(\langle \bar{\mathbf{w}}^t, \mathbf{z} \rangle) \mathbf{y} \mathbf{z}] \rangle - \lambda_{\max}(\Sigma) \tilde{\mathcal{O}}(\sqrt{d/n}).$$

Thanks to Lemma 11, we can write

$$\begin{aligned}\mathbb{E}[\phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) \mathbf{y} \mathbf{z}] &= \mathbb{E}[\phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) g(\langle \bar{\mathbf{u}}, \mathbf{z} \rangle) \mathbf{z}] \\ &= -\zeta_{\phi, g}(\langle \bar{\mathbf{w}}, \bar{\mathbf{u}} \rangle) \bar{\mathbf{u}} - \psi_{\phi, g}(\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle) \bar{\mathbf{w}},\end{aligned}$$

where  $\zeta_{\phi, g}$  and  $\psi_{\phi, g}$  were introduced in (D.2) and (D.3) respectively. Recall the definition of  $\boldsymbol{\nu}(\bar{\mathbf{w}}^t)$ ,

$$\boldsymbol{\nu}(\bar{\mathbf{w}}^t) := (\mathbf{I}_d - \bar{\mathbf{w}}^t \bar{\mathbf{w}}^{t\top}) \hat{\boldsymbol{\Sigma}} (\mathbf{I}_d - \bar{\mathbf{w}}^t \bar{\mathbf{w}}^{t\top}) \bar{\mathbf{u}}.$$

Therefore,

$$\begin{aligned}\frac{d\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle}{dt} &\geq -\zeta_{\phi, g}(\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle) \bar{\mathbf{u}}^\top (\mathbf{I}_d - \bar{\mathbf{w}}^t \bar{\mathbf{w}}^{t\top}) \hat{\boldsymbol{\Sigma}} (\mathbf{I}_d - \bar{\mathbf{w}}^t \bar{\mathbf{w}}^{t\top}) \bar{\mathbf{u}} - \lambda_{\max}(\boldsymbol{\Sigma}) \tilde{\mathcal{O}}(\sqrt{d/n}) \\ &\geq c \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle^{s-1} \langle \bar{\mathbf{u}}_\perp^t, \hat{\boldsymbol{\Sigma}} \bar{\mathbf{u}}_\perp^t \rangle - \lambda_{\max}(\boldsymbol{\Sigma}) \tilde{\mathcal{O}}(\sqrt{d/n}) \quad (\text{By Assumption 2}) \\ &\geq c \lambda_{\min}(\hat{\boldsymbol{\Sigma}}) \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle^{s-1} (1 - \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle^2) - \lambda_{\max}(\boldsymbol{\Sigma}) \tilde{\mathcal{O}}(\sqrt{d/n}),\end{aligned}$$

where  $\bar{\mathbf{u}}_\perp^t := \bar{\mathbf{u}} - \langle \bar{\mathbf{u}}, \bar{\mathbf{w}}^t \rangle \bar{\mathbf{w}}^t$ . Moreover, from Lemma 23, we have

$$\begin{aligned}\lambda_{\min}(\hat{\boldsymbol{\Sigma}}) &\geq \lambda_{\min}(\boldsymbol{\Sigma}) \left( \frac{1}{2} - \sqrt{\frac{\text{tr}(\boldsymbol{\Sigma})}{n' \lambda_{\min}(\boldsymbol{\Sigma})}} \right) \\ &\geq \lambda_{\min}(\boldsymbol{\Sigma}) \left( \frac{1}{2} - \sqrt{\frac{d\kappa(\boldsymbol{\Sigma})}{n'}} \right),\end{aligned}$$

with probability at least  $1 - e^{-n'/8}$ . Hence, for  $n' \gtrsim d\kappa(\boldsymbol{\Sigma})$  we have  $\lambda_{\min}(\hat{\boldsymbol{\Sigma}}) \gtrsim \lambda_{\min}(\boldsymbol{\Sigma})$ , and consequently,

$$\frac{d\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle}{dt} \geq c' \lambda_{\min}(\boldsymbol{\Sigma}) \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle^{s-1} (1 - \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle^2) - \lambda_{\max}(\boldsymbol{\Sigma}) \tilde{\mathcal{O}}(\sqrt{d/n}),$$

where  $c'$  is a universal constant. Notice that the first term in the RHS above denotes the signal, while the second term denotes the noise. We want to ensure the noise remains smaller than the signal throughout the trajectory, which leads to the convergence of  $\bar{\mathbf{w}}^t$  to  $\bar{\mathbf{u}}$ . Notice that the signal term is first increasing, then decreasing for  $\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle \in [0, 1]$ . Thus, it suffices to ensure the noise is smaller than the signal on the two ends of the interval, i.e. at time  $t = 0$  and at time  $t = T$  where  $\langle \bar{\mathbf{w}}^T, \bar{\mathbf{u}} \rangle = 1 - \varepsilon$ . At initialization, this condition leads to

$$n \gtrsim C d \kappa(\boldsymbol{\Sigma})^2 \langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle^{2(1-s)},$$

and at time  $t = T$ , leads to

$$n \gtrsim \frac{C d \kappa(\boldsymbol{\Sigma})^2}{\varepsilon^2},$$

where  $C$  hides constant depending only on  $s$  and at most polylogarithmic factors of  $d$ . Thus, we have established the sample complexity as presented by Theorem 5.

It remains to obtain the convergence time. With the above sample complexity, we have

$$\frac{d\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle}{dt} \geq c'' \lambda_{\min}(\boldsymbol{\Sigma}) \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle^{s-1} (1 - \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle^2),$$

where  $c''$  is a universal constant. The rest of the proof follows by integration and is identical to the proof of Proposition 4 in Appendix D.1.  $\square$

#### D.4 Proof of Corollary 6

The proof follows immediately from Theorem 5 and the following lemma which describes how  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle$  behaves under different regimes of  $r_1$  and  $r_2$ .

**Lemma 19.** Suppose  $\Sigma$  follows the  $(\kappa, \boldsymbol{\theta})$ -spiked model,  $\mathbf{w}^0$  is sampled uniformly from  $\mathbb{S}^{d-1}$ ,  $n' \gtrsim d$ , and there exist universal constants  $C_2, C'_2, C_3, C'_3 > 0$  such that

$$C_2 d^{r_2} \leq \kappa \leq C'_2 d^{r_2} \quad \text{and} \quad C_3 d^{-r_1} \leq \langle \mathbf{u}, \boldsymbol{\theta} \rangle \leq C'_3 d^{-r_1}$$

for  $r_1 \in [0, 1/2]$  and  $r_2 \in [0, 1]$ . Then, conditioned on  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle > 0$ , with any arbitrarily large constant probability  $1 - \delta$ , for sufficiently large  $d$  (that depends on  $\delta$ ) we have

$$\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle \gtrsim \begin{cases} d^{-1/2} & 0 \leq r_2 < r_1 \\ d^{r_2 - r_1 - 1/2} & r_1 < r_2 < 2r_1 \\ d^{(r_2 - 1)/2} & 2r_1 < r_2 < 1 \end{cases}. \quad (\text{D.8})$$

**Proof.** By definition,

$$\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle = \frac{\langle \hat{\Sigma}^{1/2} \mathbf{w}^0, \Sigma^{1/2} \mathbf{u} \rangle}{\|\hat{\Sigma}^{1/2} \mathbf{w}^0\| \|\Sigma^{1/2} \mathbf{u}\|}.$$

Recall that we are conditioning our arguments on  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle > 0$ , hence the numerator of the above fraction is positive. To translate the sample complexities of Theorems 5 and 7 to the spiked model, our goal is to lower bound  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle$  in terms of  $d, r_1$ , and  $r_2$ .

