
Provable Guarantees for Neural Networks via Gradient Feature Learning

Zhenmei Shi*, Junyi Wei*, Yingyu Liang

University of Wisconsin, Madison

zhmeishi@cs.wisc.edu, jwei53@wisc.edu, yliang@cs.wisc.edu

Abstract

Neural networks have achieved remarkable empirical performance, while the current theoretical analysis is not adequate for understanding their success, e.g., the Neural Tangent Kernel approach fails to capture their key feature learning ability, while recent analyses on feature learning are typically problem-specific. This work proposes a unified analysis framework for two-layer networks trained by gradient descent. The framework is centered around the principle of feature learning from gradients, and its effectiveness is demonstrated by applications in several prototypical problems such as mixtures of Gaussians and parity functions. The framework also sheds light on interesting network learning phenomena such as feature learning beyond kernels and the lottery ticket hypothesis.

1 Introduction

Neural network (NN) learning has achieved remarkable empirical success and has been a main driving force for the recent progress in machine learning and artificial intelligence. On the other hand, theoretical understandings significantly lag behind. Traditional analysis approaches are not adequate due to the overparameterization of practical networks and the non-convex optimization in the training via gradient descent. One line of work (e.g. [9, 31, 38, 60, 71, 123] and many others) shows under proper conditions, heavily overparameterized networks are approximately linear models over data-independent features, i.e., a linear function on the Neural Tangent Kernel (NTK). While making weak assumptions about the data and thus applicable to various settings, this approach requires the network learning to be approximately using fixed data-independent features (i.e., the kernel regime, or fixed feature methods). It thus fails to capture the feature learning ability of networks (i.e., to learn a feature mapping for the inputs which allow accurate prediction), which is widely believed to be the key factor to their empirical success in many applications (e.g., [54, 77, 117, 119]). To study feature learning in networks, a recent line of work (e.g. [5, 6, 14, 33, 52, 72, 76, 116] and others) shows examples where networks provably enjoy advantages over fixed feature methods (including NTK), under different settings and assumptions. While providing more insights, these studies typically focus on specific problems, and their analyses exploit the specific properties of the problems and appear to be unrelated to each other. *Is there a common principle for feature learning in networks via gradient descent? Is there a unified analysis framework that can clarify the principle and also lead to provable error guarantees for prototypical problem settings?*

In this work, we take a step toward this goal by proposing a gradient feature learning framework for analyzing two-layer network learning by gradient descent. (1) The framework makes essentially no assumption about the data distribution and can be applied to various problems. Furthermore, it is centered around features from gradients, clearly illustrating how gradient descent leads to feature learning in networks and subsequently accurate predictions. (2) It leads to error guarantees competitive with the optimal in a family of networks that use the features induced by gradients on the

*Equal contribution.

data distribution. Then for a specific problem with structured data distributions, if the optimal in the induced family is small, the framework gives a small error guarantee.

We then apply the framework to several prototypical problems: mixtures of Gaussians, parity functions, linear data, and multiple-index models. These have been used for studying network learning (in particular, for the feature learning ability), but with different and seemingly unrelated analyses. In contrast, straightforward applications of our framework give small error guarantees, where the main effort is to compute the optimal in the induced family. Furthermore, in some cases, such as parities, we can handle more general data distributions than in the existing work.

Finally, we also demonstrate that the framework sheds light on several interesting network learning phenomena or implications such as feature learning beyond the kernel regime, lottery ticket hypothesis (LTH), simplicity bias, learning over different data distributions, and new perspectives about roadmaps forward. Due to space limitations, we present implications about features beyond the kernel regime and LTH in the main body but defer the other implications in Appendix C with a brief here. (1) For simplicity bias, it is generally believed that the optimization has some *implicit regularization* effect that restricts learning dynamics to a low capacity subset of the whole hypothesis class, so can lead to good generalization [53, 90]. Our framework provides an explanation that the learning first learns simpler functions and then more sophisticated ones. (2) For learning over different data distributions, we provide data-dependent non-vacuous guarantees, as our framework can be viewed as using the optimal gradient-induced NN to measure or quantify the “complexity” of the problem. For easier problems, this quantity is smaller, and our framework can give a better error bound to derive guarantees. (3) For new perspectives about roadmaps forward, our framework suggests the strong representation power of NN is actually the key to successful learning, while traditional ones suggest strong representation power leads to vacuous generalization bounds [19, 33]. Thus, we suggest a different analysis road. Traditional analysis typically first reasons about the optimal based on the whole function class then analyzes how NN learns proper features and reaches the optimal. In contrast, our framework defines feature family first, and then reasons about the optimal based on it.

2 Related Work

Neural Networks Learning Analysis. Recently there has been an increasing interest in the analysis of network learning. One line of work connects the sufficiently over-parameterized neural network to linear methods around its initialization like NTK (e.g. [9, 11, 20, 21, 31, 38, 49, 60, 62, 69, 71, 78, 82, 91, 93, 95, 114, 121, 122] and more), so that the neural network training is a convex problem. The key idea is that it suffices to consider the first-order Taylor expansion of the neural network around the origin when the initialization is large enough. However, NTK lies in the lazy training (kernel) regime that excludes feature learning [29, 50, 68, 113]. Many studies (e.g. [2, 5, 6, 8, 12, 14, 22, 26, 33, 37, 51, 52, 57, 58, 70, 72, 73, 76, 99, 112, 115, 116] and more) show that neural networks take advantage over NTK empirically and theoretically. Another line of work is the mean-field (MF) analysis of neural networks (e.g. [27, 28, 36, 79, 80, 100, 106] and more). The insight is to see the training dynamics of a sufficiently large-width neural network as a PDE. It uses a smaller initialization than the NTK so that the parameters may move away from the initialization. However, the MF does not provide explicit convergence rates and requires an unrealistically large width of the neural network. One more line of work is neural networks max-margin analysis (e.g. [30, 47, 48, 56, 61, 63, 74, 75, 83, 85, 86, 107, 109] and more). They need a strong assumption that the convergence starts from weights having perfect training accuracy, while feature learning happens in the early stage of training. To explain the success of neural networks beyond the limitation mentioned above, some work introduces the low intrinsic dimension of data distributions [17, 18, 23, 24, 25, 44, 67, 104, 108, 124]. Another recent line of work is that a trained network can exactly recover the ground truth or optimal solution or teacher network [3, 4, 10, 39, 84, 87, 94, 96, 120], but they have strong assumptions on data distribution or model structure, e.g., Gaussian marginals. [1, 40, 55, 110, 111] show that training dynamics of neural networks have multiple phases, e.g., feature learning at the beginning, and then dynamics in convex optimization which requires proxy convexity [43] or PL condition [65] or special data structure.

Feature Learning Based on Gradient Analysis. A recent line of work is studying how features emerge from the gradient. [7, 46] consider linear separable data and show that the first few gradient steps can learn good features, and the later steps learn a good network on neurons with these features. [33, 45, 105] have similar conclusions on non-linear data (e.g., parity functions), while in their problems one feature is sufficient for accurate prediction (i.e., single-index data model).

[32] considers multiple-index with low-degree polynomials as labeling functions and shows that a one-step gradient update can learn multiple features that lead to accurate prediction. [13, 81] studies one gradient step feature improvements at different learning rates. [97] proposes Recursive Feature Machines to show the mechanism of recursively feature learning but without giving a final loss guarantee. These studies consider specific problems and exploit properties of the data to analyze the gradient delicately, while our work provides a general framework applicable to different problems.

3 Gradient Feature Learning Framework

Problem Setup. We denote $[n] := \{1, 2, \dots, n\}$ and $\mathcal{O}(\cdot)$, $\tilde{(\cdot)}$, $\hat{(\cdot)}$ to omit the log term inside. Let $X \subseteq \mathbb{R}^d$ denote the input space, $Y \subseteq \mathbb{R}$ the label space. Let D be an arbitrary data distribution over $X \times Y$. Denote the class of two-layer networks with m neurons as:

$$F_{d,m} := \{f_{(\mathbf{a}, \mathbf{W}; \mathbf{b})} \mid f_{(\mathbf{a}, \mathbf{W}; \mathbf{b})}(\mathbf{x}) := \mathbf{a}^\top [\sigma(\mathbf{W}^\top \mathbf{x} + \mathbf{b})] = \sum_{i \in [m]} \mathbf{a}_i [\sigma(\mathbf{w}_i^\top \mathbf{x} + b_i)]\}, \quad (1)$$

where $\sigma(z) = \max(z, 0)$ is the ReLU activation function, $\mathbf{a} \in \mathbb{R}^m$ is the second layer weight, $\mathbf{W} \in \mathbb{R}^{d \times m}$ is the first layer weight, \mathbf{w}_i is the i -th column of \mathbf{W} (i.e., the weight for the i -th neuron), and $\mathbf{b} \in \mathbb{R}^m$ is the bias for the neurons. For technical simplicity, we only train \mathbf{a} , \mathbf{W} but not \mathbf{b} . Let superscript (t) denote the time step, e.g., $f_{(\mathbf{a}^{(t)}, \mathbf{W}^{(t)}; \mathbf{b})}$ denote the network at time step t . Denote $\mathbf{a} := (\mathbf{a}, \mathbf{W}, \mathbf{b})$, $\mathbf{a}^{(t)} := (\mathbf{a}^{(t)}, \mathbf{W}^{(t)}, \mathbf{b})$. The goal of neural network learning is to minimize the expected risk, i.e., $L_{\mathcal{D}}(f) := \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} L_{(\mathbf{x}, y)}(f)$, where $L_{(\mathbf{x}, y)}(f) = \ell(yf(\mathbf{x}))$ is the loss on an example (\mathbf{x}, y) for some loss function $\ell(\cdot)$, e.g., the hinge loss $\ell(z) = \max\{0, 1 - z\}$, and the logistic loss $\ell(z) = \log[1 + \exp(-z)]$. We also consider ℓ_2 regularization. The regularized loss with regularization coefficient λ is $L_{\mathcal{D}}(f) := L_{\mathcal{D}}(f) + \frac{\lambda}{2}(k\mathbf{W}k_F^2 + k\mathbf{a}k_2^2)$. Given a training set with n i.i.d. samples $Z = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i \in [n]}$ from D , the empirical risk and its regularized version are:

$$\tilde{L}_{\mathcal{Z}}(f) := \frac{1}{n} \sum_{i \in [n]} L_{(\mathbf{x}^{(i)}, y^{(i)})}(f), \quad \tilde{L}_{\mathcal{Z}}(f) := \tilde{L}_{\mathcal{Z}}(f) + \frac{\lambda}{2}(k\mathbf{W}k_F^2 + k\mathbf{a}k_2^2). \quad (2)$$

Then the training process is summarized in Algorithm 1.

Algorithm 1 Network Training via Gradient Descent

```

Initialize  $(\mathbf{a}^{(0)}, \mathbf{W}^{(0)}, \mathbf{b})$ 
for  $t = 1$  to  $T$  do
  Sample  $Z^{(t-1)} \sim D^n$ 
   $\mathbf{a}^{(t)} = \mathbf{a}^{(t-1)} - \eta^{(t)} \nabla_{\mathbf{a}} \tilde{L}_{Z^{(t-1)}}(f^{(t-1)}), \quad \mathbf{W}^{(t)} = \mathbf{W}^{(t-1)} - \eta^{(t)} \nabla_{\mathbf{W}} \tilde{L}_{Z^{(t-1)}}(f^{(t-1)})$ 
end for

```

In the whole paper, we need some natural assumptions about the data and the loss.

Assumption 3.1. We assume $\mathbb{E}[k\mathbf{x}k_2] \leq B_{x1}$, $\mathbb{E}[k\mathbf{x}k_2^2] \leq B_{x2}$, $k\mathbf{x}k_2 \leq B_x$ and for any label y , we have $|y| \leq 1$. We assume the loss function $\ell(\cdot)$ is a 1-Lipschitz convex decreasing function, normalized $\ell(0) = 1$, $j\ell'(0)j \leq 1$, and $\ell(1) = 0$.

Remark 3.2. The above are natural assumptions. Most input distributions have the bounded norms required, and the typical binary classification $Y = \{-1, 1\}$ satisfies the requirement. Also, the most popular loss functions satisfy the assumption, e.g., the hinge loss and logistic loss.

3.1 Warm Up: A Simple Setting with Frozen First Layer

To illustrate some high-level intuition, we first consider a simple setting where the first layer is frozen after one gradient update, i.e., no updates to \mathbf{W} for $t \geq 2$ in Algorithm 1.

The first idea of our framework is to provide guarantees compared to the optimal in a family of networks. Here let us consider networks with specific weights for the first layer:

Definition 3.3. For some fixed $\mathbf{W} \in \mathbb{R}^{d \times m}$, $\mathbf{b} \in \mathbb{R}^m$, and a parameter B_{a2} , consider the following family of networks $F_{\mathbf{W}; \mathbf{b}; B_{a2}}$, and the optimal approximation network loss in this family:

$$F_{\mathbf{W}; \mathbf{b}; B_{a2}} := \{f_{(\mathbf{a}, \mathbf{W}; \mathbf{b})} \in F_{d,m} \mid k\mathbf{a}k_2 \leq B_{a2}\}, \quad \text{OPT}_{\mathbf{W}; \mathbf{b}; B_{a2}} := \min_{f \in F_{\mathbf{W}; \mathbf{b}; B_{a2}}} L_{\mathcal{D}}(f). \quad (3)$$

The second idea is to compare to networks using features from gradient descent. As an illustrative example, we now provide guarantees compared to networks with first layer weights $\mathbf{W}^{(1)}$ (i.e., the weights after the first gradient step):

Theorem 3.4 (Simple Setting). Assume $\tilde{L}_Z(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})})$ is L -smooth to \mathbf{a} . Let $\eta^{(t)} = \frac{1}{L}, \lambda^{(t)} = 0$, for all $t \in \{2, 3, \dots, T\}$. Training by Algorithm 1 with no updates for the first layer after the first gradient step, w.h.p., there exists $t \geq [T]$ such that

$$L_{\mathcal{D}}(f_{(\mathbf{a}^{(t)}, \mathbf{W}^{(1)}, \mathbf{b})}) \leq \text{OPT}_{\mathbf{W}^{(1)}, \mathbf{b}; B_{a2}} + O\left(\frac{L(\|\mathbf{a}^{(1)}\|_2^2 + B_{a2}^2)}{T} + \sqrt{\frac{B_{a2}^2(\|\mathbf{W}^{(1)}\|_F^2 B_x^2 + \|\mathbf{b}\|_2^2)}{n}}\right).$$

Intuitively, the theorem shows that if the weight $\mathbf{W}^{(1)}$ after a one-step gradient gives a good set of neurons in the sense that there exists a classifier on top of these neurons with low loss, then the network will learn to approximate this good classifier and achieve low loss. The proof is based on standard convex optimization and the Rademacher complexity (details in Appendix D.1).

Such an approach, while simple, has been used to obtain interesting results on network learning in existing work, which shows that $\mathbf{W}^{(1)}$ can indeed give good neurons due to the structure of the special problems considered (e.g., parities on uniform inputs [15], or polynomials on a subspace [32]). However, it is unclear whether such intuition can still yield useful guarantees for other problems. So, for our purpose of building a general framework covering more prototypical problems, the challenge is what features from gradient descent should be considered so that the family of networks for comparison can achieve a low loss on other problems. The other challenge is that we would like to consider the typical case where the first layer weights are not frozen. In the following, we will introduce the core concept of Gradient Features to address the first challenge, and stipulate proper geometric properties of Gradient Features for the second challenge.