We begin by observing that

$$\|\hat{\Sigma}^{1/2} \mathbf{w}\| \leq \|\hat{\Sigma}^{1/2} \Sigma^{-1/2}\| \|\Sigma^{1/2} \mathbf{w}\| \lesssim \|\Sigma^{1/2} \mathbf{w}\|,$$

where the last inequality holds on the event of Lemma 24, which happens with probability at least  $1 - 2e^{-d}$ . Consequently,

$$\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle \gtrsim \frac{\langle \hat{\Sigma}^{1/2} \mathbf{w}^0, \Sigma^{1/2} \mathbf{u} \rangle}{\|\Sigma^{1/2} \mathbf{w}^0\| \|\Sigma^{1/2} \mathbf{u}\|} = \frac{\langle \mathbf{w}^0, \Sigma \mathbf{u} \rangle + \langle \mathbf{w}^0, (\hat{\Sigma}^{1/2} - \Sigma^{1/2}) \Sigma^{1/2} \mathbf{u} \rangle}{\|\Sigma^{1/2} \mathbf{w}^0\| \|\Sigma^{1/2} \mathbf{u}\|}.$$

Furthermore, due to the Markov inequality,

$$\mathbb{P}\left(\langle \mathbf{w}^0, \boldsymbol{\theta} \rangle^2 \geq \frac{C_1}{d}\right) \leq 1/C_1.$$

Similarly (by conditioning on  $\hat{\Sigma}$ )

$$\mathbb{P}\left(\langle \mathbf{w}^0, (\hat{\Sigma}^{1/2} - \Sigma^{1/2}) \Sigma^{1/2} \mathbf{u} \rangle^2 \geq \frac{C_1 \|(\hat{\Sigma}^{1/2} - \Sigma^{1/2}) \Sigma^{1/2} \mathbf{u}\|^2}{d}\right) \leq 1/C_1.$$

Additionally, on the event of Lemma 24,

$$\|(\hat{\Sigma}^{1/2} - \Sigma^{1/2}) \Sigma^{1/2} \mathbf{u}\| \leq \|\hat{\Sigma}^{1/2} \Sigma^{-1/2} - \mathbf{I}_d\| \|\Sigma \mathbf{u}\| \lesssim \sqrt{d/n'} \|\Sigma \mathbf{u}\|.$$

Therefore, on the above events, for some absolute constant  $C > 0$

$$\begin{aligned} \langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle &\gtrsim \frac{\langle \mathbf{w}^0, \Sigma \mathbf{u} \rangle - C \|\Sigma \mathbf{u}\| \sqrt{1/n'}}{\sqrt{\frac{1+C'_2 C_1}{1+\kappa}} \|\Sigma \mathbf{u}\|} \\ &\gtrsim \frac{\langle \mathbf{w}^0, \mathbf{u} \rangle + \kappa \langle \mathbf{w}^0, \boldsymbol{\theta} \rangle \langle \mathbf{u}, \boldsymbol{\theta} \rangle - C(1 + \kappa |\langle \mathbf{u}, \boldsymbol{\theta} \rangle|) \sqrt{1/n'}}{\sqrt{1 + C'_2 C_1} \sqrt{1 + \kappa \langle \mathbf{u}, \boldsymbol{\theta} \rangle^2}}. \end{aligned}$$

Recall that  $C_2 d^{r_2} \leq \kappa \leq C'_2 d^{r_2}$  and  $C_3 d^{-r_1} \leq \langle \mathbf{u}, \boldsymbol{\theta} \rangle \leq C'_3 d^{-r_1}$  (notice that changing  $\boldsymbol{\theta}$  to  $-\boldsymbol{\theta}$  does not change the spiked model of Assumption 1, thus we can assume  $\langle \mathbf{u}, \boldsymbol{\theta} \rangle \geq 0$  without loss of generality). Then,

$$\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle \gtrsim \frac{\langle \mathbf{w}^0, \mathbf{u} \rangle + \kappa \langle \mathbf{w}^0, \boldsymbol{\theta} \rangle \langle \mathbf{u}, \boldsymbol{\theta} \rangle - C(1 + \kappa \langle \mathbf{u}, \boldsymbol{\theta} \rangle) \sqrt{1/n'}}{\sqrt{1 + C'_2 C_1} \sqrt{1 + C'_2 C'_3{}^2 d^{r_2 - 2r_1}}}. \quad (\text{D.9})$$

The last term in the numerator can be made arbitrarily small by sufficiently large  $n'$ , hence we focus on other terms for now. Intuitively, when  $r_2 < 2r_1$ , the denominator is of constant order. If additionally  $r_2 < r_1$ , the dominant term in the numerator is  $\langle \mathbf{w}^0, \mathbf{u} \rangle$  and  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle \asymp 1/\sqrt{d}$ , otherwise the dominant term is  $\kappa \langle \mathbf{w}^0, \boldsymbol{\theta} \rangle \langle \mathbf{u}, \boldsymbol{\theta} \rangle$  and  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle \asymp d^{r_2 - r_1 - 1/2}$ . On the other hand, when  $r_2 > 2r_1$ , the denominator is of order  $d^{r_2/2 - r_1}$ , and once again the dominant term of the numerator is  $\kappa \langle \mathbf{w}^0, \boldsymbol{\theta} \rangle \langle \mathbf{u}, \boldsymbol{\theta} \rangle$ , therefore  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle \asymp d^{(r_2 - 1)/2}$ . Using this intuition, we analyze each of the following regimes separately.

Case 1.  $0 < r_2 < r_1$ : In this case, by [BBSS22, Lemma A.7] we have  $|\langle \mathbf{w}^0, \mathbf{u} \rangle| \geq c/\sqrt{d}$  with probability at least  $1 - 4c$ . On the intersection of all considered events with  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle > 0$ , and for sufficiently large  $d$  and  $n' \gtrsim d$ , we must have  $\langle \mathbf{w}^0, \mathbf{u} \rangle > 0$  (otherwise  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle < 0$ ). Thus by plugging the values in (D.9),

$$\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle \gtrsim \frac{c - \sqrt{C_1} C_2' C_3' d^{r_2 - r_1} - C(1 + C_2' C_3' d^{r_2 - r_1}) \sqrt{d/n'}}{\sqrt{d} \sqrt{1 + C_2' C_1} \sqrt{1 + C_2' C_3'^2 d^{r_2 - 2r_1}}} \gtrsim \frac{1}{\sqrt{d}}. \quad (\text{D.10})$$

The intersection of all desired events and  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle > 0$  happens with probability at least  $\frac{1}{2} - 4c - 2/C_1 - 2e^{-d}$ , thus conditioned on  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle > 0$  the probability is at least  $1 - 8c - 4/C_1 - 4e^{-d}$ . Choosing sufficiently small  $c$ , large  $C_1$ , and respectively large  $d$  and  $n' \gtrsim d$  with sufficiently large absolute constant, we can arbitrarily increase the (constant) probability of success. Thus the analysis of this regime is complete.

Case 2.  $r_1 < r_2 < 2r_1$ : This time we use the fact that  $|\langle \mathbf{w}^0, \mathbf{u} \rangle| \leq \sqrt{C_1/d}$  with probability at least  $1 - 1/C_1$ , and  $|\langle \mathbf{w}^0, \boldsymbol{\theta} \rangle| \geq c/\sqrt{d}$  with probability at least  $1 - 4c$ . By an argument similar to the previous case, for sufficiently large  $d$  and  $n' \gtrsim d$ ,  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle > 0$  implies  $\langle \mathbf{w}^0, \boldsymbol{\theta} \rangle > 0$ , hence by (D.9)

$$\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle \gtrsim \frac{-\sqrt{C_1} + cC_2C_3d^{r_2 - r_1} - C(1 + C_2'C_3'd^{r_2 - r_1})\sqrt{d/n'}}{\sqrt{d}\sqrt{1 + C_2'C_1}\sqrt{1 + C_2'C_3'^2d^{r_2 - 2r_1}}} \gtrsim d^{r_2 - r_1 - 1/2}, \quad (\text{D.11})$$

with probability at least  $1 - 8c - 4/C_1 - 4e^{-d}$  when conditioned on  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle > 0$ .