3.2 Core Concepts in the Gradient Feature Learning Framework

Now, we will introduce the core concept in our framework, Gradient Features, and use it to build the family of networks to derive guarantees. As mentioned, we consider the setting where the first layer is not frozen. After the network learns good features, to ensure the updates in later gradient steps of the first layer are still benign for feature learning, we need some geometric conditions about the gradient features, which are measured by parameters in the definition of Gradient Features. The conditions are general enough, so that, as shown in Section 4, many prototypical problems satisfy them and the induced family of networks enjoys low loss, leading to useful guarantees. We begin by considering what features can be learned via gradients. Note that the gradient w.r.t. \mathbf{w}_i is

$$\begin{aligned} \frac{\partial L_{\mathcal{D}}(f)}{\partial \mathbf{w}_i} &= \mathbf{a}_i E_{(\mathbf{x}, y)} [\ell'(yf(\mathbf{x}))y [\sigma'(\mathbf{w}_i^T \mathbf{x} - \mathbf{b}_i)] \mathbf{x}] \\ &= \mathbf{a}_i E_{(\mathbf{x}, y)} [\ell'(yf(\mathbf{x}))y \mathbf{x} \mathbb{1}[\mathbf{w}_i^T \mathbf{x} > \mathbf{b}_i]]. \end{aligned}$$

Inspired by this, we define the following notion:

Definition 3.5 (Simplified Gradient Vector). For any $\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$, a *Simplified Gradient Vector* is

$$G(\mathbf{w}, b) := E_{(\mathbf{x}, y) \sim \mathcal{D}} [y \mathbf{x} \mathbb{1}[\mathbf{w}^T \mathbf{x} > b]]. \quad (4)$$

Remark 3.6. Note that the definition of $G(\mathbf{w}, b)$ ignores the term $\ell'(yf(\mathbf{x}))$ in the gradient, where f is the model function. In the early stage of training (or the first gradient step), $\ell'(\cdot)$ is approximately a constant, i.e., $\ell'(yf(\mathbf{x})) \approx \ell'(0)$ due to the symmetric initialization (see Equation (8)).

Definition 3.7 (Gradient Feature). For a unit vector $D \in \mathbb{R}^d$ with $kDk_2 = 1$, and a $\gamma \in (0, 1)$, a *direction neighborhood (cone)* $C_{D, \gamma}$ is defined as:

$$C_{D, \gamma} := \{\mathbf{w} \mid \langle \mathbf{w}, D \rangle / \|\mathbf{w}\|_2 > (1 - \gamma)g\}. \quad (5)$$

Gradient Feature being cones under Mixture of Gaussians data

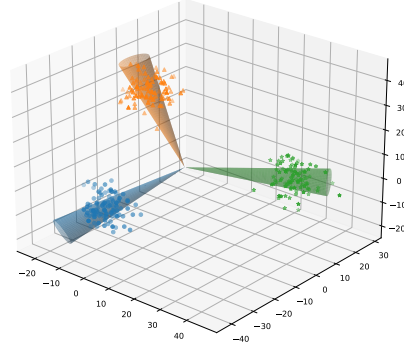


Figure 1: An illustration of Gradient Feature, i.e., Definition 3.7 with random initialization (Gaussian), under Mixture of three Gaussian clusters in 3-dimension data space with blue/green/orange color. The Gradient Feature stays in three cones, where each center of the cone aligns with the corresponding Gaussian cluster center.

Let $\mathbf{w} \in \mathbb{R}^d$, $b \in \mathbb{R}$ be random variables drawn from some distribution W, B . A Gradient Feature set with parameters p, γ, B_G is defined as:

$$S_{p, \gamma, B_G}(W, B) := \{(D, s) \mid \Pr_{\mathbf{w}; b} [G(\mathbf{w}, b) \geq c_D, \|\mathbf{g}(\mathbf{w}, b)\|_2 \leq B_G, s = b/\|\mathbf{g}\|] \geq p\}. \quad (6)$$

Remark 3.8. When clear from context, write it as S_{p, γ, B_G} . Gradient features (see Figure 1 for illustration) are simply normalized vectors D that are given (approximately) by the simplified gradient vectors. (Similarly, the normalized scalar s is given by the bias b .) To be a useful gradient feature, we require the direction to be “hit” by sufficiently large simplified gradient vectors with sufficient large probability, so as to be distinguished from noise and remain useful throughout the gradient steps. Later we will use the gradient features when W, B are the initialization distributions.

To make use of the gradient features, we consider the following family of networks using these features and with bounded norms, and will provide guarantees compared to the best in this family:

Definition 3.9 (Gradient Feature Induced Networks). *The Gradient Feature Induced Networks are:*

$$\mathcal{F}_{d; m; B_F; S} := \{f_{(\mathbf{a}; \mathbf{w}; \mathbf{b})} \in \mathcal{F}_{d; m} \mid \|\mathbf{a}_i\| \leq B_{a1}, \|\mathbf{w}_i\| \leq B_{a2}, (\mathbf{w}_i, \mathbf{b}_i/\|\mathbf{b}_i\|) \in S, \|\mathbf{b}_i\| \leq B_b\},$$

where S is some Gradient Feature set and $B_F := (B_{a1}, B_{a2}, B_b)$ are some parameters.

Remark 3.10. In above definition, the weight and bias of a neuron are simply the scalings of some item in the feature set S (for simplicity the scaling of \mathbf{w}_i is absorbed into the scaling of \mathbf{a}_i and \mathbf{b}_i).

Definition 3.11 (Optimal Approximation via Gradient Features). *The optimal approximation network and loss using Gradient Feature Induced Networks $\mathcal{F}_{d; r; B_F; S}$ are defined as:*

$$f^* := \operatorname{argmin}_{f \in \mathcal{F}_{d; r; B_F; S}} L_{\mathcal{D}}(f), \quad \text{OPT}_{d; r; B_F; S} := \min_{f \in \mathcal{F}_{d; r; B_F; S}} L_{\mathcal{D}}(f). \quad (7)$$

3.3 Provable Guarantee via Gradient Feature Learning

To obtain the guarantees, we first specify the symmetric initialization. It is convenient for the analysis and is typical in existing analysis (e.g., [7, 32, 33, 105]), though some other initialization can also work. Formally, we train a two-layer network with $4m$ neurons, $f_{(\mathbf{a}; \mathbf{w}; \mathbf{b})} \in \mathcal{F}_{d; 4m}$. We initialize $\mathbf{a}_i^{(0)}, \mathbf{w}_i^{(0)}$ from Gaussians and \mathbf{b}_i from a constant for $i \in \{1, \dots, mg\}$, and initialize the parameters for $i \in \{fm + 1, \dots, 4mg\}$ accordingly to get a zero output initial network. Specifically:

$$\begin{aligned} \text{for } i \in \{1, \dots, mg\}: \quad & \mathbf{a}_i^{(0)} \sim N(0, \sigma_a^2), \mathbf{w}_i^{(0)} \sim N(0, \sigma_w^2 I), \mathbf{b}_i = \mathbf{b}, \\ \text{for } i \in \{fm + 1, \dots, 2mg\}: \quad & \mathbf{a}_i^{(0)} = \mathbf{a}_{i-m}^{(0)}, \mathbf{w}_i^{(0)} = \mathbf{w}_{i-m}^{(0)}, \mathbf{b}_i = \mathbf{b}_{i-m}, \\ \text{for } i \in \{2m + 1, \dots, 4mg\}: \quad & \mathbf{a}_i^{(0)} = \mathbf{a}_{i-2m}^{(0)}, \mathbf{w}_i^{(0)} = \mathbf{w}_{i-2m}^{(0)}, \mathbf{b}_i = \mathbf{b}_{i-2m}, \end{aligned} \quad (8)$$

where $\sigma_a^2, \sigma_w^2, \mathbf{b} > 0$ are hyper-parameters. After initialization, \mathbf{a}, \mathbf{W} are updated as in Algorithm 1. We are now ready to present our main result in the framework.

Theorem 3.12 (Main Result). *Assume Assumption 3.1. For any $\epsilon, \delta \in (0, 1)$, if $m \geq e^d$ and*

$$\begin{aligned} m &= \left(\frac{1}{p\epsilon^4} \left(r B_{a1} B_{x1} \sqrt{\frac{B_b}{B_G}} \right)^4 + \frac{1}{\delta} + \frac{1}{p} \left(\log \left(\frac{r}{\delta} \right) \right)^2 \right), \\ T &= \left(\frac{1}{\epsilon} \left(\frac{p}{r B_{a2} B_b B_{x1}} + m \mathbf{b} \right) \left(\frac{p}{B_b B_G} + \frac{1}{B_{x1} (mp)^{\frac{1}{4}}} \right) \right), \\ \frac{n}{\log n} &= \left(\frac{m^3 p B_x^2 B_{a2}^4 B_b}{\epsilon^2 r^2 B_{a1}^2 B_G} + \frac{(mp)^{\frac{1}{2}} B_{x2}}{B_b B_G} + \frac{B_x^2}{B_{x2}} + \frac{1}{p} + \left(\frac{1}{B_G^2} + \frac{1}{B_{x1}^2} \right) \frac{B_{x2}}{j\ell'(0)^2} + \frac{Tm}{\delta} \right), \end{aligned}$$

then with initialization (8) and proper hyper-parameter values, we have with probability $1 - \delta$ over the initialization and training samples, there exists $t \geq [T]$ in Algorithm 1 with:

$$\Pr[\operatorname{sign}(f^{(t)}(\mathbf{x})) \neq y] \leq L_{\mathcal{D}}(f^{(t)})$$

$$\text{OPT}_{d; r; B_F; S_{p, \gamma, B_G}} + r B_{a1} B_{x1} \sqrt{2\gamma + O\left(\frac{p}{B_G j\ell'(0) j n^{\frac{1}{2}}}\right)} + \epsilon.$$

Intuitively, the theorem shows when a data distribution admits a small approximation error by some “ground-truth” network with r neurons using gradient features from $S_{p; :B_G}$ (i.e., a small optimal approximate loss $\text{OPT}_{d;r;B_F;S_{p; :B_G}}$), the gradient descent training can successfully learn good neural networks with sufficiently many m neurons.

Now we discuss the requirements and the error guarantee. Viewing boundedness parameters $B_{a1}, B_{\chi1}$ etc. as constants, then the number m of neurons learned is roughly $\sim \left(\frac{r^4}{p}\right)$, a polynomial overparameterization compared to the “ground-truth” network. The proof shows that such an overparameterization is needed such that some neurons can capture the gradient features given by gradient descent. This is consistent with existing analysis about overparameterization network learning, and also consistent with existing empirical observations.

The error bound consists of three terms. The last term ϵ can be made arbitrarily small, while the other two depend on the concrete data distribution. Specifically, with larger r and γ , the second term increases. While the first term (the optimal approximation loss) decreases, since a larger r means a larger “ground-truth” network family, and a larger γ means a larger Gradient Feature set $S_{p; :B_G}$. So, there is a trade-off between these two terms. When we later apply the framework to concrete problems (e.g., mixtures of Gaussians, parity functions), we will show that depending on the specific data distribution, we can choose the proper values for r, γ to make the error small. This then leads to error guarantees for the concrete problems and demonstrates the unifying power of the framework. Please refer to Appendix D.3 for more discussion about our problem setup and our core concept, e.g., parameter choice, early stopping, the role of s , activation functions, and so on.

Proof Sketch. The intuition in the proof of Theorem 3.12 is closely related to the notion of Gradient Features. First, the gradient descent will produce gradients that approximate the features in $S_{p; :B_G}$. Then, the gradient descent update gives a good set of neurons, such that there exists an accurate classifier using these neurons with loss comparable to the optimal approximation loss. Finally, the training will learn to approximate the accurate classifier, resulting in the desired error guarantee. The complete proof is in Appendix D (the population version in Appendix D.2 and the empirical version in Appendix D.4), including the proper values for hyper-parameters such as $\eta^{(t)}$ in Theorem D.17. Below, we briefly sketch the key ideas and omit the technical details.

We first show that a large subset of neurons has gradients at the first step as good features. (The claim can be extended to multiple steps; for simplicity, we follow existing work (e.g., [33, 105]) and present only the first step.) Let r_i denote the gradient of the i -th neuron $r_{\mathbf{w}_i, L_{\mathcal{D}}(f_{(0)})}$. Denote the subset of neurons with nice gradients approximating feature (D, s) as:

$$G_{(D,s):Nice} := \left\{ i \in [2m] : s = \mathbf{b}_i / \|\mathbf{b}_i\|, \|\mathbf{r}_i\| > (1 - \gamma) k r_i k_2, k r_i k_2 \leq \left| \mathbf{a}_i^{(0)} \right|_{B_G} \right\}. \quad (9)$$

Lemma 3.13 (Feature Emergence). *For any r size subset $f(D_1, s_1), \dots, (D_r, s_r) \in S_{p; :B_G}$, with probability at least $1 - re^{-\frac{mp}{4}}$, for all $j \in [r]$, we have $|G_{(D_j, s_j):Nice}| \geq \frac{mp}{4}$.*

This is because $r_i = \ell'(0) \mathbf{a}_i^{(0)} \mathbb{E}_{(\mathbf{x}, y)} \left[y \sigma' \left[\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right] \mathbf{x} \right] = \ell'(0) \mathbf{a}_i^{(0)} G(\mathbf{w}_i^{(0)}, \mathbf{b}_i)$. Now consider $s_j = +1$ (the case -1 is similar). Since \mathbf{w}_i is initialized by Gaussians, by r_i 's connection to Gradient Features, we can see that for all $i \in [m]$, $\Pr \left[i \in G_{(D_j, +1):Nice} \right] \geq \frac{\rho}{2}$. The lemma follows from concentration via a large enough m , i.e., sufficient overparameterization. The gradients allow obtaining a set of neurons approximating the “ground-truth” network with comparable loss:

Lemma 3.14 (Existence of Good Networks). *For any $\delta \in (0, 1)$, with proper hyper-parameter values, with probability at least $1 - \delta$, there is \mathbf{a} such that $\|\mathbf{a}\|_0 = O\left(r^{\frac{\rho}{mp}}\right)$ and $f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}(\mathbf{x}) = \sum_{i=1}^{4m} \mathbf{a}_i \sigma \left(\left\langle \mathbf{w}_i^{(1)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right)$ satisfies*

$$L_{\mathcal{D}}(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}) \leq \text{OPT}_{d;r;B_F;S_{p; :B_G}} + \frac{\rho}{2r} B_{a1} B_{\chi1} \left(\frac{\rho}{\gamma} + \sqrt{\frac{2B_b}{\rho \frac{mp}{B_G}}} \right).$$

Given the good set of neurons, we finally show that the remaining gradient steps can learn an accurate classifier. Intuitively, with small step sizes $\eta^{(t)}$, the weights of the first layer \mathbf{w}_i do not change too much (stay in a neighborhood) while the second layer weights grow, and thus the learning is similar to convex learning using the good set of neurons. Technically, we adopt the online convex optimization analysis (Theorem D.5) in [33] to get the final loss guarantee in Theorem 3.12.

4 Applications in Special Cases

In this section we will apply the gradient feature learning framework to some specific problems, corresponding to concrete data distributions D . We primarily focus on prototypical problems for analyzing feature learning in networks. We will present here the results for mixtures of Gaussians and parity functions, and include the complete proofs and some other results in Appendix E.

4.1 Mixtures of Gaussians

Mixtures of Gaussians are among the most fundamental and widely used statistical models. Recently, it has been used to study neural network learning, in particular, the effect of gradient descent for feature learning of two-layer neural networks and the advantage over fixed feature methods [46, 99].