Case 3.  $2r_1 < r_2 < 1$ : Once again recall (D.9). To bound the numerator, we repeat the exact same argument as in the previous case, thus

$$\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle \geq \frac{-\sqrt{C_1} + cC_2C_3d^{r_2 - r_1} - C(1 + C_2'C_3'd^{r_2 - 1})\sqrt{d/n'}}{\sqrt{(1 + C_2'C_1)C_2'C_3'^2d^{\frac{1+r_2-2r_1}{2}}}} \gtrsim d^{(r_2 - 1)/2}, \quad (\text{D.12})$$

which finishes the proof of the lemma.  $\square$

## E Proofs of Section 4

### E.1 Proof of Theorem 7

We recall from (D.5) that

$$\frac{d\bar{\mathbf{w}}^t}{dt} = \frac{(\mathbf{I}_d - \bar{\mathbf{w}}^t \bar{\mathbf{w}}^{t\top}) \hat{\boldsymbol{\Sigma}}^{1/2} d\mathbf{w}^t}{\|\hat{\boldsymbol{\Sigma}}^{1/2} \mathbf{w}\|}.$$

Furthermore, the preconditioned dynamics of  $\mathbf{w}^t$  given by (4.1) reads

$$\frac{d\mathbf{w}^t}{dt} = \frac{\eta(\mathbf{w}^t) \hat{\boldsymbol{\Sigma}}^{-1/2} (\mathbf{I}_d - \bar{\mathbf{w}}^t \bar{\mathbf{w}}^{t\top})}{\|\hat{\boldsymbol{\Sigma}}^{1/2} \mathbf{w}\|} \left\{ \frac{1}{n} \sum_{i=1}^n \phi'(\langle \bar{\mathbf{w}}^t, \tilde{\mathbf{z}}^{(i)} \rangle) y^{(i)} \tilde{\mathbf{z}}^{(i)} \right\},$$

where we recall  $\tilde{\mathbf{z}}^{(i)} := \hat{\boldsymbol{\Sigma}}^{-1/2} \mathbf{x}^{(i)}$ . Plugging in  $\eta(\mathbf{w}^t) = \|\hat{\boldsymbol{\Sigma}}^{1/2} \mathbf{w}\|^2$  yields

$$\begin{aligned} \frac{d\bar{\mathbf{w}}^t}{dt} &= (\mathbf{I}_d - \bar{\mathbf{w}}^t \bar{\mathbf{w}}^{t\top})^2 \left\{ \frac{1}{n} \sum_{i=1}^n \phi'(\langle \bar{\mathbf{w}}^t, \tilde{\mathbf{z}}^{(i)} \rangle) y^{(i)} \tilde{\mathbf{z}}^{(i)} \right\} \\ &= (\mathbf{I}_d - \bar{\mathbf{w}}^t \bar{\mathbf{w}}^{t\top}) \left\{ \frac{1}{n} \sum_{i=1}^n \phi'(\langle \bar{\mathbf{w}}^t, \tilde{\mathbf{z}}^{(i)} \rangle) y^{(i)} \tilde{\mathbf{z}}^{(i)} \right\} \end{aligned}$$

The rest of the analysis is identical to that of the proof of Theorem 5 in Appendix D.3. Specifically, by defining

$$\bar{\mathbf{u}}_{\perp}^t := \bar{\mathbf{u}} - \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle \bar{\mathbf{w}}^t,$$

we have

$$\begin{aligned} \frac{d\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle}{dt} &= \langle \bar{\mathbf{u}}_{\perp}^t, \mathbb{E}_{\mathbf{z}, y} [\phi'(\langle \bar{\mathbf{w}}^t, \mathbf{z} \rangle) y \mathbf{z}] \rangle \\ &\quad + \underbrace{\left\langle \bar{\mathbf{u}}_{\perp}^t, \frac{1}{n} \sum_{i=1}^n \phi'(\langle \bar{\mathbf{w}}^t, \mathbf{z}^{(i)} \rangle) y^{(i)} \mathbf{z}^{(i)} - \mathbb{E}_{\mathbf{z}, y} [\phi'(\langle \bar{\mathbf{w}}^t, \mathbf{z} \rangle) y \mathbf{z}] \right\rangle}_{=: \mathcal{E}_1} \\ &\quad + \underbrace{\frac{1}{n} \sum_{i=1}^n \left\{ \phi'(\langle \bar{\mathbf{w}}^t, \tilde{\mathbf{z}}^{(i)} \rangle) - \phi'(\langle \bar{\mathbf{w}}^t, \mathbf{z}^{(i)} \rangle) \right\} y^{(i)} \langle \mathbf{z}^{(i)}, \bar{\mathbf{u}}_{\perp}^t \rangle}_{=: \mathcal{E}_2} \\ &\quad + \underbrace{\frac{1}{n} \sum_{i=1}^n \phi'(\langle \bar{\mathbf{w}}, \tilde{\mathbf{z}}^{(i)} \rangle) y^{(i)} \langle \tilde{\mathbf{z}}^{(i)} - \mathbf{z}^{(i)}, \bar{\mathbf{u}}_{\perp}^t \rangle}_{=: \mathcal{E}_3}. \end{aligned}$$

As long as  $n \gtrsim d$ ,  $n' = n$  for the smooth case, and  $n' \gtrsim n^2$  for the ReLU case, the first two steps of the proof of Theorem 5 in Appendix D.3 implies that

$$\mathcal{E}_1 \wedge \mathcal{E}_2 \wedge \mathcal{E}_3 \geq -\|\bar{\mathbf{u}}_{\perp}^t\| \tilde{\mathcal{O}}(\sqrt{d/n}) \geq -\tilde{\mathcal{O}}(\sqrt{d/n}).$$

Once again, we apply Lemma 11 to obtain

$$\mathbb{E}[\phi'(\langle \bar{\mathbf{w}}, \mathbf{z} \rangle) y \mathbf{z}] = -\zeta_{\phi, g}(\langle \bar{\mathbf{w}}, \bar{\mathbf{u}} \rangle) \bar{\mathbf{u}} - \psi_{\phi, g}(\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle) \bar{\mathbf{w}},$$

with  $\zeta_{\phi, g}$  and  $\psi_{\phi, g}$  given in (D.2) and (D.3) respectively. As a result,

$$\begin{aligned} \frac{d\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle}{dt} &\geq -\zeta_{\phi, g}(\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle) \|\bar{\mathbf{u}}_{\perp}^t\|^2 - \tilde{\mathcal{O}}(\sqrt{d/n}) \\ &\geq c \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle^{s-1} (1 - \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle^2) - \tilde{\mathcal{O}}(\sqrt{d/n}) \quad (\text{By Assumption 2}). \end{aligned}$$

We need to ensure the noise term, i.e. the second term on the RHS remains smaller than the signal, i.e. the first term. The signal term attains its minimum at either initialization  $t = 0$  or at the end of the trajectory  $t = T$  where  $\langle \bar{\mathbf{w}}^T, \bar{\mathbf{u}} \rangle = 1 - \varepsilon$ , which imposes the following sufficient conditions on  $n$ . Namely, at initialization we require

$$n \gtrsim C d \langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle^{2(1-s)},$$

while at  $t = T$  we require

$$n \gtrsim C d / \varepsilon^2,$$

where  $C$  hides constant that only depend on  $s$  and polylogarithmic factors of  $d$ . Hence, we obtain

$$\frac{d\langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle}{dt} \geq c' \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle^{s-1} (1 - \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle^2).$$

for some universal constant  $c' > 0$ . Via integration (similar to the proof of Proposition 4 in Appendix D.1), for

$$T_1 := \sup\{t > 0 : \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle < 1/2\}$$

we obtain

$$T_1 \lesssim \tau_s(\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle),$$

and for

$$T_2 := \sup\{t > 0 : \langle \bar{\mathbf{w}}^t, \bar{\mathbf{u}} \rangle < 1 - \varepsilon\}$$

we obtain  $T_2 - T_1 \lesssim \ln(1/\varepsilon)$ , which completes the proof. We conclude by remarking that the proof of Corollary 8 is immediate given Theorem 7 and Lemma 19.  $\square$

## E.2 Preliminary Lemmas for Proving Theorem 9

We will adapt the following lemma from [MHPG<sup>+</sup>23], which provides a non-parametric approximation of  $g$  via random biases.

**Lemma 20.** [MHPG<sup>+</sup>23, Lemma 22] *For any smooth  $g : \mathbb{R} \rightarrow \mathbb{R}$  and  $\Delta > 0$ , let  $\tilde{g} : \mathbb{R} \rightarrow \mathbb{R}$  be a smooth function such that  $\tilde{g}(z) = g(z)$  for  $|z| \leq \Delta$  and  $\tilde{g}(-2\Delta) = \tilde{g}'(-2\Delta) = 0$ . Suppose  $\{b_j\}_{j=1}^m \stackrel{i.i.d.}{\sim} \text{Unif}(-2\Delta, 2\Delta)$ , and let  $\Delta_* := \Delta \sup_{|z| \leq 2\Delta} |\tilde{g}''(z)|$ . Then, there exist second layer weights  $\{a_j(b_j)\}_{j=1}^m$  with  $\|\mathbf{a}\| \lesssim \Delta_*/\sqrt{m}$ , such that for any fixed  $z \in [-\Delta, \Delta]$  and any  $\delta > 0$ , with probability at least  $1 - \delta$  over the random biases,*

$$\left| \sum_{j=1}^m a(b_j) \phi(z + b_j) - g(b_j) \right| \lesssim \Delta \Delta_* \sqrt{\frac{\ln(1/\delta)}{m}},$$

where  $\phi$  is the ReLU activation.

We use the above lemma to show the existence of a second layer with  $\tilde{\mathcal{O}}(1/\sqrt{m})$  norm with training error of order  $\tilde{\mathcal{O}}(1/m)$ .

**Lemma 21.** *For any  $\varepsilon < 1$ , suppose  $\langle \bar{\mathbf{w}}, \bar{\mathbf{u}} \rangle \geq 1 - \varepsilon$ . Then for any  $q > 0$ , sufficiently large  $d$ ,  $n \gtrsim d/\varepsilon^2$ , with probability at least  $1 - \mathcal{O}(d^{-q})$  over the random biases and the dataset, there exists a second layer  $\mathbf{a}$  with  $\|\mathbf{a}\| \leq \tilde{\mathcal{O}}(1/\sqrt{m})$  described by Lemma 20 such that*

$$\frac{1}{n} \sum_{i=1}^n \left( \sum_{j=1}^m a_j \phi(\langle \bar{\mathbf{w}}, \tilde{\mathbf{z}}^{(i)} \rangle + b_j) - y^{(i)} \right)^2 \lesssim \mathbb{E}[\varepsilon^2] + \tilde{\mathcal{O}}(1/m + \varepsilon),$$

where  $\phi$  is the ReLU activation.

**Proof.** We begin by replacing  $\bar{\mathbf{w}}$  and  $\tilde{\mathbf{z}}^{(i)}$  with  $\bar{\mathbf{u}}$  and  $\mathbf{z}^{(i)}$ . Specifically, via Jensen's inequality,

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{j=1}^m a_j \phi(\langle \bar{\mathbf{w}}, \tilde{\mathbf{z}}^{(i)} \rangle + b_j) - y^{(i)} \right\}^2 &\leq \underbrace{\frac{4}{n} \sum_{i=1}^n \left\{ \sum_{j=1}^m a_j \phi(\langle \bar{\mathbf{u}}, \mathbf{z}^{(i)} \rangle + b_j) - g(\langle \bar{\mathbf{u}}, \mathbf{z}^{(i)} \rangle) \right\}^2}_{=:\mathcal{E}_1} \\ &\leq \underbrace{\frac{4}{n} \sum_{i=1}^n \left\{ y^{(i)} - g(\langle \bar{\mathbf{u}}, \mathbf{z}^{(i)} \rangle) \right\}^2}_{=:\mathcal{E}_2} \\ &\leq \underbrace{\frac{4}{n} \sum_{i=1}^n \left( \sum_{j=1}^m a_j \left\{ \phi(\langle \bar{\mathbf{w}}, \tilde{\mathbf{z}}^{(i)} \rangle + b_j) - \phi(\langle \bar{\mathbf{w}}, \mathbf{z}^{(i)} \rangle + b_j) \right\} \right)^2}_{=:\mathcal{E}_3} \\ &\leq \underbrace{\frac{4}{n} \sum_{i=1}^n \left( \sum_{j=1}^m a_j \left\{ \phi(\langle \bar{\mathbf{w}}, \mathbf{z}^{(i)} \rangle + b_j) - \phi(\langle \bar{\mathbf{u}}, \mathbf{z}^{(i)} \rangle + b_j) \right\} \right)^2}_{=:\mathcal{E}_4}. \end{aligned}$$

We bound each term separately. For  $\mathcal{E}_1$ , we can invoke Lemma 20 which implies that each term in the sum can be bounded by  $\tilde{\mathcal{O}}(1/m)$  with probability at least  $1 - 1/(nd^q)$ , thus by a union bound, with probability at least  $1 - d^{-q}$  over the random biases,

$$\mathcal{E}_1 \leq \tilde{\mathcal{O}}(1/m).$$

By sub-Gaussianity of  $\varepsilon^{(i)}$  (hence sub-exponentiality of  $\varepsilon^{(i)2}$ ), for  $n \gtrsim d$  (with a sufficiently large constant) we have

$$\mathcal{E}_2 \lesssim \mathbb{E}[\varepsilon^2] + \sqrt{d/n},$$

with probability at least  $1 - e^{-d}$ .

For  $\mathcal{E}_3$ , via the Lipschitzness of ReLU and the Cauchy-Schwartz inequality we can write

$$\mathcal{E}_3 \leq \frac{\tilde{\mathcal{O}}(1)}{n} \sum_{i=1}^n \left\langle \bar{\mathbf{w}}, \tilde{\mathbf{z}}^{(i)} - \mathbf{z}^{(i)} \right\rangle^2 \leq \tilde{\mathcal{O}}(1) \|\mathbf{I} - \hat{\Sigma}^{-1/2} \Sigma^{1/2}\|^2 \leq \tilde{\mathcal{O}}(d/n'),$$

where we used the event of Lemma 24 which happens with probability at least  $1 - 2e^{-d}$ , and  $\tilde{\mathcal{O}}(1)$  represents a constant that depends at most polylogarithmically on  $d$ .

Finally, we bound the last term. Once again via the Lipschitzness of the ReLU activation and the Cauchy-Schwartz inequality

$$\mathcal{E}_4 \leq \frac{\tilde{\mathcal{O}}(1)}{n} \sum_{i=1}^n \left\langle \bar{\mathbf{w}} - \bar{\mathbf{u}}, \mathbf{z}^{(i)} \right\rangle^2 \leq \tilde{\mathcal{O}}(\|\bar{\mathbf{w}} - \bar{\mathbf{u}}\|^2) \leq \tilde{\mathcal{O}}(\varepsilon),$$

where once again we used the event of Lemma 24. On the intersection of all desired events, we have

$$\frac{1}{n} \sum_{i=1}^n \left( \sum_{j=1}^m a_j \phi(\langle \bar{\mathbf{w}}, \tilde{\mathbf{z}}^{(i)} \rangle + b_j) - y^{(i)} \right) \lesssim \mathbb{E}[\varepsilon^2] + \tilde{\mathcal{O}}(1/m + \sqrt{d/n} + d/n' + \varepsilon).$$

We conclude the proof by noticing that  $n' \gtrsim n^2$  and  $n \gtrsim d\varepsilon^{-2}$ .  $\square$

Additionally, we will use the following standard Lemma on the Rademacher complexity of two-layer neural networks, which in particular is a restatement of [MHPG<sup>+</sup>23, Lemma 18] in a way suitable for our analysis.