Data Distributions. We follow notations from [99]. The data are from a mixture of r high-dimensional Gaussians, and each Gaussian is assigned to one of two possible labels in $Y = \{-1, +1\}$. Let $S(y) \subseteq [r]$ denote the set of indices of Gaussians associated with the label y . The data distribution is then: $q(\mathbf{x}, y) = q(y)q(\mathbf{x}|y)$, $q(\mathbf{x}|y) = \sum_{j \in S(y)} p_j \mathcal{N}_j(\mathbf{x})$, where $\mathcal{N}_j(\mathbf{x})$ is a multivariate normal distribution with mean μ_j , covariance Σ_j , and p_j are chosen such that $q(\mathbf{x}, y)$ is correctly normalized. We will make some assumptions about the Gaussians, for which we first introduce some notations.

$$D_j := \frac{\mu_j}{k_{\mu_j} k_2}, \quad \mu_j := \mu_j / \sqrt{d}, \quad B_1 := \min_{j \in [r]} k_{\mu_j} k_2, \quad B_2 := \max_{j \in [r]} k_{\mu_j} k_2, \quad p_B := \min_{j \in [r]} p_j.$$

Assumption 4.1. Let $\delta, \tau \in (0, 1]$ be a parameter that will control our final error guarantee. Assume

- Equiprobable labels: $q(-1) = q(+1) = 1/2$.
- For all $j \in [r]$, $\Sigma_j = \sigma_j I_{d \times d}$. Let $\sigma_B := \max_{j \in [r]} \sigma_j$ and $\sigma_{B^+} := \max_{j \in [r]} \sigma_j$.
- $r \geq 2d$, $p_B \geq \frac{1}{2d}$, $\left(1/d + \sqrt{\tau \sigma_{B^+}^2 \log d/d}\right) \leq B_1 \leq B_2 \leq d$.
- The Gaussians are well-separated: for all $i \neq j \in [r]$, we have $\|\mu_i - \mu_j\| \geq \theta$, where $\theta \geq \min \left\{ \frac{1}{2r}, \frac{\sigma_{B^+}}{B_2} \sqrt{\frac{\log d}{d}} \right\}$.

Remark 4.2. The first two assumptions are for simplicity; they can be relaxed. We can generalize our analysis to the mixture of Gaussians with unbalanced label probabilities and general covariances. The third assumption is to make sure that each Gaussian has a good amount of probability mass to be learned. The remaining assumptions are to make sure that the Gaussians are well-separated and can be distinguished by the learning algorithm.

We are now ready to apply the framework to these data distributions, for which we only need to compute the Gradient Feature set and the corresponding optimal approximation loss.

Lemma 4.3 (Mixtures of Gaussians: Gradient Features). $(D_j, +1) \in S_{p_j, \sigma_{B^+}}$ for all $j \in [r]$, where

$$p = \frac{B_1}{\tau \log d \sigma_{B^+} \sqrt{d} \sqrt{B_2^2 - B_1^2}}, \quad \gamma = \frac{1}{d^{0.9 - 1.5}}, \quad B_G = p_B B_1 \sqrt{d} = O\left(\frac{\sigma_{B^+}}{d^{0.9}}\right).$$

Let $f^*(\mathbf{x}) = \sum_{j=1}^r \frac{y(j)}{\sqrt{\log d}} \frac{1}{\sigma_{B^+}} \left[\sigma \left(\frac{1}{\tau \log d \sigma_{B^+}} \langle \mathbf{x}, \mu_j \rangle \right) \right]$ whose hinge loss is at most $\frac{3}{d} + \frac{4}{d^{0.9} \sqrt{\log d}}$.

Given the values on gradient feature parameters p, γ, B_G and the optimal approximation loss $\text{OPT}_{d,r;B_F;S_{p_j; \sigma_{B^+}}; B_G}$, the framework immediately leads to the following guarantee:

Theorem 4.4 (Mixtures of Gaussians: Main Result). Assume Assumption 4.1. For any $\epsilon, \delta \in (0, 1)$, when Algorithm 1 uses hinge loss with

$$m = \text{poly} \left(\frac{1}{\delta}, \frac{1}{\epsilon}, d \sqrt{B_2^2 - B_1^2}, r, \frac{1}{p_B} \right) e^d, \quad T = \text{poly}(m), \quad n = \text{poly}(m)$$

and proper hyper-parameters, then with probability at least $1 - \delta$, there exists $t \in [T]$ such that

$$\Pr[\text{sign}(f^{(t)}(\mathbf{x})) \neq y] \leq \frac{\sqrt{2r}}{d^{0.4 - 0.8}} + \epsilon.$$

The theorem shows that gradient descent can learn to a small error via learning the gradient features, given proper hyper-parameters. In particular, we need sufficient overparameterization (a sufficiently large number m of neurons). When σ_{B+}^2/B^2_1 is a constant which is the prototypical interesting case, and we choose a constant τ , then m is polynomial in the key parameters $\frac{1}{\rho_B}, \frac{1}{\rho_B}, d, r, \frac{1}{\rho_B}$, and the error bound is inverse polynomial in d . The complete proof is given in Appendix E.2.

[46] studies (almost) linear separable cases while our setting includes non-linear separable cases, e.g., XOR. [99] mainly studies neural network classification on 4 Gaussian clusters with XOR structured labels, while our setting is much more general, e.g., our cluster number can extend up to $2d$.

4.1.1 Mixtures of Gaussians: Beyond the Kernel Regime

As discussed in the introduction, it is important for the analysis to go beyond fixed feature methods such as NTK (i.e., the kernel regime), so as to capture the feature learning ability which is believed to be the key factor for the empirical success. We first review the fixed feature methods. Following [33], suppose \mathcal{X} is a data-independent feature mapping of dimension N with bounded features, i.e., $\mathcal{X} : \mathcal{X} \rightarrow [1, 1]^N$. For $B > 0$, the family of linear models on \mathcal{X} with bounded norm B is $H_B = \{h(\mathbf{x}) : h(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle, \|\mathbf{w}\|_2 \leq B\}$. This can capture linear models on fixed finite-dimensional feature maps, e.g., NTK, and also infinite dimensional feature maps, e.g., kernels like RBF, that can be approximated by feature maps of polynomial dimensions [64, 98, 105].

Our framework indeed goes beyond fixed features and shows features from gradients are more powerful than features from random initialization, e.g., NTK. Our framework can show the advantage of network learning over kernel methods under the setting of [99] (4 Gaussian clusters with XOR structured labels). For large enough d , our framework only needs roughly $(\log d)$ neurons and $(\log d)^2$ samples to achieve arbitrary small constant error (see Theorem E.18 when $\sigma_B = 1$), while fixed feature methods need (d^2) features and (d^2) samples to achieve nontrivial errors (as proved in [99]). Moreover, [99] uses ODE to simulate the optimization process for the 2-layer networks learning XOR-shaped Gaussian mixture with (1) neurons and gives convincing evidence that (d) samples is enough to learn it, yet they do not give a rigorous convergence guarantee for this problem. We successfully derive a convergence guarantee and we require a much smaller sample size $(\log d)^2$. For the proof (detailed in Appendix E.3), we only need to calculate the p, γ, B_G of the data distribution carefully and then inject these numbers into Theorem 3.12.

4.2 Parity Functions

Parity functions are a canonical family of learning problems in computational learning theory, usually for showing theoretical computational barriers [103]. The typical sparse parties over d -dim binary inputs $\phi \in \{0, 1\}^d$ are $\prod_{i \in A} \phi_i$ where $A \subseteq [d]$ is a subset of dimensions. Recent studies have shown that when the distribution of inputs ϕ has structures rather than uniform, neural networks can perform feature learning and finally learn parity functions with a small error, while methods without feature learning, e.g. NTK, cannot achieve as good results [33, 76, 105]. Thus, this has been a prototypical setting for studying feature learning phenomena in networks. Here we consider a generalization of this problem and show that our framework can show successful learning via gradient descent.

Data Distributions. Suppose $\mathbf{M} \in \mathbb{R}^{d \times D}$ is an unknown dictionary with D columns that can be regarded as patterns. For simplicity, assume $d = D$ and \mathbf{M} is orthonormal. Let $\phi \in \mathbb{R}^d$ be a hidden representation vector. Let $A \subseteq [D]$ be a subset of size rk corresponding to the class relevant patterns and r is an odd number. Then the input is generated by $\mathbf{M}\phi$, and some function on ϕ_A generates the label. WLOG, let $A = \{r1, \dots, rk\}, A^c = \{rk+1, \dots, dg\}$. Also, we split A such that for all $j \in [r], A_j = \{(j-1)k+1, \dots, jk\}$. Then the input \mathbf{x} and the class label y are given by:

$$\mathbf{x} = \mathbf{M}\phi, y = g^*(\phi_A) = \text{sign}\left(\sum_{j \in [r]} \text{XOR}(\phi_{A_j})\right), \quad (10)$$

where g^* is the ground-truth labeling function mapping from \mathbb{R}^{rk} to $Y = \{-1, 1\}$, ϕ_A is the sub-vector of ϕ with indices in A , and $\text{XOR}(\phi_{A_j}) = \prod_{i \in A_j} \phi_i$ is the parity function. We still need to specify the distribution X of ϕ , which determines the structure of the input distribution:

$$X := (1 - 2rp_A)X_U + \sum_{j \in [r]} p_A(X_{j,+} + X_{j,-}). \quad (11)$$

For all corresponding ϕ_{A^c} in X , we have $\phi_{A^c} \perp A^\perp$, independently: $\phi_{A^c} = \begin{cases} +1, & \text{w.p. } p_o \\ 1, & \text{w.p. } p_o \\ 0, & \text{w.p. } 1 - 2p_o \end{cases}$,

where p_o controls the signal noise ratio: if p_o is large, then there are many nonzero entries in A^\perp which are noise interfering with the learning of the ground-truth labeling function on A . For corresponding ϕ_A , any $j \in [r]$, we have

- In $X_{j,+}$, $\phi_{A_j} = [+1, +1, \dots, +1]^\top$ and $\phi_{A \setminus A_j}$ only have zero elements.
- In $X_{j,-}$, $\phi_{A_j} = [-1, -1, \dots, -1]^\top$ and $\phi_{A \setminus A_j}$ only have zero elements.
- In X_U , we have ϕ_A draw from $\mathcal{F}+1, \quad 1/d^{r \cdot k}$ uniformly.

In short, we have r parity functions each corresponding to a block of k dimensions; $X_{j,+}$ and $X_{j,-}$ stands for the component providing a strong signal for the j -th parity; X_U corresponds to uniform distribution unrelated to any parity and providing weak learning signal; A^\perp is the noise part. The label depends on the sum of the r parity functions.

Assumption 4.5. Let $\delta = \tau/d$ be a parameter that will control our final error guarantee. Assume k is an odd number and: $k \geq (\tau \log d)$, $d \geq rk + (\tau r \log d)$, $p_o = O\left(\frac{rk}{d-rk}\right)$, $p_A \geq \frac{1}{d}$.

Remark 4.6. We set up the problem to be more general than the parity function learning in existing work. If $r = 1$, the labeling function reduces to the traditional k -sparse parties of d bits. The assumptions require k, d , and p_A to be sufficiently large so as to provide enough large signals for learning. Note that when $k = \frac{d}{16}, r = 1, p_o = \frac{1}{2}$, our analysis also holds, which shows our framework is beyond the kernel regime (discuss in detail in Section 4.2.1).

To apply our framework, again we only need to compute the Gradient Feature set and the corresponding optimal loss. We first define the Gradient Features: For all $j \in [r]$, let $D_j = \frac{\sum_{i \in A_j} \mathbf{M}_i}{\|\sum_{i \in A_j} \mathbf{M}_i\|_2}$.

Lemma 4.7 (Parity Functions: Gradient Features). We have $(D_j, +1), (D_j, -1) \in S_{p_o, \delta, B_G}$ for all $j \in [r]$, where

$$p = \left(\frac{1}{\tau r \log d} \frac{1}{d} \right)^{\gamma}, \quad \gamma = \frac{1}{d-2}, \quad B_G = \frac{\rho_{\bar{k}}}{k p_A} = O\left(\frac{\rho_{\bar{k}}}{d}\right). \quad (12)$$

With gradient features from S_{p_o, δ, B_G} , let $f^*(\mathbf{x}) = \sum_{j=1}^r \sum_{i=0}^k (-1)^{i+1} \frac{\rho_{\bar{k}}}{k} \left[\sigma\left(h D_j, \mathbf{x} \cdot \frac{2i-k-1}{\sqrt{k}}\right) + \sigma\left(h D_j, \mathbf{x} \cdot \frac{2i-k+1}{\sqrt{k}}\right) \right]$ whose hinge loss is 0.

Above, we show that D_j is the ‘‘indicator function’’ for the subset A_j so that we can build the optimal neural network based on such directions. Given the values on gradient feature parameters and the optimal approximation loss, the framework immediately leads to the following guarantee:

Theorem 4.8 (Parity Functions: Main Result). Assume Assumption 4.5. For any $\epsilon, \delta \in (0, 1)$, when Algorithm 1 uses hinge loss with

$$m = \text{poly}\left(\frac{1}{\delta}, \frac{1}{\epsilon}, d^{\gamma}, k, \frac{1}{p_A}\right) e^d, \quad T = \text{poly}(m), \quad n = \text{poly}(m)$$

and proper hyper-parameters, then with probability at least $1 - \delta$, there exists $t \geq T$ such that

$$\Pr[\text{sign}(f_t(\mathbf{x})) \neq y] \leq \frac{3r \rho_{\bar{k}}}{d^{(-3) \cdot 2}} + \epsilon.$$

The theorem shows that gradient descent can learn to a small error in this problem. We also need sufficient overparameterization: When r is a constant (e.g., $r = 1$ in existing work), and we choose a constant τ , m is polynomial in $\frac{1}{\delta}, \frac{1}{\epsilon}, d, k, \frac{1}{p_A}$, and the error bound is inverse polynomial in d . The proof is in Appendix E.4. Our setting is more general than that in [33, 76] which corresponds to $\mathbf{M} = I, r = 1, p_A = \frac{1}{4}, p_o = \frac{1}{2}$. [105] study single index learning, where one feature direction is enough for a two-layer network to recover the label, while our setting considers r directions D_1, \dots, D_r , so the network needs to learn multiple directions to get a small error.

4.2.1 Parity Functions: Beyond the Kernel Regime

Again, we show that our framework indeed goes beyond fixed features under parity functions. Our problem setting in Section 4.2 is general enough to include the problem setting in [33]. Their lower bound for fixed feature methods directly applies to our case and leads to the following:

Proposition 4.9. *There exists a data distribution in the parity learning setting in Section 4.2 with $\mathbf{M} = I, r = 1, p_A = \frac{1}{4}, k = \frac{d}{16}, p_o = \frac{1}{2}$, such that all $h \geq H_B$ have hinge-loss at least $\frac{1}{2} - \frac{\sqrt{NB}}{2^k\sqrt{2}}$.*

This means to get an inverse-polynomially small loss, fixed feature models need to have an exponentially large size, i.e., either the number of features N or the norm B needs to be exponential in k . In contrast, Theorem 4.8 shows our framework guarantees a small loss with a polynomially large model, runtime, and sample complexity. Clearly, our framework is beyond the fixed feature methods.

Parities on Uniform Inputs. When $r = 1, p_A = 0$, our problem setting will degenerate to the classic sparse parity function on a uniform input distribution. This has also been used for analyzing network learning [16]. For this case, our framework can get a $k2^{O(k)} \log(k)$ network width bound and a $O(d^k)$ sample complexity bound, matching those in [16]. This then again confirms the advantage of network learning over kernel methods that requires $d^{O(k)}$ dimensions as shown in [16]. See the full statement in Theorem E.31, details in Appendix E.5, and alternative analysis in Appendix E.6.

5 Further Implications and Conclusion

Our general framework sheds light on several interesting phenomena in NN learning observed in practice. Feature learning beyond the kernel regime has been discussed in Section 4.1.1 and Section 4.2.1. Here we discuss the LTH and defer more implications such as simplicity bias, learning over different data distributions, and new perspectives about roadmaps forward in Appendix C.

Lottery Ticket Hypothesis (LTH). Another interesting phenomenon is the LTH [41]: randomly-initialized networks contain subnetworks that when trained in isolation reach test accuracy comparable to the original network in a similar number of iterations. Later studies (e.g., [42]) show that LTH is more stable when subnetworks are found in the network after a few gradient steps.

Our framework provides an explanation for two-layer networks: the lottery ticket subnetwork contains exactly those neurons whose gradient feature approximates the weights of the “ground-truth” network f^* ; they may not exist at initialization but can be found after the first gradient step. More precisely, Lemma 3.14 shows that after the first gradient step, there is a *sparse* second-layer weight \mathbf{a} with $k\mathbf{a}k_0 = O\left(r^{\frac{1}{m}}\frac{1}{mp}\right)$, such that using this weight on the hidden neurons gives a network with a small loss. Let U be the support of \mathbf{a} . Equivalently, there is a small-loss subnetwork f^U with only neurons in U and with second-layer weight \mathbf{a}_U on these neurons. Following the same proof of Theorem 3.12:

Proposition 5.1. *In the same setting of Theorem 3.12 but only considering the subnetwork supported on U after the first gradient step, with the same requirements on m and T , with proper hyperparameter values, we have the same guarantee: with probability $1 - \delta$, there is $t \geq [T]$ with*

$$\Pr[\text{sign}(f_{(t)}^U)(\mathbf{x}) \neq y] \leq \text{OPT}_{d;r;B_F;S_P;B_G} + rB_{a1}B_{x1} \sqrt{2\gamma + O\left(\frac{\sqrt{B_{x2} \log n}}{B_G \sqrt{n}}\right)} + \epsilon.$$

This essentially formally proves LTH for two-layer networks, showing (a) the existence of the winning lottery subnetwork and (b) that gradient descent on the subnetwork can learn to similar loss in similar runtime as on the whole network. In particular, (b) is novel and not analyzed in existing work.

We provide our work’s broader impacts and limitations (e.g., statement of recovering existing results and some failure cases beyond our framework) in Appendix A and Appendix B respectively.

Conclusion. We propose a general framework for analyzing two-layer neural network learning by gradient descent and show that it can lead to provable guarantees for several prototypical problem settings for analyzing network learning. In particular, our framework goes beyond fixed feature methods, e.g., NTK. It sheds light on several interesting phenomena in NN learning, e.g., the lottery ticket hypothesis and simplicity bias. Future directions include: (1) How to extend the framework to deeper networks? (2) While the current framework focuses on the gradient features in the early gradient steps, whether feature learning also happens in later steps and if so how to formalize that?

Acknowledgements

The work is partially supported by Air Force Grant FA9550-18-1-0166, the National Science Foundation (NSF) Grants 2008559-IIS, 2023239-DMS, and CCF-2046710.

References

- [1] 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. In *Conference on Learning Theory*. PMLR, 2022.
- [2] Emmanuel Abbe, Samy Bengio, Elisabetta Cornacchia, Jon Kleinberg, Aryo Lotfi, Maithra Raghu, and Chiyuan Zhang. Learning to reason with neural networks: Generalization, unseen data and boolean measures. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [3] Shunta Akiyama and Taiji Suzuki. On learnability via gradient method for two-layer relu neural networks in teacher-student setting. In *International Conference on Machine Learning*, pages 152–162. PMLR, 2021.
- [4] Shunta Akiyama and Taiji Suzuki. Excess risk of two-layer reLU neural networks in teacher-student settings and its superiority to kernel methods. In *The Eleventh International Conference on Learning Representations*, 2023.
- [5] Zeyuan Allen-Zhu and Yuanzhi Li. What can resnet learn efficiently, going beyond kernels? In *Advances in Neural Information Processing Systems*, 2019.
- [6] Zeyuan Allen-Zhu and Yuanzhi Li. Backward feature correction: How deep learning performs deep learning. *arXiv preprint arXiv:2001.04413*, 2020.
- [7] Zeyuan Allen-Zhu and Yuanzhi Li. Feature purification: How adversarial training performs robust deep learning. In *2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS)*, pages 977–988. IEEE, 2022.
- [8] Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparameterized neural networks, going beyond two layers. In *Advances in neural information processing systems*, 2019.
- [9] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In *International Conference on Machine Learning*, 2019.
- [10] Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. A convergence analysis of gradient descent for deep linear neural networks. In *International Conference on Learning Representations*, 2018.
- [11] Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *International Conference on Machine Learning*, pages 322–332. PMLR, 2019.
- [12] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Ruslan Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. *arXiv preprint arXiv:1904.11955*, 2019.
- [13] 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.
- [14] Yu Bai and Jason D Lee. Beyond linearization: On quadratic and higher-order approximation of wide neural networks. In *International Conference on Learning Representations*, 2019.
- [15] 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.

- [16] Boaz Barak, Benjamin L Edelman, Surbhi Goel, Sham M Kakade, Cyril Zhang, et al. Hidden progress in deep learning: Sgd learns parities near the computational limit. In *Advances in Neural Information Processing Systems*, 2022.
- [17] Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 2020.
- [18] 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.
- [19] Avrim Blum and Ronald L Rivest. Training a 3-node neural network is np-complete. In *Advances in neural information processing systems*, pages 494–501, 1989.
- [20] Yuan Cao and Quanquan Gu. Generalization bounds of stochastic gradient descent for wide and deep neural networks. *Advances in Neural Information Processing Systems*, 2019.
- [21] Yuan Cao, Zhiying Fang, Yue Wu, Ding-Xuan Zhou, and Quanquan Gu. Towards understanding the spectral bias of deep learning, 2020.
- [22] Yuan Cao, Zixiang Chen, Misha Belkin, and Quanquan Gu. Benign overfitting in two-layer convolutional neural networks. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [23] Niladri S Chatterji, Philip M Long, and Peter L Bartlett. When does gradient descent with logistic loss find interpolating two-layer networks? *Journal of Machine Learning Research*, pages 1–48, 2021.
- [24] Minshuo Chen, Haoming Jiang, Wenjing Liao, and Tuo Zhao. Efficient approximation of deep relu networks for functions on low dimensional manifolds. *Advances in neural information processing systems*, 32:8174–8184, 2019.
- [25] Minshuo Chen, Haoming Jiang, Wenjing Liao, and Tuo Zhao. Nonparametric regression on low-dimensional manifolds using deep relu networks: Function approximation and statistical recovery. *arXiv preprint arXiv:1908.01842*, 2019.
- [26] Minshuo Chen, Yu Bai, Jason D Lee, Tuo Zhao, Huan Wang, Caiming Xiong, and Richard Socher. Towards understanding hierarchical learning: Benefits of neural representations. *arXiv preprint arXiv:2006.13436*, 2020.
- [27] Zhengdao Chen, Eric Vanden-Eijnden, and Joan Bruna. On feature learning in neural networks with global convergence guarantees. In *International Conference on Learning Representations*, 2022.
- [28] 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*, 31, 2018.
- [29] Lenaic Chizat and Francis Bach. A note on lazy training in supervised differentiable programming. *arXiv preprint arXiv:1812.07956*, 2018.
- [30] Lenaic Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss. In *Conference on Learning Theory*. PMLR, 2020.
- [31] Lenaic Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. In *Advances in Neural Information Processing Systems*, 2019.
- [32] Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi. Neural networks can learn representations with gradient descent. In *Conference on Learning Theory*. PMLR, 2022.
- [33] Amit Daniely and Eran Malach. Learning parities with neural networks. *Advances in Neural Information Processing Systems*, 33, 2020.
- [34] Amit Daniely and Gal Vardi. Hardness of learning neural networks with natural weights. *Advances in Neural Information Processing Systems*, 33:930–940, 2020.

- [35] Amit Daniely, Nathan Srebro, and Gal Vardi. Efficiently learning neural networks: What assumptions may suffice? *arXiv preprint arXiv:2302.07426*, 2023.
- [36] Zhiyan Ding, Shi Chen, Qin Li, and Stephen J Wright. Overparameterization of deep resnet: zero loss and mean-field analysis. *The Journal of Machine Learning Research*, 2022.
- [37] Xialiang Dou and Tengyuan Liang. Training neural networks as learning data-adaptive kernels: Provable representation and approximation benefits. *Journal of the American Statistical Association*, 2020.
- [38] Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In *International Conference on Machine Learning*, 2019.
- [39] Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. In *International Conference on Learning Representations*, 2018.
- [40] Yu Feng and Yuhai Tu. Phases of learning dynamics in artificial neural networks: in the absence or presence of mislabeled data. *Machine Learning: Science and Technology*, 2021.
- [41] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2018.
- [42] Jonathan Frankle, Gintare Karolina Dziugaite, Daniel M Roy, and Michael Carbin. Stabilizing the lottery ticket hypothesis. *arXiv preprint arXiv:1903.01611*, 2019.
- [43] Spencer Frei and Quanquan Gu. Proxy convexity: A unified framework for the analysis of neural networks trained by gradient descent. *Advances in Neural Information Processing Systems*, 34, 2021.
- [44] Spencer Frei, Yuan Cao, and Quanquan Gu. Provable generalization of sgd-trained neural networks of any width in the presence of adversarial label noise. *arXiv preprint arXiv:2101.01152*, 2021.
- [45] Spencer Frei, Niladri S Chatterji, and Peter L Bartlett. Random feature amplification: Feature learning and generalization in neural networks. *arXiv preprint arXiv:2202.07626*, 2022.
- [46] Spencer Frei, Gal Vardi, Peter L Bartlett, Nathan Srebro, and Wei Hu. Implicit bias in leaky relu networks trained on high-dimensional data. *arXiv preprint arXiv:2210.07082*, 2022.
- [47] Spencer Frei, Gal Vardi, Peter L Bartlett, and Nathan Srebro. Benign overfitting in linear classifiers and leaky relu networks from kkt conditions for margin maximization. *arXiv preprint arXiv:2303.01462*, 2023.
- [48] Spencer Frei, Gal Vardi, Peter L Bartlett, and Nathan Srebro. The double-edged sword of implicit bias: Generalization vs. robustness in relu networks. *arXiv preprint arXiv:2303.01456*, 2023.
- [49] Mario Geiger, Stefano Spigler, Arthur Jacot, and Matthieu Wyart. Disentangling feature and lazy training in deep neural networks. *Journal of Statistical Mechanics: Theory and Experiment*, page 113301, 2020.
- [50] Mario Geiger, Leonardo Petrini, and Matthieu Wyart. Landscape and training regimes in deep learning. *Physics Reports*, 924:1–18, 2021.
- [51] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Limitations of lazy training of two-layers neural networks. *arXiv preprint arXiv:1906.08899*, 2019.
- [52] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. When do neural networks outperform kernel methods? In *Advances in Neural Information Processing Systems*, 2020.

- [53] Gauthier Gidel, Francis Bach, and Simon Lacoste-Julien. Implicit regularization of discrete gradient dynamics in linear neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- [54] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Computer Vision and Pattern Recognition*, 2014.
- [55] Sebastian Goldt, Madhu Advani, Andrew M Saxe, Florent Krzakala, and Lenka Zdeborová. Dynamics of stochastic gradient descent for two-layer neural networks in the teacher-student setup. *Advances in neural information processing systems*, 32, 2019.
- [56] Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias in terms of optimization geometry. In *International Conference on Machine Learning*, pages 1832–1841. PMLR, 2018.
- [57] Boris Hanin and Mihai Nica. Finite depth and width corrections to the neural tangent kernel. In *International Conference on Learning Representations*, 2019.
- [58] Jiaoyang Huang and Horng-Tzer Yau. Dynamics of deep neural networks and neural tangent hierarchy. In *International conference on machine learning*, pages 4542–4551. PMLR, 2020.
- [59] Arthur Jacot. Implicit bias of large depth networks: a notion of rank for nonlinear functions. In *The Eleventh International Conference on Learning Representations*, 2023.
- [60] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, 2018.
- [61] Ziwei Ji and Matus Telgarsky. The implicit bias of gradient descent on nonseparable data. In *Conference on Learning Theory*, pages 1772–1798. PMLR, 2019.
- [62] Ziwei Ji and Matus Telgarsky. Polylogarithmic width suffices for gradient descent to achieve arbitrarily small test error with shallow relu networks. In *International Conference on Learning Representations*, 2019.
- [63] Ziwei Ji and Matus Telgarsky. Directional convergence and alignment in deep learning. *Advances in Neural Information Processing Systems*, 33:17176–17186, 2020.
- [64] Pritish Kamath, Omar Montasser, and Nathan Srebro. Approximate is good enough: Probabilistic variants of dimensional and margin complexity. In *Conference on Learning Theory*, 2020.
- [65] Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear convergence of gradient and proximal-gradient methods under the polyak-łojasiewicz condition. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 795–811. Springer, 2016.
- [66] Jonas Moritz Kohler and Aurelien Lucchi. Sub-sampled cubic regularization for non-convex optimization. In *International Conference on Machine Learning*. PMLR, 2017.
- [67] Guy Kornowski, Gilad Yehudai, and Ohad Shamir. From tempered to benign overfitting in relu neural networks. *arXiv preprint arXiv:2305.15141*, 2023.
- [68] Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S Schoenholz, Jeffrey Pennington, and Jascha Sohl-Dickstein. Deep neural networks as gaussian processes. In *International Conference on Learning Representations*, 2018.
- [69] Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. *Advances in neural information processing systems*, 2019.
- [70] Jaehoon Lee, Samuel Schoenholz, Jeffrey Pennington, Ben Adlam, Lechao Xiao, Roman Novak, and Jascha Sohl-Dickstein. Finite versus infinite neural networks: an empirical study. *Advances in Neural Information Processing Systems*, 33:15156–15172, 2020.