**Lemma 22.** *Let  $\mathcal{F}$  be a class of real-valued functions on  $(\mathbf{z}, y)$ . Given  $n$  samples  $\{\mathbf{z}^{(i)}, y\}_{i=1}^n$ , define the empirical Rademacher complexity of  $\mathcal{F}$  as*

$$\hat{\mathfrak{R}}_n(\mathcal{F}) := \mathbb{E}_{(\varsigma_i)_{i=1}^n} \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \varsigma_i f(\mathbf{z}^{(i)}, y^{(i)}) \right],$$

where  $(\varsigma_i)$  are i.i.d. Rademacher random variables (i.e.  $\pm 1$  with equal probability). Suppose  $\mathcal{F}$  is given by

$$\mathcal{F} := \left\{ (\mathbf{z}, y) \mapsto \left( \sum_{j=1}^m a_j \phi(\langle \bar{\mathbf{u}}, \mathbf{z} \rangle + b_j) - y \right)^2 \wedge C : \|\mathbf{a}\| \leq r_a/\sqrt{m}, \quad |b_j| \leq r_b, \forall 1 \leq j \leq m \right\},$$

for some fixed  $\bar{\mathbf{u}} \in \mathbb{S}^{d-1}$ . Suppose  $\{\mathbf{z}^{(i)}\}_{i=1}^n \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d)$ , and suppose  $|\phi'| \leq 1$ . Then,

$$\mathbb{E}_{(\mathbf{z}^{(i)}, y^{(i)})_{i=1}^n} \left[ \hat{\mathfrak{R}}_n(\mathcal{F}) \right] \leq \frac{2\sqrt{2C}(1 + r_b)r_a}{\sqrt{n}}.$$

**Proof.** See the proof of [MHPG<sup>+</sup>23, Lemma 18].  $\square$

### E.3 Proof of Theorem 9

Throughout the proof, we will assume  $\langle \bar{\mathbf{w}}, \bar{\mathbf{u}} \rangle \geq 1 - \varepsilon$  where we recall

$$\bar{\mathbf{w}} := \frac{\hat{\Sigma}^{1/2} \mathbf{w}}{\|\hat{\Sigma}^{1/2} \mathbf{w}\|} \quad \text{and} \quad \bar{\mathbf{u}} := \frac{\Sigma^{1/2} \mathbf{u}}{\|\Sigma^{1/2} \mathbf{u}\|}.$$

From either Theorem 5 or Theorem 7, we can assume  $\langle \bar{\mathbf{w}}, \bar{\mathbf{u}} \rangle \geq 1 - \varepsilon$  with probability at least  $1 - \mathcal{O}(d^{-q})$  for any fixed  $q > 0$ . For simplicity, let

$$\hat{y}(\tilde{\mathbf{z}}; \bar{\mathbf{w}}) = \sum_{j=1}^m a_j \phi(\langle \bar{\mathbf{w}}, \tilde{\mathbf{z}} \rangle + b_j),$$

and similarly define  $\hat{y}(\mathbf{z}; \bar{\mathbf{u}})$ . We define the following quantities,

$$\mathcal{R}(\bar{\mathbf{w}}) := \mathbb{E}_{\mathbf{z}, y} \left[ (\hat{y}(\tilde{\mathbf{z}}; \bar{\mathbf{w}}) - y)^2 \right] \quad \text{and} \quad R(\bar{\mathbf{u}}) := \mathbb{E}_{\mathbf{z}, y} \left[ (\hat{y}(\mathbf{z}; \bar{\mathbf{u}}) - y)^2 \right], \quad (\text{E.1})$$

and similarly define their empirical counterparts,

$$\hat{\mathcal{R}}(\bar{\mathbf{w}}) := \frac{1}{n} \sum_{i=1}^n \left( \hat{y}(\tilde{\mathbf{z}}^{(i)}; \bar{\mathbf{w}}) - y^{(i)} \right)^2 \quad \text{and} \quad \hat{R}(\bar{\mathbf{u}}) := \frac{1}{n} \sum_{i=1}^n \left( \hat{y}(\mathbf{z}^{(i)}; \bar{\mathbf{u}}) - y^{(i)} \right)^2. \quad (\text{E.2})$$

Notice that ultimately, we are interested in bounding  $\mathcal{R}(\bar{\mathbf{w}})$ . We break down the proof into three steps. In the first step, we show that  $\mathcal{R}(\bar{\mathbf{w}})$  can be upper bounded by  $R(\bar{\mathbf{u}})$ . Then, via a generalization bound, we show that the  $R(\bar{\mathbf{u}})$  can be upper bounded by  $\hat{R}(\bar{\mathbf{u}})$ . Finally, we show that  $\hat{R}(\bar{\mathbf{u}})$  can be upper bounded by the training error, i.e.  $\hat{\mathcal{R}}(\bar{\mathbf{w}})$ , and convex optimization of the last layer can attain the near-optimal value of this training error which is bounded by Lemma 21.

**Step 1. Bounding  $\mathcal{R}(\bar{\mathbf{w}})$  via  $R(\bar{\mathbf{u}})$ .** By Jensen's inequality,

$$\mathbb{E}_{\mathbf{z}, y} \left[ (\hat{y}(\tilde{\mathbf{z}}; \bar{\mathbf{w}}) - y)^2 \right] \leq 3 \mathbb{E}_{\mathbf{z}} \left[ (\hat{y}(\tilde{\mathbf{z}}; \bar{\mathbf{w}}) - \hat{y}(\mathbf{z}; \bar{\mathbf{w}}))^2 \right] + 3 \mathbb{E}_{\mathbf{z}} \left[ (\hat{y}(\mathbf{z}; \bar{\mathbf{w}}) - \hat{y}(\mathbf{z}; \bar{\mathbf{u}}))^2 \right] + 3 \mathbb{E}_{\mathbf{z}, y} \left[ (\hat{y}(\mathbf{z}; \bar{\mathbf{u}}) - y)^2 \right].$$

Suppose  $\|\mathbf{a}\| \leq r_a / \sqrt{m}$ . For the first term, by Lipschitzness of  $\phi$  and the Cauchy-Schwartz inequality

$$\begin{aligned} \mathbb{E}_{\mathbf{z}} \left[ (\hat{y}(\tilde{\mathbf{z}}; \bar{\mathbf{w}}) - \hat{y}(\mathbf{z}; \bar{\mathbf{w}}))^2 \right] &= \mathbb{E}_{\mathbf{z}} \left[ \left( \sum_{j=1}^m a_j \left\{ \phi \left( \langle \bar{\mathbf{w}}, \tilde{\mathbf{z}}^{(i)} \rangle + b_j \right) - \phi \left( \langle \bar{\mathbf{w}}, \mathbf{z}^{(i)} \rangle + b_j \right) \right\} \right)^2 \right] \\ &\leq r_a^2 \mathbb{E}_{\mathbf{z}} \left[ \langle \bar{\mathbf{w}}, \tilde{\mathbf{z}} - \mathbf{z} \rangle^2 \right] \\ &\leq r_a^2 \|\mathbf{I}_d - \hat{\Sigma}^{-1/2} \Sigma^{1/2}\|^2 \lesssim r_a^2 d / n', \end{aligned}$$

where the last inequality holds with probability at least  $1 - 2e^{-d}$  on the event of Lemma 24.