- [71] Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. In *Advances in Neural Information Processing Systems*, 2018.
- [72] Yuanzhi Li, Tengyu Ma, and Hongyang R Zhang. Learning over-parametrized two-layer neural networks beyond ntk. In *Conference on Learning Theory*, 2020.
- [73] Tao Luo, Zhi-Qin John Xu, Zheng Ma, and Yaoyu Zhang. Phase diagram for two-layer relu neural networks at infinite-width limit. *Journal of Machine Learning Research*, 2021.
- [74] Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural networks. In *International Conference on Learning Representations*, 2019.
- [75] Kaifeng Lyu, Zhiyuan Li, Runzhe Wang, and Sanjeev Arora. Gradient descent on two-layer nets: Margin maximization and simplicity bias. *Advances in Neural Information Processing Systems*, 34:12978–12991, 2021.
- [76] Eran Malach, Pritish Kamath, Emmanuel Abbe, and Nathan Srebro. Quantifying the benefit of using differentiable learning over tangent kernels. *arXiv preprint arXiv:2103.01210*, 2021.
- [77] Christopher D Manning, Kevin Clark, John Hewitt, Urvashi Khandelwal, and Omer Levy. Emergent linguistic structure in artificial neural networks trained by self-supervision. *Proceedings of the National Academy of Sciences*, pages 30046–30054, 2020.
- [78] Alexander G de G Matthews, Mark Rowland, Jiri Hron, Richard E Turner, and Zoubin Ghahramani. Gaussian process behaviour in wide deep neural networks. In *International Conference on Learning Representations*, 2018.
- [79] 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*, 2018.
- [80] Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit. In *Conference on Learning Theory*, pages 2388–2464. PMLR, 2019.
- [81] Behrad Moniri, Donghwan Lee, Hamed Hassani, and Edgar Dobriban. A theory of non-linear feature learning with one gradient step in two-layer neural networks. *arXiv preprint arXiv:2310.07891*, 2023.
- [82] Andrea Montanari and Yiqiao Zhong. The interpolation phase transition in neural networks: Memorization and generalization under lazy training. *The Annals of Statistics*, 2022.
- [83] Edward Moroshko, Blake E Woodworth, Suriya Gunasekar, Jason D Lee, Nati Srebro, and Daniel Soudry. Implicit bias in deep linear classification: Initialization scale vs training accuracy. *Advances in Neural Information Processing Systems*, 33, 2020.
- [84] Alireza Mousavi-Hosseini, Sejun Park, Manuela Girotti, Ioannis Mitliagkas, and Murat A Erdogdu. Neural networks efficiently learn low-dimensional representations with sgd. *arXiv preprint arXiv:2209.14863*, 2022.
- [85] Mor Shpigel Nacson, Suriya Gunasekar, Jason Lee, Nathan Srebro, and Daniel Soudry. Lexicographic and depth-sensitive margins in homogeneous and non-homogeneous deep models. In *International Conference on Machine Learning*, pages 4683–4692. PMLR, 2019.
- [86] Mor Shpigel Nacson, Jason Lee, Suriya Gunasekar, Pedro Henrique Pamplona Savarese, Nathan Srebro, and Daniel Soudry. Convergence of gradient descent on separable data. In *The 22nd International Conference on Artificial Intelligence and Statistics*. PMLR, 2019.
- [87] Mor Shpigel Nacson, Nathan Srebro, and Daniel Soudry. Stochastic gradient descent on separable data: Exact convergence with a fixed learning rate. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 3051–3059. PMLR, 2019.
- [88] Vaishnavh Nagarajan and J Zico Kolter. Uniform convergence may be unable to explain generalization in deep learning. *Advances in Neural Information Processing Systems*, 2019.

- [89] Preetum Nakkiran, Gal Kaplun, Dimitris Kalimeris, Tristan Yang, Benjamin L Edelman, Fred Zhang, and Boaz Barak. Sgd on neural networks learns functions of increasing complexity. *arXiv preprint arXiv:1905.11604*, 2019.
- [90] Behnam Neyshabur. Implicit regularization in deep learning. *arXiv preprint arXiv:1709.01953*, 2017.
- [91] Roman Novak, Lechao Xiao, Jaehoon Lee, Yasaman Bahri, Daniel A Abolafia, Jeffrey Pennington, and Jascha Sohl-Dickstein. Bayesian convolutional neural networks with many channels are gaussian processes. In *International Conference on Learning Representations*, 2019.
- [92] Ryan O’Donnell. *Analysis of boolean functions*. Cambridge University Press, 2014.
- [93] Samet Oymak and Mahdi Soltanolkotabi. Overparameterized nonlinear learning: Gradient descent takes the shortest path? In *International Conference on Machine Learning*, pages 4951–4960. PMLR, 2019.
- [94] Samet Oymak and Mahdi Soltanolkotabi. Toward moderate overparameterization: Global convergence guarantees for training shallow neural networks. *IEEE Journal on Selected Areas in Information Theory*, pages 84–105, 2020.
- [95] Samet Oymak, Zalan Fabian, Mingchen Li, and Mahdi Soltanolkotabi. Generalization guarantees for neural networks via harnessing the low-rank structure of the jacobian. *arXiv preprint arXiv:1906.05392*, 2019.
- [96] Vardan Papyan, XY Han, and David L Donoho. Prevalence of neural collapse during the terminal phase of deep learning training. *Proceedings of the National Academy of Sciences*, pages 24652–24663, 2020.
- [97] Adityanarayanan Radhakrishnan, Daniel Beaglehole, Parthe Pandit, and Mikhail Belkin. Mechanism of feature learning in deep fully connected networks and kernel machines that recursively learn features, 2023.
- [98] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems*, 2008.
- [99] Maria Refinetti, Sebastian Goldt, Florent Krzakala, and Lenka Zdeborov. Classifying high-dimensional gaussian mixtures: Where kernel methods fail and neural networks succeed. In *International Conference on Machine Learning*, pages 8936–8947. PMLR, 2021.
- [100] Yunwei Ren, Mo Zhou, and Rong Ge. Depth separation with multilayer mean-field networks. In *The Eleventh International Conference on Learning Representations*, 2023.
- [101] Itay Safran, Ronen Eldan, and Ohad Shamir. Depth separations in neural networks: what is actually being separated? In *Conference on Learning Theory*, pages 2664–2666. PMLR, 2019.
- [102] Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli. The pitfalls of simplicity bias in neural networks. In *NeurIPS*, 2020.
- [103] Shai Shalev-Shwartz, Ohad Shamir, and Shaked Shammah. Failures of gradient-based deep learning. In *International Conference on Machine Learning*, pages 3067–3075. PMLR, 2017.
- [104] Zhenmei Shi, Yifei Ming, Ying Fan, Frederic Sala, and Yingyu Liang. Domain generalization with nuclear norm regularization. In *NeurIPS 2022 Workshop on Distribution Shifts: Connecting Methods and Applications*, 2022.
- [105] Zhenmei Shi, Junyi Wei, and Yingyu Liang. A theoretical analysis on feature learning in neural networks: Emergence from inputs and advantage over fixed features. In *International Conference on Learning Representations*, 2022.
- [106] Justin Sirignano and Konstantinos Spiliopoulos. Mean field analysis of neural networks: A central limit theorem. *Stochastic Processes and their Applications*, pages 1820–1852, 2020.

- [107] Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*, pages 2822–2878, 2018.
- [108] Dominik Stöger and Mahdi Soltanolkotabi. Small random initialization is akin to spectral learning: Optimization and generalization guarantees for overparameterized low-rank matrix reconstruction. *Advances in Neural Information Processing Systems*, 34:23831–23843, 2021.
- [109] Matus Telgarsky. Feature selection with gradient descent on two-layer networks in low-rotation regimes. *arXiv preprint arXiv:2208.02789*, 2022.
- [110] Rodrigo Veiga, Ludovic Stephan, Bruno Loureiro, Florent Krzakala, and Lenka Zdeborová. Phase diagram of stochastic gradient descent in high-dimensional two-layer neural networks. *arXiv preprint arXiv:2202.00293*, 2022.
- [111] Yifei Wang, Jonathan Lacotte, and Mert Pilanci. The hidden convex optimization landscape of two-layer relu neural networks: an exact characterization of the optimal solutions. *arXiv e-prints*, pages arXiv–2006, 2020.
- [112] 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.
- [113] Blake Woodworth, Suriya Gunasekar, Jason D Lee, Edward Moroshko, Pedro Savarese, Itay Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models. In *Conference on Learning Theory*, 2020.
- [114] Greg Yang. Scaling limits of wide neural networks with weight sharing: Gaussian process behavior, gradient independence, and neural tangent kernel derivation. *arXiv preprint arXiv:1902.04760*, 2019.
- [115] Greg Yang and Edward J Hu. Feature learning in infinite-width neural networks. *arXiv preprint arXiv:2011.14522*, 2020.
- [116] Gilad Yehudai and Ohad Shamir. On the power and limitations of random features for understanding neural networks. *Advances in Neural Information Processing Systems*, 2019.
- [117] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision*, 2014.
- [118] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *International Conference on Learning Representations*, 2017.
- [119] Chiyuan Zhang, Samy Bengio, and Yoram Singer. Are all layers created equal? *arXiv preprint arXiv:1902.01996*, 2019.
- [120] Mo Zhou, Rong Ge, and Chi Jin. A local convergence theory for mildly over-parameterized two-layer neural network. In *COLT*, 2021.
- [121] Difan Zou and Quanquan Gu. An improved analysis of training over-parameterized deep neural networks. *Advances in neural information processing systems*, 32, 2019.
- [122] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Stochastic gradient descent optimizes over-parameterized deep relu networks. *arXiv preprint arXiv:1811.08888*, 2018.
- [123] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Gradient descent optimizes over-parameterized deep ReLU networks. *Machine Learning*, 109:467–492, 2020.
- [124] Difan Zou, Yuan Cao, Yuanzhi Li, and Quanquan Gu. The benefits of mixup for feature learning. In *Proceedings of the 40th International Conference on Machine Learning*, Proceedings of Machine Learning Research, pages 43423–43479. PMLR, 2023.

Appendix

Contents

1	Introduction	1
2	Related Work	2
3	Gradient Feature Learning Framework	3
3.1	Warm Up: A Simple Setting with Frozen First Layer	3
3.2	Core Concepts in the Gradient Feature Learning Framework	4
3.3	Provable Guarantee via Gradient Feature Learning	5
4	Applications in Special Cases	7
4.1	Mixtures of Gaussians	7
4.1.1	Mixtures of Gaussians: Beyond the Kernel Regime	8
4.2	Parity Functions	8
4.2.1	Parity Functions: Beyond the Kernel Regime	10
5	Further Implications and Conclusion	10
A	Broader Impacts	19
B	Limitations	19
C	More Further Implications	20
D	Gradient Feature Learning Framework	21
D.1	Simplified Gradient Feature Learning Framework	21
D.2	Gradient Feature Learning Framework under Expected Risk	22
D.2.1	Feature Learning	23
D.2.2	Good Network Exists	24
D.2.3	Learning an Accurate Classifier	26
D.3	More Discussion about Setting	30
D.4	Gradient Feature Learning Framework under Empirical Risk with Sample Complexity	31
E	Applications in Special Cases	38
E.1	Linear Data	38
E.2	Mixture of Gaussians	40
E.2.1	Problem Setup	40
E.2.2	Mixture of Gaussians: Feature Learning	41
E.2.3	Mixture of Gaussians: Final Guarantee	46
E.3	Mixture of Gaussians - XOR	48
E.3.1	Mixture of Gaussians - XOR: Feature Learning	48

E.3.2	Mixture of Gaussians - XOR: Final Guarantee	52
E.4	Parity Functions	52
E.4.1	Problem Setup	52
E.4.2	Parity Functions: Feature Learning	54
E.4.3	Parity Functions: Final Guarantee	58
E.5	Uniform Parity Functions	59
E.5.1	Uniform Parity Functions: Final Guarantee	60
E.6	Uniform Parity Functions: Alternative Analysis	61
E.6.1	Modified General Feature Learning Framework for Uniform Parity Functions	61
E.6.2	Feature Learning of Uniform Parity Functions	66
E.7	Multiple Index Model with Low Degree Polynomial	68
E.7.1	Problem Setup	68
E.7.2	Multiple Index Model: Final Guarantee	68
F	Auxiliary Lemmas	70

Appendix A discusses the potential societal impact of our work. Appendix B describes the limitations of our work. In Appendix C, we present our framework implications about simplicity bias. The complete proof of our main results is given in Appendix D. We present the case study of linear data in Appendix E.1, mixtures of Gaussians in Appendix E.2 and Appendix E.3, parity functions in Appendix E.4, Appendix E.5 and Appendix E.6, and multiple-index models in Appendix E.7. We put the auxiliary lemmas in Appendix F.

A Broader Impacts

Our paper is purely theoretical in nature, and thus we do not anticipate an immediate negative ethical impact. We provide a unified theoretical framework that can be applied to different theoretical problems. We propose the two key ideas of gradient feature and gradient feature-induced neural networks not only to show their ability to unify several current works but also to open a new direction of thinking with respect to the learning process. These notations have the potential to be extended to multi-layer gradient features and multi-step learning, and this work is only our first step.

On the other hand, this work may lead to a better understanding and inspire the development of improved network learning methods, which may have a positive impact on the theoretical machine-learning community. It may also be beneficial to engineering-inclined machine-learning researchers.

B Limitations

Recover Existing Results. The framework may or may not recover the width or sample complexity bounds in existing work.

1. The framework can give matching bounds as the existing work in some cases, like parities over uniform inputs (Appendix E.5).
2. In some other cases, it gives polynomial error bounds not the same as those in the existing work (e.g., for parities over structured inputs). This is because our work is analyzing general cases, and thus may not give better than or the same bounds as those in special cases, since special cases have more properties that can be exploited to get potentially better bounds. On the other hand, our bounds can already show the advantage over kernel methods (e.g., Proposition 4.9).

We would like to emphasize that our contribution is providing an analysis framework that can (1) formalize the unifying principles of learning features from gradients in network training, and (2) give

polynomial error bounds for prototypical problems. Our focus is not to recover the guarantees in existing work.

Failure Cases. There are some failure cases that gradient feature learning framework cannot cover:

1. In [101], they constructed a function that is easy to approximate using a 3-layer network but not approximable by any 2-layer network. Since the function is not approximable by any 2-layer network, it cannot be approximated by the gradient-induced networks as well, so OPT will be large. As a result, the final error will be large.
2. In uniform parity data distribution, considering an odd number of features rather than even, i.e., k is an odd number in Assumption E.28, we can show that our gradient feature set is empty even when p in Equation (6) is exponentially small, thus the OPT is a positive constant since the gradient induced network can only be constants. Meanwhile, the neural network won't be able to learn this data distribution because its gradient is always 0 through the training, and the final error equals OPT.

The first case corresponds to the approximation hardness of 2-layer networks, while the second case gives a learning hardness example. The above two cases show that if there is an approximation or learning hardness, our gradient feature learning framework may be vacuous because the optimal model in the gradient feature class has a large risk, then the ground-truth mapping from inputs to labels is not learnable by gradient descent. These analyses are consistent with previous works [15, 101].

C More Further Implications

Our general framework also sheds some light on several interesting phenomena in neural network (NN) learning observed in practice. Feature learning beyond the kernel regime has been discussed in Section 4.1.1 and Section 4.2.1. The lottery ticket hypothesis (LTH) has been discussed in Section 5. Below we discuss other implications.

Implicit Regularization/Simplicity Bias. It is now well known that practical NN are overparameterized and traditional uniform convergence bounds cannot adequately explain their generalization performance [59, 88, 118]. It is generally believed that the optimization has some *implicit regularization* effect that restricts learning dynamics to a subset of the whole hypothesis class, which is not of high capacity so can lead to good generalization [53, 90]. Furthermore, learning dynamics tend to first learn simple functions and then learn more and more sophisticated ones (referred to as simplicity bias) [89, 102]. However, it remains elusive to formalize such simplicity bias.

Our framework provides a candidate explanation: the learning dynamics first learn to approximate the best network in a smaller family of gradient feature induced networks $F_{d,r;B_F;S}$ and then learn to approximate the best in a larger family. Consider the number of neurons r for illustration. Let $r_1 < r_2$, and let T_1 and T_2 be their corresponding runtime bounds for T in the main Theorem 3.12. Clearly, $T_1 < T_2$. Then, at time T_1 , the theorem guarantees the learning dynamics learn to approximate the best in the family $F_{d,r_1;B_F;S}$ with r_1 neurons, but not for the larger family $F_{d,r_2;B_F;S}$. Later, at time T_2 , the learning dynamics learn to approximate the best in the larger family $F_{d,r_2;B_F;S}$. That is, the learning first learns simpler functions and then more sophisticated ones where the simplicity bias is measured by the size of the family of gradient feature-induced networks. The implicit regularization is then restricting to networks approximating smaller families of gradient feature-induced networks. Furthermore, we can also conclude that for an SGD-optimized NN, its actual representation power is from the subset of NN based on gradient features, instead of the whole set of NN. This view helps explain the simplicity bias/implicit regularization phenomenon of NN learning in practice.

Learning over Different Data Distributions. Our framework articulates the following key principles (pointed out for specific problems in existing work but not articulated more generally):

- Role of gradient: the gradient leads to the emergence of good features, which is useful for the learning of upper layers in later stages.
- From features to solutions: learned features in early steps will not be distorted, if not improved, in later stages. The training dynamic for upper layers will eventually learn a good combination of hidden neurons based on gradient features, giving a good solution.

Then, more interesting insights are obtained from the generality of the framework. To build a general framework, the meaningful error guarantees should be data-dependent, since NN learning on general data distributions is hard and data-independent guarantees will be vacuous [34, 35]. Comparing the optimal in a family of “ground-truth” functions (inspired by agnostic learning in learning theory) is a useful method to obtain the data-dependent bound. We further construct the “ground-truth” functions using properties of the training dynamics, i.e., gradient features. This greatly facilitates the analysis of the training dynamics and is the key to obtaining the final guarantees. On the other hand, the framework can also be viewed as using the optimal by gradient-induced NN to measure or quantify the “complexity” of the problem. For easier problems, this quantity is smaller, and our framework can give a better error bound. So this provides a united way to derive guarantees for specific problems.

New Perspectives about Roadmaps Forward. We argue a new perspective about the connection between the strong representation power and the successful learning of NN. Traditionally, the strong representation power of NN is the key reason for hardness results of NN learning: NN has strong representation power and can encode hard learning questions, so they are hard to learn. See the proof in SQ bound from [33] or NP-hardness from [19]. The strong representation power also causes trouble for the statistical aspect: it leads to vacuous generalization bounds when traditional uniform convergence tools are used.

Our framework suggests a perspective in sharp contrast: the strong representation power of NN with gradient features is actually the key to successful learning. More concretely, the optimal error of the gradient feature-induced NN being small (i.e., strong representation power for a given data distribution) can lead to a small guarantee, which is the key to successful learning. The above new perspective suggests a different analysis road than traditional ones. Traditional analysis typically first reasons about the optimal based on the whole function class, i.e. the ground truth, then analyze how NN learns proper features and reaches the optimal. In contrast, our framework defines feature family first, and then reasons about the optimal based on it.

Our framework provides the foundation for future work on analyzing gradient-based NN learning, which may inspire future directions including but not limited to (1) defining a new feature family for 2-layer NN rather than gradient feature, (2) considering deep NN and introducing new gradient features (e.g., gradient feature notion for upper layers), (3) defining different gradient feature family at different training stages (e.g., gradient feature notion for later stages). In particular, the challenges in the later-stage analysis are: (a) the weights in the later stage will not be as normal as the initialization, and we need new tools to analyze their properties; (b) to show that the later-stage features eventually lead to a good solution, we may need new analysis tools for the non-convex optimization due to the changes in the first layer weights.

D Gradient Feature Learning Framework

We first prove a Simplified Gradient Feature Learning Framework in Appendix D.1, which only considers one-step gradient feature learning. Then, we prove our Gradient Feature Learning Framework, e.g., no freezing of the first layer. In Appendix D.2, we consider population loss to simplify the proof. Then, we provide more discussion about our problem setup and our core concept in Appendix D.3. Finally, we prove our Gradient Feature Learning Framework under empirical loss considering sample complexity in Appendix D.4.

D.1 Simplified Gradient Feature Learning Framework

Algorithm 2 Training by Algorithm 1 with no updates for the first layer after the first gradient step

Initialize $f_{(\mathbf{a}^{(0)}, \mathbf{W}^{(0)}, \mathbf{b})} \geq F_{d,m}$; Sample $Z \sim D^n$
 Get $(\mathbf{a}^{(1)}, \mathbf{W}^{(1)}, \mathbf{b})$ by one gradient step update and fix $\mathbf{W}^{(1)}, \mathbf{b}$
for $t = 2$ **to** T **do**
 $\mathbf{a}^{(t)} = \mathbf{a}^{(t-1)} - \eta^{(t)} r_{\mathbf{a}} \tilde{\mathcal{L}}_Z(f_{(\mathbf{a}^{(t-1)}, \mathbf{W}^{(1)}, \mathbf{b})})$
end for

Theorem 3.4 (Simple Setting). *Assume $\tilde{\mathcal{L}}_Z(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})})$ is L -smooth to \mathbf{a} . Let $\eta^{(t)} = \frac{1}{L}, \lambda^{(t)} = 0$, for all $t \geq 2, 3, \dots, T$. Training by Algorithm 1 with no updates for the first layer after the*

first gradient step, w.h.p., there exists $t \geq [T]$ such that $L_{\mathcal{D}}(f_{(\mathbf{a}^{(t)}, \mathbf{W}^{(1)}, \mathbf{b})}) \leq \text{OPT}_{\mathbf{W}^{(1)}, \mathbf{b}; B_{a2}} + O\left(\frac{L(\|\mathbf{a}^{(1)}\|_2^2 + B_{a2}^2)}{T} + \sqrt{\frac{B_{a2}^2(\|\mathbf{W}^{(1)}\|_F^2 B_x^2 + \|\mathbf{b}\|_2^2)}{n}}\right)$.

Proof of Theorem 3.4. Recall that

$$\mathcal{F}_{\mathbf{W}, \mathbf{b}; B_{a2}} := \{f_{(\mathbf{a}, \mathbf{W}; \mathbf{b})} \in \mathcal{F}_{d; m} \mid \|\mathbf{a}\|_2 \leq B_{a2}\}, \quad \text{OPT}_{\mathbf{W}, \mathbf{b}; B_{a2}} := \min_{f \in \mathcal{F}_{\mathbf{W}, \mathbf{b}; B_{a2}}} L_{\mathcal{D}}(f). \quad (13)$$

We denote $f^* = \text{argmin}_{f \in \mathcal{F}_{\mathbf{W}, \mathbf{b}; B_{a2}}} L_{\mathcal{D}}(f)$ and $\tilde{f}^* = \text{argmin}_{f \in \mathcal{F}_{\mathbf{W}, \mathbf{b}; B_{a2}}} \tilde{L}_{\mathcal{Z}}(f)$. We use \mathbf{a}^* and \mathbf{a}^* to denote their second layer weights respectively. Then, we have

$$L_{\mathcal{D}}(f_{(\mathbf{a}^{(t)}, \mathbf{W}^{(1)}, \mathbf{b})}) = L_{\mathcal{D}}(f_{(\mathbf{a}^{(t)}, \mathbf{W}^{(1)}, \mathbf{b})}) - \tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}^{(t)}, \mathbf{W}^{(1)}, \mathbf{b})}) \quad (14)$$

$$+ \tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}^{(t)}, \mathbf{W}^{(1)}, \mathbf{b})}) - \tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}^*, \mathbf{W}^{(1)}, \mathbf{b})}) \quad (15)$$

$$+ \tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}^*, \mathbf{W}^{(1)}, \mathbf{b})}) - \tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}^*, \mathbf{W}^{(1)}, \mathbf{b})}) \quad (16)$$

$$+ \tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}^*, \mathbf{W}^{(1)}, \mathbf{b})}) - L_{\mathcal{D}}(f_{(\mathbf{a}^*, \mathbf{W}^{(1)}, \mathbf{b})}) \quad (17)$$

$$+ L_{\mathcal{D}}(f_{(\mathbf{a}^*, \mathbf{W}^{(1)}, \mathbf{b})}) \quad (18)$$

$$\left| L_{\mathcal{D}}(f_{(\mathbf{a}^{(t)}, \mathbf{W}^{(1)}, \mathbf{b})}) - \tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}^{(t)}, \mathbf{W}^{(1)}, \mathbf{b})}) \right| \quad (19)$$

$$+ \left| \tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}^{(t)}, \mathbf{W}^{(1)}, \mathbf{b})}) - \tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}^*, \mathbf{W}^{(1)}, \mathbf{b})}) \right| \quad (20)$$

$$+ 0 \quad (21)$$

$$+ \left| \tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}^*, \mathbf{W}^{(1)}, \mathbf{b})}) - L_{\mathcal{D}}(f_{(\mathbf{a}^*, \mathbf{W}^{(1)}, \mathbf{b})}) \right| \quad (22)$$

$$+ \text{OPT}_{\mathbf{W}^{(1)}, \mathbf{b}; B_{a2}}. \quad (23)$$

Fixing $\mathbf{W}^{(1)}$, \mathbf{b} and optimizing \mathbf{a} only is a convex optimization problem. Note that $\eta \leq \frac{1}{L}$, where $\tilde{L}_{\mathcal{Z}}$ is L -smooth to \mathbf{a} . Thus with gradient descent, we have

$$\frac{1}{T} \sum_{t=1}^T \tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}^{(t)}, \mathbf{W}^{(1)}, \mathbf{b})}) - \tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}^*, \mathbf{W}^{(1)}, \mathbf{b})}) \leq \frac{k_{\mathbf{a}}^{(1)} \|\mathbf{a}^*\|_2^2}{2T\eta}. \quad (24)$$

Then our theorem gets proved by Lemma F.9 and generalization bounds based on Rademacher complexity. \square

D.2 Gradient Feature Learning Framework under Expected Risk

We consider the following training process under population loss to simplify the proof. We prove our Gradient Feature Learning Framework under empirical loss considering sample complexity in Appendix D.4.

Algorithm 3 Network Training via Gradient Descent

Initialize $(\mathbf{a}^{(0)}, \mathbf{W}^{(0)}, \mathbf{b})$ as in Equation (8)
for $t = 1$ **to** T **do**
 $\mathbf{a}^{(t)} = \mathbf{a}^{(t-1)} - \eta^{(t)} \nabla_{\mathbf{a}} L_{\mathcal{D}}^{(t)}(f_{(\cdot, \mathbf{W}^{(t-1)}, \mathbf{b})})$
 $\mathbf{W}^{(t)} = \mathbf{W}^{(t-1)} - \eta^{(t)} \nabla_{\mathbf{W}} L_{\mathcal{D}}^{(t)}(f_{(\cdot, \cdot, \mathbf{b})})$
end for

Given an input distribution, we can get a Gradient Feature set $S_{p; B_G}$ and $f^*(\mathbf{x}) = \sum_{j=1}^r \mathbf{a}_j^* \sigma(\langle \mathbf{w}_j^*, \mathbf{x} \rangle - \mathbf{b}_j^*)$, where $f^* \in \mathcal{F}_{d; r; B_F; S_{p; B_G}}$ is a Gradient Feature Induced networks defined in Definition 3.11. Considering training by Algorithm 3, we have the following results.

Theorem D.1 (Gradient Feature Learning Framework under Expected Risk). *Assume Assumption 3.1. For any $\epsilon, \delta \geq (0, 1)$, if $m \geq e^d$ and*

$$m = \left(\frac{1}{p} \left(\frac{rB_{a1}B_{x1}}{\epsilon} \sqrt{\frac{B_b}{B_G}} \right)^4 + \frac{1}{\delta} + \frac{1}{p} \left(\log \left(\frac{r}{\delta} \right) \right)^2 \right), \quad (25)$$

$$T = \left(\frac{1}{\epsilon} \left(\frac{rB_{a2}B_bB_{x1}}{(mp)^{\frac{1}{4}}} + mb \right) \left(\frac{r \log m}{B_b B_G} + \frac{1}{B_{x1}(mp)^{\frac{1}{4}}} \right) \right), \quad (26)$$

then with proper hyper-parameter values, we have with probability $1 - \delta$, there exists $t \geq [T]$ in Algorithm 3 with

$$\Pr[\text{sign}(f^{(t)}(\mathbf{x})) \neq y] \leq L_{\mathcal{D}}(f^{(t)}) - \text{OPT}_{d;r;B_F;S_{p_i};B_G} + rB_{a1}B_{x1}\sqrt{2\gamma} + \epsilon. \quad (27)$$

See the full statement and proof in Theorem D.9. Below, we show some lemmas used in the analysis of population loss.

D.2.1 Feature Learning

We first show that a large subset of neurons has gradients at the first step as good features.

Definition D.2 (Nice Gradients Set. Equivalent to Equation (9)). *We define*

$$G_{(D_j;+1);Nice} := \left\{ i \geq [m] : \langle \mathbf{w}_i^{(1)}, D_j \rangle > (1 - \gamma) \|\mathbf{w}_i^{(1)}\|_2, \|\mathbf{w}_i^{(1)}\|_2 \geq \left| \eta^{(1)} \ell'(0) \mathbf{a}_i^{(0)} \right|_{B_G} \right\}$$

$$G_{(D_j;-1);Nice} := \left\{ i \geq [2m] \cap [m] : \langle \mathbf{w}_i^{(1)}, D_j \rangle > (1 - \gamma) \|\mathbf{w}_i^{(1)}\|_2, \|\mathbf{w}_i^{(1)}\|_2 \geq \left| \eta^{(1)} \ell'(0) \mathbf{a}_i^{(0)} \right|_{B_G} \right\}$$

where γ, B_G is the same in the Definition 3.7.

Lemma D.3 (Feature Emergence. Full Statement of Lemma 3.13). *Let $\lambda^{(1)} = \frac{1}{(1-\gamma)}$. For any r size subset $f(D_1, s_1), \dots, (D_r, s_r) \in S_{p_i; B_G}$, with probability at least $1 - 2re^{-cmp}$ where $c > 0$ is a universal constant, we have that for all $j \geq [r]$, $|G_{(D_j; s_j); Nice}| \geq \frac{mp}{4}$.*

Proof of Lemma D.3. By symmetric initialization and Lemma F.1, we have for all $i \geq [2m]$

$$\mathbf{w}_i^{(1)} = \eta^{(1)} \ell'(0) \mathbf{a}_i^{(0)} \mathbb{E}_{(\mathbf{x}; y)} \left[y \sigma' \left[\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle - \mathbf{b}_i \right] \mathbf{x} \right] \quad (28)$$

$$= \eta^{(1)} \ell'(0) \mathbf{a}_i^{(0)} G(\mathbf{w}_i^{(0)}, \mathbf{b}_i). \quad (29)$$

For all $j \geq [r]$, as $(D_j, s_j) \in S_{p_i; B_G}$, by Lemma F.3,

(1) if $s_j = +1$, for all $i \geq [m]$, we have

$$\Pr \left[i \geq G_{(D_j; s_j); Nice} \right] \quad (30)$$

$$= \Pr \left[\frac{\langle \mathbf{w}_i^{(1)}, D_j \rangle}{\|\mathbf{w}_i^{(1)}\|_2} > (1 - \gamma), \|\mathbf{w}_i^{(1)}\|_2 \geq \left| \eta^{(1)} \ell'(0) \mathbf{a}_i^{(0)} \right|_{B_G} \right] \quad (31)$$

$$= \Pr \left[\frac{\langle \mathbf{w}_i^{(1)}, D_j \rangle}{\|\mathbf{w}_i^{(1)}\|_2} > (1 - \gamma), \|\mathbf{w}_i^{(1)}\|_2 \geq \left| \eta^{(1)} \ell'(0) \mathbf{a}_i^{(0)} \right|_{B_G}, \frac{\mathbf{b}_i}{\|\mathbf{b}_i\|} = s_j \right] \quad (32)$$

$$\Pr \left[G(\mathbf{w}_i^{(0)}, \mathbf{b}_i) \geq C_{D_j}; \|\mathbf{w}_i^{(0)}\|_2 \geq B_G, \frac{\mathbf{b}_i}{\|\mathbf{b}_i\|} = s_j, \mathbf{a}_i^{(0)} \langle G(\mathbf{w}_i^{(0)}, \mathbf{b}_i), D_j \rangle > 0 \right] \geq \frac{p}{2}, \quad (33)$$

(2) if $s_j = -1$, for all $i \geq [2m] \cap [m]$, similarly we have

$$\Pr \left[i \geq G_{(D_j; s_j); Nice} \right] \geq \frac{p}{2}. \quad (34)$$

By concentration inequality, (Chernoff's inequality under small deviations), we have

$$\Pr \left[|G_{(D_j; s_j); Nice}| < \frac{mp}{4} \right] \leq 2e^{-cmp}. \quad (35)$$

We complete the proof by union bound. \square

D.2.2 Good Network Exists

Then, the gradients allow for obtaining a set of neurons approximating the ‘‘ground-truth’’ network with comparable loss.