For the middle term, via a similar argument,

$$\mathbb{E}_{\mathbf{z}} \left[ (\hat{y}(\mathbf{z}; \bar{\mathbf{w}}) - \hat{y}(\mathbf{z}; \bar{\mathbf{u}}))^2 \right] \leq r_a^2 \mathbb{E}_{\mathbf{z}} \left[ \langle \bar{\mathbf{w}} - \bar{\mathbf{u}}, \mathbf{z} \rangle^2 \right] \leq 2r_a \varepsilon.$$

In what follows, we will restrict the analysis to the case where  $r_a = \tilde{\mathcal{O}}(1)$ . Therefore, we have

$$\mathbb{E}_{\mathbf{z}, y} \left[ (\hat{y}(\tilde{\mathbf{z}}; \bar{\mathbf{w}}) - y)^2 \right] \leq 3 \mathbb{E}_{\mathbf{z}, y} \left[ (\hat{y}(\mathbf{z}; \bar{\mathbf{u}}) - y)^2 \right] + \tilde{\mathcal{O}}(d/n' + \varepsilon).$$

**Step 2. Generalization: Bounding  $R(\bar{\mathbf{u}})$  via  $\hat{R}(\bar{\mathbf{u}})$ .**

Define the event

$$E := \left\{ |\langle \bar{\mathbf{u}}, \mathbf{z} \rangle| \vee |\epsilon| \leq \sqrt{2 \ln(nd^q)} \right\}.$$

and similarly define  $E^{(i)}$  by replacing  $\mathbf{z}$  and  $\epsilon$  with  $\mathbf{z}^{(i)}$  and  $\epsilon^{(i)}$  respectively. Via the Cauchy-Schwartz and Jensen inequalities

$$\begin{aligned} \mathbb{E}_{\mathbf{z}, y} \left[ (\hat{y}(\mathbf{z}; \bar{\mathbf{u}}) - y)^2 \right] &= \mathbb{E}_{\mathbf{z}, y} \left[ (\hat{y}(\mathbf{z}; \bar{\mathbf{u}}) - y)^2 \mathbf{1}(E) \right] + \mathbb{E}_{\mathbf{z}, y} \left[ (\hat{y}(\mathbf{z}; \bar{\mathbf{u}}) - y)^2 \mathbf{1}(E^C) \right] \\ &\leq \mathbb{E}_{\mathbf{z}, y} \left[ (\hat{y}(\mathbf{z}; \bar{\mathbf{u}}) - y)^2 \mathbf{1}(E) \right] + \sqrt{8} (\mathbb{E}[\hat{y}(\mathbf{z}; \bar{\mathbf{u}})^4] + \mathbb{E}[y^4])^{1/2} \mathbb{P}(E^C)^{1/2} \end{aligned}$$

Moreover,  $\mathbb{E}[y^4] \lesssim 1$ ,  $\mathbb{E}_{\mathbf{z}, y}[\hat{y}(\mathbf{z}; \bar{\mathbf{u}})^4] \leq \tilde{\mathcal{O}}(1)$ , and  $\mathbb{P}(E^C) \leq 4/(nd^q)$  (via a standard sub-Gaussian tail bound). Consequently,

$$\mathbb{E}_{\mathbf{z}, y} \left[ (\hat{y}(\mathbf{z}; \bar{\mathbf{u}}) - y)^2 \right] \leq \mathbb{E}_{\mathbf{z}, y} \left[ (\hat{y}(\mathbf{z}; \bar{\mathbf{u}}) - y)^2 \mathbf{1}(E) \right] + \tilde{\mathcal{O}}(n^{-1} d^{-q}),$$

Let

$$\ell(\mathbf{z}^{(i)}, y^{(i)}; \mathbf{a}, \mathbf{b}) := \left( \sum_{j=1}^m a_j \phi \left( \langle \bar{\mathbf{u}}, \mathbf{z}^{(i)} \rangle + b_j \right) - y^{(i)} \right)^2 \mathbf{1}(E).$$

Notice that  $\bar{\mathbf{u}}$  is fixed. Then, by a standard symmetrization argument (see e.g. [VH16, Lemma 7.4]) and Lemma 22

$$\mathbb{E} \left[ \sup_{\|\mathbf{a}\| \leq r_a / \sqrt{m}, |b_j| \leq r_b} \mathbb{E}_{\mathbf{z}, y} [\ell(\mathbf{z}, y; \mathbf{a}, \mathbf{b})] - \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{z}^{(i)}, y^{(i)}; \mathbf{a}, \mathbf{b}) \right] \leq 2 \mathbb{E} [\hat{\mathfrak{R}}_n(\mathcal{F})] \leq \tilde{\mathcal{O}}(\sqrt{1/n}).$$

where  $\hat{\mathfrak{R}}_n(\mathcal{F})$ . As the loss is bounded, we can apply McDiarmid's inequality to turn the above bound in expectation into a bound in probability, in particular

$$\sup_{\|\mathbf{a}\| \leq r_a / \sqrt{m}, |b_j| \leq r_b} \mathbb{E}_{\mathbf{z}, y} [\ell(\mathbf{z}, y; \mathbf{a}, \mathbf{b})] - \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{z}^{(i)}, y^{(i)}; \mathbf{a}, \mathbf{b}) \leq \tilde{\mathcal{O}}(\sqrt{d/n})$$

with probability at least  $1 - 2e^{-d}$ . Therefore, we conclude this step by noticing that

$$\begin{aligned} \mathbb{E}_{\mathbf{z}, y} \left[ (\hat{y}(\mathbf{z}; \bar{\mathbf{u}}) - y)^2 \right] &\leq \frac{1}{n} \sum_{i=1}^n \left( \hat{y}(\mathbf{z}^{(i)}; \bar{\mathbf{u}}) - y^{(i)} \right)^2 \mathbf{1}(E^{(i)}) + \tilde{\mathcal{O}}(\sqrt{d/n} + n^{-1}d^{-q}) \\ &\leq \frac{1}{n} \sum_{i=1}^n \left( \hat{y}(\mathbf{z}^{(i)}; \bar{\mathbf{u}}) - y^{(i)} \right)^2 + \tilde{\mathcal{O}}(\sqrt{d/n}), \end{aligned}$$

with probability at least  $1 - 2e^{-d}$ .

**Step 3. Bounding the training error and finishing the proof.** This step is similar to the proof of [MHPG<sup>+</sup>23, Theorem 4]. For conciseness, define

$$\hat{\mathcal{R}}(\mathbf{a}) := \frac{1}{n} \sum_{i=1}^n \left( \hat{y}(\mathbf{x}^{(i)}; \mathbf{W}, \mathbf{a}, \mathbf{b}) - y^{(i)} \right)^2.$$

and  $\hat{\mathcal{R}}_\lambda(\mathbf{a}) := \hat{\mathcal{R}}(\mathbf{a}) + \lambda \|\mathbf{a}\|^2/2$ . Our goal is to choose suitable  $\lambda$  such that the minimizer

$$\mathbf{a}^* := \arg \min_{\mathbf{a} \in \mathbb{R}^m} \hat{\mathcal{R}}(\mathbf{a}) + \lambda \|\mathbf{a}\|^2/2,$$

satisfies  $\|\mathbf{a}^*\| \leq r_a / \sqrt{m}$  while the value of the above minimization problem which we denote with  $\hat{\mathcal{R}}_\lambda^*$  does not significantly exceed

$$\min_{\|\mathbf{a}\| \leq r_a / \sqrt{m}} \hat{\mathcal{R}}(\mathbf{a}).$$

We argue that the suitable choice for  $\lambda$  is

$$\lambda \asymp \frac{m \mathbb{E}[\epsilon^2] + m\varepsilon + 1}{r_a^2} = \tilde{\Theta}(m \mathbb{E}[\epsilon^2] + m\varepsilon + 1). \quad (\text{E.3})$$

Let  $\hat{\mathcal{R}}^*$  denote the minimizer of the regularized problem and  $\tilde{\mathbf{a}} := \arg \min_{\|\mathbf{a}\| \leq r_a / \sqrt{m}} \hat{\mathcal{R}}(\mathbf{a})$ . From Lemma 21, with a proper choice of  $r_a = \tilde{\Theta}(1)$ , we have

$$\hat{\mathcal{R}}(\tilde{\mathbf{a}}) \lesssim \mathbb{E}[\epsilon^2] + \tilde{\mathcal{O}}(1/m + \varepsilon).$$

with probability at least  $1 - \mathcal{O}(d^{-q})$  over the biases and the dataset. Note that as  $\mathbf{a}^*$  is the minimizer of  $\hat{\mathcal{R}}_\lambda$ , we have

$$\hat{\mathcal{R}}(\mathbf{a}^*) + \frac{\lambda \|\mathbf{a}^*\|^2}{2} \leq \hat{\mathcal{R}}(\tilde{\mathbf{a}}) + \frac{\lambda \|\tilde{\mathbf{a}}\|^2}{2},$$

and in particular

$$\frac{\lambda \|\mathbf{a}^*\|^2}{2} \leq \hat{\mathcal{R}}(\tilde{\mathbf{a}}) + \frac{\lambda \|\tilde{\mathbf{a}}\|^2}{2} \implies \|\mathbf{a}^*\| \leq \tilde{\mathcal{O}}(1/\sqrt{m}).$$

and

$$\hat{\mathcal{R}}(\mathbf{a}^*) \leq \hat{\mathcal{R}}(\tilde{\mathbf{a}}) + \frac{\lambda \|\tilde{\mathbf{a}}\|^2}{2} \lesssim \mathbb{E}[\epsilon^2] + \tilde{\mathcal{O}}(1/m + \varepsilon).$$

Let  $\{\mathbf{a}^t\}_{t \geq 0}$  be the solution to the gradient flow of  $\mathbf{a}$ . Then,

$$\frac{d\|\mathbf{a}^t - \mathbf{a}^*\|^2}{dt} = -2\langle \mathbf{a}^t - \mathbf{a}^*, \nabla \hat{\mathcal{R}}_\lambda(\mathbf{a}^t) \rangle,$$

and by the first-order condition of strong convexity

$$\langle \mathbf{a}^t - \mathbf{a}^*, \nabla \hat{\mathcal{R}}_\lambda(\mathbf{a}^t) \rangle \geq \lambda \|\mathbf{a}^t - \mathbf{a}^*\|^2,$$

therefore

$$\|\mathbf{a}^{T'} - \mathbf{a}^*\|^2 \leq e^{-2\lambda T'} \|\mathbf{a}^0 - \mathbf{a}^*\|^2.$$

As the training error (of the regularized problem) is  $\lambda$ -strongly convex in  $\mathbf{a}$ , by applying the standard Polyak-Łojasiewicz condition, gradient flow for training  $\mathbf{a}$  obtains

$$\hat{\mathcal{R}}_\lambda(\mathbf{a}^{T'}) - \hat{\mathcal{R}}_\lambda^* \leq (\hat{\mathcal{R}}_\lambda(\mathbf{a}^0) - \hat{\mathcal{R}}_\lambda^*) e^{-2\lambda T'}.$$

Furthermore, since

$$\|\mathbf{a}^*\|^2 - \|\mathbf{a}^{T'}\|^2 \leq 2\|\mathbf{a}^*\| \|\mathbf{a}^{T'} - \mathbf{a}^*\| - \|\mathbf{a}^{T'} - \mathbf{a}^*\|^2 \leq 2\|\mathbf{a}^{T'} - \mathbf{a}^*\| \|\mathbf{a}^*\|,$$

we have

$$\hat{\mathcal{R}}(\mathbf{a}^{T'}) - \hat{\mathcal{R}}^* \leq 2\|\mathbf{a}^0 - \mathbf{a}^*\| \|\mathbf{a}^*\| e^{-\lambda T'} + (\hat{\mathcal{R}}_\lambda(\mathbf{a}^0) - \hat{\mathcal{R}}_\lambda^*) e^{-2\lambda T'}.$$

Consequently, choosing

$$T' \geq \frac{\ln\left(\frac{\|\mathbf{a}^0 - \mathbf{a}^*\|}{\|\mathbf{a}^*\|}\right)}{\lambda} \vee \frac{\ln\left(\frac{4\|\mathbf{a}^0 - \mathbf{a}^*\| \|\mathbf{a}^*\|}{\varepsilon}\right)}{\lambda} \vee \frac{\ln\left(\frac{2(\hat{\mathcal{R}}_\lambda(\mathbf{a}^0) - \hat{\mathcal{R}}_\lambda^*)}{\varepsilon}\right)}{2\lambda}, \quad (\text{E.4})$$

implies

$$\hat{\mathcal{R}}(\mathbf{a}^{T'}) \leq \hat{\mathcal{R}}^* + \varepsilon \quad \text{and} \quad \|\mathbf{a}^{T'}\| \leq 2\|\mathbf{a}^*\| \lesssim \tilde{\mathcal{O}}(1/\sqrt{m}).$$

Therefore

$$\hat{\mathcal{R}}(\mathbf{a}^{T'}) \lesssim \mathbb{E}[\varepsilon^2] + \tilde{\mathcal{O}}(1/m + \varepsilon).$$

Recall that

$$\hat{\mathcal{R}}(\mathbf{a}^{T'}) = \frac{1}{n} \sum_{i=1}^n \left( \hat{y}(\tilde{\mathbf{z}}^{(i)}; \bar{\mathbf{w}}) - y^{(i)} \right)^2,$$

is the final training error which we also denoted by  $\hat{\mathcal{R}}(\bar{\mathbf{w}})$  earlier in this section when were not focusing on the second layer. From the previous two steps, we know how to bound  $\mathcal{R}(\bar{\mathbf{w}})$  via  $\hat{\mathcal{R}}(\bar{\mathbf{u}})$ . Thus the last step is to upper bound  $\hat{\mathcal{R}}(\bar{\mathbf{u}})$  via  $\hat{\mathcal{R}}(\bar{\mathbf{w}})$ . To that end, via Jensen's inequality

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \left( \hat{y}(\mathbf{z}^{(i)}; \bar{\mathbf{u}}) - y^{(i)} \right)^2 &\leq \frac{3}{n} \sum_{i=1}^n \left( \hat{y}(\tilde{\mathbf{z}}^{(i)}; \bar{\mathbf{w}}) - y^{(i)} \right)^2 \\ &\quad + \frac{3}{n} \sum_{i=1}^n \left( \hat{y}(\mathbf{z}^{(i)}; \bar{\mathbf{w}}) - \hat{y}(\tilde{\mathbf{z}}^{(i)}; \bar{\mathbf{w}}) \right)^2 \\ &\quad + \frac{3}{n} \sum_{i=1}^n \left( \hat{y}(\mathbf{z}^{(i)}; \bar{\mathbf{w}}) - \hat{y}(\mathbf{z}^{(i)}; \bar{\mathbf{u}}) \right)^2. \end{aligned}$$

The first term on the RHS is  $\hat{\mathcal{R}}(\bar{\mathbf{w}})$  for which we developed a bound earlier in this step. Bounding the latter two terms can be performed similarly to the arguments in the previous sections. In particular,

$$\frac{1}{n} \sum_{i=1}^n \left( \hat{y}(\mathbf{z}^{(i)}; \bar{\mathbf{w}}) - \hat{y}(\tilde{\mathbf{z}}^{(i)}; \bar{\mathbf{w}}) \right)^2 \leq r_a^2 \|\mathbf{I}_d - \hat{\Sigma}^{-1/2} \Sigma^{1/2}\|^2 \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{z}^{(i)} \mathbf{z}^{(i)\top} \right\|^2 \leq \tilde{\mathcal{O}}(d/n'),$$

where the last inequality holds with probability at least  $1 - 2e^{-d}$  (over the event of Lemma 24). Similarly,

$$\frac{1}{n} \sum_{i=1}^n \left( \hat{y}(\mathbf{z}^{(i)}; \bar{\mathbf{w}}) - \hat{y}(\mathbf{z}^{(i)}; \bar{\mathbf{u}}) \right)^2 \leq r_a \|\bar{\mathbf{w}} - \bar{\mathbf{u}}\|^2 \leq \tilde{\mathcal{O}}(\varepsilon).$$