Lemma D.4 (Existence of Good Networks. Full Statement of Lemma 3.14). *Let $\lambda^{(1)} = \frac{1}{(1)}$. For any $B \geq (0, B_b)$, let $\sigma_a = \left(\frac{B_{x1} B_b}{(1) B_G B} \right)$ and $\delta = 2re^{-\sqrt{mp}}$. Then, with probability at least $1 - \delta$ over the initialization, there exists \mathbf{a}_i 's such that $f_{(\mathbf{a}, \mathbf{w}^{(1)}, \mathbf{b})}(\mathbf{x}) = \sum_{i=1}^{4m} \mathbf{a}_i \sigma \left(\langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle - \mathbf{b}_i \right)$ satisfies*

$$\mathcal{L}_{\mathcal{D}}(f_{(\mathbf{a}, \mathbf{w}^{(1)}, \mathbf{b})}) \leq r B_{a1} \left(\frac{B_{x1}^2 B_b}{mp B_G B} + B_{x1} \sqrt{2\gamma} + B \right) + \text{OPT}_{d,r;B_F;S_{p_i};B_G}, \quad (36)$$

$$\text{and } \|\mathbf{a}\|_0 = O\left(r(mp)^{\frac{1}{2}}\right), \|\mathbf{a}\|_2 = O\left(\frac{B_{a2} B_b}{b(mp)^{\frac{1}{4}}}\right), \|\mathbf{a}\|_{\infty} = O\left(\frac{B_{a1} B_b}{b(mp)^{\frac{1}{2}}}\right).$$

Proof of Lemma D.4. Recall $f^*(\mathbf{x}) = \sum_{j=1}^r \mathbf{a}_j^* \sigma(\langle \mathbf{w}_j^*, \mathbf{x} \rangle - \mathbf{b}_j^*)$, where $f^* \in \mathcal{F}_{d,r;B_F;S_{p_i};B_G}$ is defined in Definition 3.11 and let $s_j^* = \frac{\mathbf{b}_j^*}{|\mathbf{b}_j^*|}$. By Lemma D.3, with probability at least $1 - \delta_1$, $\delta_1 = 2re^{-cmp}$, for all $j \in [r]$, we have $jG_{(\mathbf{w}_j; s_j); \text{Nice}} \leq \frac{mp}{4}$. Then for all $i \in G_{(\mathbf{w}_j; s_j); \text{Nice}} \subseteq [2m]$, we have $\ell'(0) \eta^{(1)} G(\mathbf{w}_i^{(0)}, \mathbf{b}_i) \frac{\mathbf{b}_j^*}{b}$ only depend on $\mathbf{w}_i^{(0)}$ and \mathbf{b}_i , which is independent of $\mathbf{a}_i^{(0)}$. Given Definition 3.7, we have

$$\ell'(0) \eta^{(1)} kG(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \frac{\mathbf{b}_j^*}{b} \geq \left[\ell'(0) \eta^{(1)} B_{x1} \frac{B_b}{b}, \ell'(0) \eta^{(1)} B_{x1} \frac{B_b}{b} \right]. \quad (37)$$

We split $[r]$ into $\mathcal{I} = \{j \in [r] : \|\mathbf{b}_j^*\| < B\}$, $\mathcal{J} = \{j \in [r] : \|\mathbf{b}_j^*\| \geq B\}$ and $\mathcal{K} = \{j \in [r] : \|\mathbf{b}_j^*\| \geq B\}$. Let $\epsilon_a = \frac{B_{x1} B_b}{\sqrt{mp} B_G B}$. Then we know that for all $j \in \mathcal{I} \cup \mathcal{K}$, for all $i \in G_{(\mathbf{w}_j; s_j); \text{Nice}}$, we have

$$\Pr_{\mathbf{a}_i^{(0)} \sim \mathcal{N}(0; \frac{2}{b})} \left[\left| \mathbf{a}_i^{(0)} \ell'(0) \eta^{(1)} kG(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \frac{\mathbf{b}_j^*}{b} - 1 \right| \leq \epsilon_a \right] \quad (38)$$

$$= \Pr_{\mathbf{a}_i^{(0)} \sim \mathcal{N}(0; \frac{2}{b})} \left[1 - \epsilon_a \leq \mathbf{a}_i^{(0)} \ell'(0) \eta^{(1)} kG(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \frac{\mathbf{b}_j^*}{b} \leq 1 + \epsilon_a \right] \quad (39)$$

$$= \Pr_{g \sim \mathcal{N}(0;1)} \left[1 - \epsilon_a \leq g \left(\frac{kG(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \mathbf{b}_j^*}{B_G B} \right) \leq 1 + \epsilon_a \right] \quad (40)$$

$$= \Pr_{g \sim \mathcal{N}(0;1)} \left[(1 - \epsilon_a) \left(\frac{B_G B}{kG(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \mathbf{b}_j^*} \right) \leq g \leq (1 + \epsilon_a) \left(\frac{B_G B}{kG(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \mathbf{b}_j^*} \right) \right]$$

$$= \left(\frac{\epsilon_a B_G B}{kG(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \mathbf{b}_j^*} \right) \quad (41)$$

$$\left(\frac{\epsilon_a B_G B}{B_{x1} B_b} \right) \quad (42)$$

$$= \left(\frac{1}{mp} \right). \quad (43)$$

Thus, with probability $\left(\frac{1}{\sqrt{mp}} \right)$ over $\mathbf{a}_i^{(0)}$, we have

$$\left| \mathbf{a}_i^{(0)} \ell'(0) \eta^{(1)} kG(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \frac{\mathbf{b}_j^*}{b} - 1 \right| \leq \epsilon_a, \quad \left| \mathbf{a}_i^{(0)} \right| = O\left(\frac{b}{\ell'(0) \eta^{(1)} B_G B} \right). \quad (44)$$

Similarly, for $j \in \mathcal{J}$, for all $i \in G_{(\mathbf{w}_j; s_j); \text{Nice}}$, with probability $\left(\frac{1}{\sqrt{mp}} \right)$ over $\mathbf{a}_i^{(0)}$, we have

$$\left| \mathbf{a}_i^{(0)} \ell'(0) \eta^{(1)} kG(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \frac{B}{b} - 1 \right| \leq \epsilon_a, \quad \left| \mathbf{a}_i^{(0)} \right| = O\left(\frac{b}{\ell'(0) \eta^{(1)} B_G B} \right). \quad (45)$$

For all $j \geq [r]$, let $\mathcal{J}_j = G(\mathbf{w}_j; \mathcal{S}_j)_{\text{Nice}}$ be the set of i 's such that condition Equation (44) or Equation (45) are satisfied. By Chernoff bound and union bound, with probability at least $1 - \delta_2$, $\delta_2 = re^{-\sqrt{mp}}$, for all $j \geq [r]$ we have $|\mathcal{J}_j| \geq \frac{r}{2} \left(\frac{p}{mp} \right)$.

We have for $\delta_j \geq \frac{1}{2} + \frac{1}{2} \epsilon$, $\delta_i \geq \frac{1}{2} - \frac{1}{2} \epsilon$,

$$\left| \frac{j \mathbf{b}_j^*}{b} \langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle - \langle \mathbf{w}_j^*, \mathbf{x} \rangle \right| \quad (46)$$

$$\left\| \mathbf{a}_i^{(0)} \ell'(0) \eta^{(1)} k G(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \frac{j \mathbf{b}_j^*}{b} \frac{\mathbf{w}_i^{(1)}}{k \mathbf{w}_i^{(1)} k_2} - \frac{\mathbf{w}_i^{(1)}}{k \mathbf{w}_i^{(1)} k_2} + \frac{\mathbf{w}_i^{(1)}}{k \mathbf{w}_i^{(1)} k_2} \langle \mathbf{w}_j^*, \mathbf{x} \rangle \right\| k \mathbf{x} k_2 \quad (47)$$

$$(\epsilon_a + \sqrt{2\gamma}) k \mathbf{x} k_2. \quad (48)$$

Similarly, for $\delta_j \geq \frac{1}{2}$, $\delta_i \geq \frac{1}{2} - \frac{1}{2} \epsilon$,

$$\left| \frac{B}{b} \langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle - \langle \mathbf{w}_j^*, \mathbf{x} \rangle \right| \leq (\epsilon_a + \sqrt{2\gamma}) k \mathbf{x} k_2. \quad (49)$$

If $i \geq j$, $j \geq \frac{1}{2} + \frac{1}{2} \epsilon$, set $\mathbf{a}_i = \mathbf{a}_j^* \frac{|\mathbf{b}_j|}{|j|b}$, if $i \geq j$, $j \geq \frac{1}{2}$, set $\mathbf{a}_i = \mathbf{a}_j^* \frac{B}{|j|b}$, otherwise set $\mathbf{a}_i = 0$, we have $k \mathbf{a} k_0 = O\left(r(mp)^{\frac{1}{2}}\right)$, $k \mathbf{a} k_2 = O\left(\frac{B a_2 B_0}{b(mp)^{\frac{1}{4}}}\right)$, $k \mathbf{a} k_\infty = O\left(\frac{B a_1 B_0}{b(mp)^{\frac{1}{2}}}\right)$.

Finally, we have

$$L_{\mathcal{D}}(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}) \quad (50)$$

$$= L_{\mathcal{D}}(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}) - L_{\mathcal{D}}(f^*) + L_{\mathcal{D}}(f^*) \quad (51)$$

$$E_{(\mathbf{x}; \mathbf{y})} \left[\left| f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}(\mathbf{x}) - f^*(\mathbf{x}) \right| \right] + L_{\mathcal{D}}(f^*) \quad (52)$$

$$E_{(\mathbf{x}; \mathbf{y})} \left[\left| \sum_{i=1}^m \mathbf{a}_i \sigma(\langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle + b) + \sum_{i=m+1}^{2m} \mathbf{a}_i \sigma(\langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle + b) - \sum_{j=1}^r \mathbf{a}_j^* \sigma(\langle \mathbf{w}_j^*, \mathbf{x} \rangle + b_j^*) \right| \right] + L_{\mathcal{D}}(f^*) \quad (53)$$

$$E_{(\mathbf{x}; \mathbf{y})} \left[\left| \sum_{j \in \mathcal{J}_+} \sum_{i \in \mathcal{J}_j} \mathbf{a}_j^* \frac{1}{|j|} \left| \frac{j \mathbf{b}_j^*}{b} \sigma(\langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle + b) - \sigma(\langle \mathbf{w}_j^*, \mathbf{x} \rangle + b_j^*) \right| \right| \right] \quad (54)$$

$$+ E_{(\mathbf{x}; \mathbf{y})} \left[\left| \sum_{j \in \mathcal{J}_+} \sum_{i \in \mathcal{J}_j} \mathbf{a}_j^* \frac{1}{|j|} \left| \frac{j \mathbf{b}_j^*}{b} \sigma(\langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle + b) - \sigma(\langle \mathbf{w}_j^*, \mathbf{x} \rangle + b_j^*) \right| \right| \right] \quad (55)$$

$$+ E_{(\mathbf{x}; \mathbf{y})} \left[\left| \sum_{j \in \mathcal{J}_+} \sum_{i \in \mathcal{J}_j} \mathbf{a}_j^* \frac{1}{|j|} \left| \frac{B}{b} \sigma(\langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle + b) - \sigma(\langle \mathbf{w}_j^*, \mathbf{x} \rangle + b_j^*) \right| \right| \right] + L_{\mathcal{D}}(f^*) \quad (56)$$

$$E_{(\mathbf{x}; \mathbf{y})} \left[\left| \sum_{j \in \mathcal{J}_+} \sum_{i \in \mathcal{J}_j} \mathbf{a}_j^* \frac{1}{|j|} \left| \frac{j \mathbf{b}_j^*}{b} \langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle - \langle \mathbf{w}_j^*, \mathbf{x} \rangle \right| \right| \right] \quad (57)$$

$$+ E_{(\mathbf{x}; \mathbf{y})} \left[\left| \sum_{j \in \mathcal{J}_+} \sum_{i \in \mathcal{J}_j} \mathbf{a}_j^* \frac{1}{|j|} \left| \frac{j \mathbf{b}_j^*}{b} \langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle - \langle \mathbf{w}_j^*, \mathbf{x} \rangle \right| \right| \right] \quad (58)$$

$$+ E_{(\mathbf{x}; \mathbf{y})} \left[\left| \sum_{j \in \mathcal{J}_+} \sum_{i \in \mathcal{J}_j} \mathbf{a}_j^* \frac{1}{|j|} \left| \frac{B}{b} \langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle + B - \langle \mathbf{w}_j^*, \mathbf{x} \rangle \right| \right| \right] + L_{\mathcal{D}}(f^*) \quad (59)$$

$$r k \mathbf{a}^* k_\infty (\epsilon_a + \sqrt{2\gamma}) E_{(\mathbf{x}; \mathbf{y})} k \mathbf{x} k_2 + j k \mathbf{a}^* k_\infty B + L_{\mathcal{D}}(f^*) \quad (60)$$

$$r B_{\mathcal{X}1} B_{a1} (\epsilon_a + \sqrt{2\gamma}) + j B_{a1} B + \text{OPT}_{d; r; B_F; \mathcal{S}_P; \mathcal{S}_G}. \quad (61)$$

We finish the proof by union bound and $\delta = \delta_1 + \delta_2$. \square

D.2.3 Learning an Accurate Classifier

We will use the following theorem from existing work to prove that gradient descent learns a good classifier (Theorem D.9). Theorem D.1 is simply a direct corollary of Theorem D.9.

Theorem D.5 (Theorem 13 in [33]). *Fix some η , and let f_1, \dots, f_T be some sequence of convex functions. Fix some θ_1 , and assume we update $\theta_{t+1} = \theta_t - \eta \nabla f_t(\theta_t)$. Then for every θ^* the following holds:*

$$\frac{1}{T} \sum_{t=1}^T f_t(\theta_t) \leq \frac{1}{T} \sum_{t=1}^T f_t(\theta^*) + \frac{1}{2\eta T} k_{\theta^*}^2 + k_{\theta_1} k_2 \frac{1}{T} \sum_{t=1}^T k_T f_t(\theta_t) k_2 + \eta \frac{1}{T} \sum_{t=1}^T k_T f_t(\theta_t) k_2^2.$$

To apply the theorem we first present a few lemmas bounding the change in the network during steps.