Putting the bounds back together (recall  $n' \geq n \gtrsim d\varepsilon^{-2}$ ), we arrive at

$$\hat{\mathcal{R}}(\bar{\mathbf{u}}) \lesssim \hat{\mathcal{R}}(\bar{\mathbf{w}}) + \tilde{\mathcal{O}}(\varepsilon + d/n') \lesssim \hat{\mathcal{R}}(\bar{\mathbf{w}}) + \tilde{\mathcal{O}}(\varepsilon).$$

Combining the result of this step with the two previous steps implies

$$\mathcal{R}(\bar{\mathbf{w}}) \lesssim \mathbb{E}[\varepsilon^2] + \tilde{\mathcal{O}}(1/m + \varepsilon),$$

with probability at least  $1 - \mathcal{O}(d^{-q})$  (when conditioned on  $\langle \bar{\mathbf{w}}^0, \bar{\mathbf{u}} \rangle > 0$ ) which completes the proof of Theorem 9.  $\square$

## F Auxiliary Lemmas

In this section, we recall a number of standard lemmas which we employ in various parts of our proofs.

**Lemma 23.** [Wai19, Theorem 6.1]. *Suppose  $\{\mathbf{x}^{(i)}\}_{i=1}^{n'} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \Sigma)$ . Let  $\hat{\Sigma} := \frac{1}{n'} \sum_{i=1}^{n'} \mathbf{x}^{(i)} \mathbf{x}^{(i)\top}$ . Then, for  $n' \geq \text{tr}(\Sigma)/\lambda_{\max}(\Sigma)$ ,*

$$\lambda_{\max}(\hat{\Sigma}) \leq \lambda_{\max}(\Sigma) \left( 4 + 5 \sqrt{\frac{\text{tr}(\Sigma)}{n' \lambda_{\max}(\Sigma)}} \right)$$

with probability at least  $1 - e^{-n'/2}$ . Furthermore, for  $n' \geq d$ ,

$$\lambda_{\min}(\hat{\Sigma}) \geq \lambda_{\min}(\Sigma) \left( \frac{1}{4} - \sqrt{\frac{\text{tr}(\Sigma)}{n \lambda_{\min}(\Sigma)}} \right)$$

with probability at least  $1 - e^{-n'/8}$ .

**Lemma 24.** *Suppose  $\{\mathbf{z}^{(i)}\}_{i=1}^{n'} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d)$ , let  $\mathbf{x}^{(i)} := \Sigma^{1/2} \mathbf{z}^{(i)}$  for some invertible  $\Sigma$ , and define*

$$\hat{\Sigma} := \frac{1}{n'} \sum_{i=1}^{n'} \mathbf{x}^{(i)} \mathbf{x}^{(i)\top}.$$

Then

$$\|\mathbf{I}_d - \hat{\Sigma}^{1/2} \Sigma^{-1/2}\| \vee \|\mathbf{I}_d - \hat{\Sigma}^{-1/2} \Sigma^{1/2}\| \lesssim \sqrt{\frac{d}{n'}}$$

with probability at least  $1 - 2e^{-d}$ .

**Proof.** We have

$$\begin{aligned} \|\mathbf{I}_d - \hat{\Sigma}^{-1/2} \Sigma^{1/2}\| &= \left\{ \lambda_{\max}(\hat{\Sigma}^{-1/2} \Sigma^{1/2}) - 1 \right\} \vee \left\{ 1 - \lambda_{\min}(\hat{\Sigma}^{-1/2} \Sigma^{1/2}) \right\} \\ &= \left\{ \lambda_{\max}(\Sigma^{1/2} \hat{\Sigma}^{-1} \Sigma^{1/2})^{1/2} - 1 \right\} \vee \left\{ 1 - \lambda_{\min}(\Sigma^{1/2} \hat{\Sigma}^{-1} \Sigma^{1/2}) \right\} \\ &= \left\{ \lambda_{\min}(\Sigma^{-1/2} \hat{\Sigma} \Sigma^{-1/2})^{-1/2} - 1 \right\} \vee \left\{ 1 - \lambda_{\max}(\Sigma^{-1/2} \hat{\Sigma} \Sigma^{-1/2})^{-1/2} \right\} \\ &= \left\{ \lambda_{\min} \left( \frac{1}{n'} \sum_{i=1}^{n'} \mathbf{z}^{(i)} \mathbf{z}^{(i)\top} \right)^{-1/2} - 1 \right\} \vee \left\{ 1 - \lambda_{\max} \left( \frac{1}{n'} \sum_{i=1}^{n'} \mathbf{z}^{(i)} \mathbf{z}^{(i)\top} \right)^{-1/2} \right\}. \end{aligned}$$

Similarly,

$$\|\mathbf{I}_d - \hat{\Sigma}^{1/2} \Sigma^{-1/2}\| = \left\{ \lambda_{\max} \left( \frac{1}{n'} \sum_{i=1}^{n'} \mathbf{z}^{(i)} \mathbf{z}^{(i)\top} \right)^{1/2} - 1 \right\} \vee \left\{ 1 - \lambda_{\min} \left( \frac{1}{n'} \sum_{i=1}^{n'} \mathbf{z}^{(i)} \mathbf{z}^{(i)\top} \right)^{1/2} \right\}$$

Moreover, by [Wai19, Example 6.2], we have with probability at least  $1 - 2e^{-d}$ ,

$$\lambda_{\max} \left( \frac{1}{n'} \sum_{i=1}^{n'} \mathbf{z}^{(i)} \mathbf{z}^{(i)\top} \right) \leq 1 + (\sqrt{2}+1) \sqrt{\frac{d}{n'}} \quad \text{and} \quad \lambda_{\min} \left( \frac{1}{n'} \sum_{i=1}^{n'} \mathbf{z}^{(i)} \mathbf{z}^{(i)\top} \right) \geq 1 - (\sqrt{2}+1) \sqrt{\frac{d}{n'}}.$$

Thus, for  $n' \gtrsim d$  (with a sufficiently large absolute constant), we have

$$\|\mathbf{I}_d - \hat{\Sigma}^{1/2} \hat{\Sigma}^{-1/2}\| \vee \|\mathbf{I}_d - \hat{\Sigma}^{-1/2} \Sigma^{1/2}\| \lesssim \sqrt{\frac{d}{n'}}$$

with probability at least  $1 - 2e^{-d}$ . □

**Lemma 25** (Chernoff's Inequality). *Suppose  $X_1, \dots, X_n$  are i.i.d. Bernoulli random variables, and further assume that  $\mathbb{E}[\sum_i X_i] \leq \mu$ . Then, for any  $\delta \geq 1$ ,*

$$\mathbb{P} \left( \sum_{i=1}^n X_i \geq \mu(1 + \delta) \right) \leq e^{-\mu\delta/3}. \quad (\text{F.1})$$

**Proof.** The proof follows from a standard Chernoff bound. From [Ver18, Theorem 2.3.1]

$$\mathbb{P} \left( \sum_i X_i \geq \mu(1 + \delta) \right) \leq e^{\mu(\delta - (1+\delta) \ln(1+\delta))},$$

(notice that the statement of [Ver18, Theorem 2.3.1] holds true even when  $\mathbb{E}[\sum_i X_i] = \mu$  is replaced with  $\mathbb{E}[\sum_i X_i] \leq \mu$ ). We conclude by remarking that  $\delta - (1 + \delta) \ln(1 + \delta) \leq -\delta/3$  for  $\delta \geq 1$ . □