Lemma D.6 (Bound of $\mathbf{w}_i^{(0)}$, $\mathbf{w}_i^{(1)}$). *Assume the same conditions as in Lemma D.4, and $d \geq \log m$, with probability at least $1 - \delta - \frac{1}{m^2}$ over the initialization, $k_{\mathbf{a}^{(0)}} k_{\infty} = O\left(\frac{b\sqrt{\log m}}{\sigma_w^{(0)} B_G B}\right)$, and for all $i \geq [4m]$, we have $k_{\mathbf{w}_i^{(0)}} k_2 = O\left(\sigma_w \frac{D}{d}\right)$. Finally, $k_{\mathbf{a}^{(1)}} k_{\infty} = O\left(\eta^{(1)} \ell'(0) (B_{\mathcal{X}1} \sigma_w \frac{D}{d} + \mathfrak{b})\right)$, and for all $i \geq [4m]$, $k_{\mathbf{w}_i^{(1)}} k_2 = O\left(\frac{b\sqrt{\log m} B_{\mathcal{X}1}}{B_G B}\right)$.*

Proof of Lemma D.6. By Lemma F.4, we have $k_{\mathbf{a}^{(0)}} k_{\infty} = O\left(\frac{b\sqrt{\log m}}{\sigma_w^{(0)} B_G B}\right)$ with probability at least $1 - \frac{1}{2m^2}$ by property of maximum i.i.d Gaussians. For any $i \geq [4m]$, by Lemma F.5 and $d \geq \log m$, we have

$$\Pr\left(\frac{1}{\sigma_w^2} \left\| \mathbf{w}_i^{(0)} \right\|_2^2 \geq d + 2\sqrt{4d \log(m)} + 8 \log(m)\right) = O\left(\frac{1}{m^4}\right). \quad (62)$$

Thus, by union bound, with probability at least $1 - \frac{1}{2m^2}$, for all $i \geq [4m]$, we have $k_{\mathbf{w}_i^{(0)}} k_2 = O\left(\sigma_w \frac{D}{d}\right)$.

For all $i \geq [4m]$, we have

$$j \mathbf{a}_i^{(1)} j = \left| \eta^{(1)} \ell'(0) \left[\mathbb{E}_{(\mathbf{x}, y)} \left[y \left[\sigma \left(\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle - \mathfrak{b}_i \right) \right] \right] \right] \right| \quad (63)$$

$$\leq \eta^{(1)} \ell'(0) (k_{\mathbf{w}_i^{(0)}} k_2 \mathbb{E}_{(\mathbf{x}, y)} [k_{\mathbf{x}} k_2] + \mathfrak{b}) \quad (64)$$

$$= O\left(\eta^{(1)} \ell'(0) (B_{\mathcal{X}1} \sigma_w \frac{D}{d} + \mathfrak{b})\right). \quad (65)$$

$$k_{\mathbf{w}_i^{(1)}} k_2 = \eta^{(1)} \ell'(0) \left\| \mathbf{a}_i^{(0)} \mathbb{E}_{(\mathbf{x}, y)} \left[y \sigma' \left[\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle - \mathfrak{b}_i \right] \mathbf{x} \right] \right\|_2 \quad (66)$$

$$= O\left(\frac{b \sqrt{\log m} B_{\mathcal{X}1}}{B_G B}\right). \quad (67)$$

□

Lemma D.7 (Bound of $\mathbf{w}_i^{(t)}$). *Assume the same conditions as in Lemma D.6, and let $\eta = \eta^{(t)}$ for all $t \geq 2, 3, \dots, T$, $0 < T \eta B_{\mathcal{X}1} = o(1)$, and $0 = \lambda = \lambda^{(t)}$ for all $t \geq 2, 3, \dots, T$, for all $i \geq [4m]$, we have*

$$j \mathbf{a}_i^{(t)} j = O\left(j \mathbf{a}_i^{(1)} j + k_{\mathbf{w}_i^{(1)}} k_2 + \frac{\mathfrak{b}}{B_{\mathcal{X}1}} + \eta \mathfrak{b}\right) \quad (68)$$

$$k_{\mathbf{w}_i^{(t)}} k_2 = O\left(t \eta B_{\mathcal{X}1} j \mathbf{a}_i^{(1)} j + t \eta^2 B_{\mathcal{X}1}^2 k_{\mathbf{w}_i^{(1)}} k_2 + t \eta^2 B_{\mathcal{X}1} \mathfrak{b}\right). \quad (69)$$

Proof of Lemma D.7. For all $i \in [4m]$, by Lemma D.6,

$$j\mathbf{a}_i^{(t)}j = \left| (1 - \eta\lambda)\mathbf{a}_i^{(t-1)} - \eta E_{(\mathbf{x};y)} \left[\ell'(y f_{(t-1)}(\mathbf{x})) y \left[\sigma \left(\langle \mathbf{w}_i^{(t-1)}, \mathbf{x} \rangle - \mathbf{b}_i \right) \right] \right] \right| \quad (70)$$

$$\left| (1 - \eta\lambda)\mathbf{a}_i^{(t-1)} \right| + \eta \left| E_{(\mathbf{x};y)} \left[\left[\sigma \left(\langle \mathbf{w}_i^{(t-1)}, \mathbf{x} \rangle - \mathbf{b}_i \right) \right] \right] \right| \quad (71)$$

$$\left| \mathbf{a}_i^{(t-1)} \right| + \eta(B_{x1}k\mathbf{w}_i^{(t-1)}k_2 + \mathfrak{b}) \quad (72)$$

$$\left| \mathbf{a}_i^{(t-1)} \right| + \eta B_{x1}k\mathbf{w}_i^{(t-1)}k_2 + \eta B_{x1}k\mathbf{w}_i^{(1)}k_2 + \eta\mathfrak{b} \quad (73)$$

$$= \left| \mathbf{a}_i^{(t-1)} \right| + \eta B_{x1}k\mathbf{w}_i^{(t-1)}k_2 + \eta Z_i, \quad (74)$$

where we denote $Z_i = B_{x1}k\mathbf{w}_i^{(1)}k_2 + \mathfrak{b}$. Then we give a bound of the first layer's weights change,

$$k\mathbf{w}_i^{(t)}k_2 - \mathbf{w}_i^{(1)}k_2 \quad (75)$$

$$= \left\| (1 - \eta\lambda)\mathbf{w}_i^{(t-1)} - \eta\mathbf{a}_i^{(t-1)} E_{(\mathbf{x};y)} \left[\ell'(y f_{(t-1)}(\mathbf{x})) y \sigma' \left[\langle \mathbf{w}_i^{(t-1)}, \mathbf{x} \rangle - \mathbf{b}_i \right] \mathbf{x} \right] - \mathbf{w}_i^{(1)} \right\|_2 \quad (76)$$

$$k\mathbf{w}_i^{(t-1)}k_2 + \eta B_{x1}j\mathbf{a}_i^{(t-1)}j. \quad (77)$$

Combine two bounds, we can get

$$j\mathbf{a}_i^{(t)}j - j\mathbf{a}_i^{(t-1)}j + \eta Z_i + (\eta B_{x1})^2 \sum_{l=1}^{t-2} j\mathbf{a}_i^{(l)}j \quad (78)$$

$$\leq \sum_{l=1}^t j\mathbf{a}_i^{(l)}j - 2 \left(\sum_{l=1}^{t-1} j\mathbf{a}_i^{(l)}j \right) + (1 - (\eta B_{x1})^2) \left(\sum_{l=1}^{t-2} j\mathbf{a}_i^{(l)}j \right) + \eta Z_i. \quad (79)$$

Let $h(1) = j\mathbf{a}_i^{(1)}j$, $h(2) = 2j\mathbf{a}_i^{(1)}j + \eta Z_i$ and $h(t+2) = 2h(t+1) - (1 - (\eta B_{x1})^2)h(t) + \eta Z_i$ for $n \in \mathbb{N}_+$, by Lemma F.8, we have

$$h(t) = \frac{Z_i}{\eta B_{x1}^2} + c_1(1 - \eta B_{x1})^{(t-1)} + c_2(1 + \eta B_{x1})^{(t-1)} \quad (80)$$

$$c_1 = \frac{1}{2} \left(j\mathbf{a}_i^{(1)}j + \frac{Z_i}{\eta B_{x1}^2} - \frac{j\mathbf{a}_i^{(1)}j + \eta Z_i}{\eta B_{x1}} \right) \quad (81)$$

$$c_2 = \frac{1}{2} \left(j\mathbf{a}_i^{(1)}j + \frac{Z_i}{\eta B_{x1}^2} + \frac{j\mathbf{a}_i^{(1)}j + \eta Z_i}{\eta B_{x1}} \right). \quad (82)$$

Thus, by $j\mathbf{a}_i^{(t)}j \leq c_2$, and $0 < T\eta B_{x1} \ll o(1)$, we have

$$j\mathbf{a}_i^{(t)}j - h(t) \leq h(t-1) \quad (83)$$

$$= \eta B_{x1}c_1(1 - \eta B_{x1})^{(t-2)} + \eta B_{x1}c_2(1 + \eta B_{x1})^{(t-2)} \quad (84)$$

$$2\eta B_{x1}c_2(1 + \eta B_{x1})^t \quad (85)$$

$$O(2\eta B_{x1}c_2). \quad (86)$$

Similarly, by binomial approximation, we also have

$$k\mathbf{w}_i^{(t)}k_2 - \mathbf{w}_i^{(1)}k_2 \leq \eta B_{x1}h(t-1) \quad (87)$$

$$= \eta B_{x1} \left(\frac{Z_i}{\eta B_{x1}^2} + c_1(1 - \eta B_{x1})^{(t-2)} + c_2(1 + \eta B_{x1})^{(t-2)} \right) \quad (88)$$

$$\leq \eta B_{x1} O \left(\frac{Z_i}{\eta B_{x1}^2} + c_1(1 - (t-2)\eta B_{x1}) + c_2(1 + (t-2)\eta B_{x1}) \right) \quad (89)$$

$$\leq \eta B_{x1} O \left(\frac{Z_i}{\eta B_{x1}^2} + c_1 + c_2 + (c_2 - c_1)t\eta B_{x1} \right) \quad (90)$$

$$\leq \eta B_{x1} O \left(j\mathbf{a}_i^{(1)}j + \frac{j\mathbf{a}_i^{(1)}j + \eta Z_i}{\eta B_{x1}} t\eta B_{x1} \right) \quad (91)$$

$$= O \left((\eta j\mathbf{a}_i^{(1)}j + \eta^2 Z_i)t\eta B_{x1} \right). \quad (92)$$

We finish the proof by plugging Z_i, c_2 into the bound. \square

Lemma D.8 (Bound of Loss Gap and Gradient). *Assume the same conditions as in Lemma D.7, for all $t \geq [T]$, we have*

$$|jL_{\mathcal{D}}(f_{(\mathbf{a}; \mathbf{w}^{(t); \mathbf{b}})}) - L_{\mathcal{D}}(f_{(\mathbf{a}; \mathbf{w}^{(1); \mathbf{b}})})| \leq B_{x1} k \mathbf{a} k_2 \sqrt{k \mathbf{a} k_0} \max_{i \in [4m]} k \mathbf{w}_i^{(t)} - \mathbf{w}_i^{(1)} k_2 \quad (93)$$

and for all $t \geq [T]$, for all $i \in [4m]$, we have

$$\left| \frac{\partial L_{\mathcal{D}}(f_{(\cdot; \cdot; \cdot)})}{\partial \mathbf{a}_i^{(t)}} \right| \leq B_{x1} (k \mathbf{w}_i^{(t)} - \mathbf{w}_i^{(1)} k_2 + k \mathbf{w}_i^{(1)} k_2) + \mathfrak{b}. \quad (94)$$

Proof of Lemma D.8. It follows from that

$$|jL_{\mathcal{D}}(f_{(\mathbf{a}; \mathbf{w}^{(t); \mathbf{b}})}) - L_{\mathcal{D}}(f_{(\mathbf{a}; \mathbf{w}^{(1); \mathbf{b}})})| \quad (95)$$

$$= \mathbb{E}_{(\mathbf{x}; y)} [f_{(\mathbf{a}; \mathbf{w}^{(t); \mathbf{b}})}(\mathbf{x}) - f_{(\mathbf{a}; \mathbf{w}^{(1); \mathbf{b}})}(\mathbf{x})] \quad (96)$$

$$= \mathbb{E}_{(\mathbf{x}; y)} \left[k \mathbf{a} k_2 \sqrt{k \mathbf{a} k_0} \max_{i \in [4m]} \left| \sigma \left[\langle \mathbf{w}_i^{(t)}, \mathbf{x} \rangle - \mathbf{b}_i \right] - \sigma \left[\langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle - \mathbf{b}_i \right] \right| \right] \quad (97)$$

$$\leq B_{x1} k \mathbf{a} k_2 \sqrt{k \mathbf{a} k_0} \max_{i \in [4m]} k \mathbf{w}_i^{(t)} - \mathbf{w}_i^{(1)} k_2. \quad (98)$$

Also, we have

$$\left| \frac{\partial L_{\mathcal{D}}(f_{(\cdot; \cdot; \cdot)})}{\partial \mathbf{a}_i^{(t)}} \right| = \left| \mathbb{E}_{(\mathbf{x}; y)} \left[\ell'(y f_{(\cdot; \cdot; \cdot)}(\mathbf{x})) y \left[\sigma \left(\langle \mathbf{w}_i^{(t)}, \mathbf{x} \rangle - \mathbf{b}_i \right) \right] \right] \right| \quad (99)$$

$$\leq B_{x1} k \mathbf{w}_i^{(t)} k_2 + \mathfrak{b} \quad (100)$$

$$\leq B_{x1} (k \mathbf{w}_i^{(t)} - \mathbf{w}_i^{(1)} k_2 + k \mathbf{w}_i^{(1)} k_2) + \mathfrak{b}. \quad (101)$$

\square

We are now ready to prove the main theorem.

Theorem D.9 (Online Convex Optimization. Full Statement of Theorem D.1). *Consider training by Algorithm 3, and any $\delta \in (0, 1)$. Assume $d = \log m$. Set*

$$\sigma_w > 0, \quad \mathfrak{b} > 0, \quad \eta^{(t)} = \eta, \quad \lambda^{(t)} = 0 \text{ for all } t \in \{2, 3, \dots, T\},$$

$$\eta^{(1)} = \left(\frac{\min\{O(\eta), O(\eta \mathfrak{b})\} g}{\ell'(0)(B_{x1} \sigma_w \rho^{\frac{1}{d}} + \mathfrak{b})} \right), \quad \lambda^{(1)} = \frac{1}{\eta^{(1)}}, \quad \sigma_a = \left(\frac{\mathfrak{b}(mp)^{\frac{1}{4}}}{\ell'(0) \eta^{(1)} B_{x1} \rho^{\frac{1}{B_G B_b}}} \right).$$

Let $0 < T \eta B_{x1} = o(1)$, $m = \left(\frac{1}{\sqrt{\epsilon}} + \frac{1}{\rho} (\log \frac{1}{\epsilon})^2 \right)$. With probability at least $1 - \delta$ over the initialization, there exists $t \geq [T]$ such that

$$\begin{aligned} L_{\mathcal{D}}(f_{(\cdot; \cdot; \cdot)}) - \text{OPT}_{d; r; B_F; S_{\rho}; B_G} &+ r B_{a1} \left(\frac{2 B_{x1}}{(mp)^{\frac{1}{4}}} \sqrt{\frac{B_b}{B_G}} + B_{x1} \sqrt{2\gamma} \right) \\ &+ \eta \left(\frac{\rho}{r B_{a2} B_b T \eta B_{x1}^2} + m \mathfrak{b} \right) O \left(\frac{\rho \log m B_{x1} (mp)^{\frac{1}{4}}}{B_b B_G} + 1 \right) + O \left(\frac{B_{a2}^2 B_b^2}{\eta T \mathfrak{b}^2 (mp)^{\frac{1}{2}}} \right). \end{aligned} \quad (102)$$

Furthermore, for any $\epsilon \in (0, 1)$, set

$$\mathfrak{b} = \left(\frac{B_G^{\frac{1}{4}} B_{a2} B_b^{\frac{3}{4}}}{r B_{a1}} \right), \quad m = \left(\frac{1}{\rho \epsilon^4} \left(r B_{a1} B_{x1} \sqrt{\frac{B_b}{B_G}} \right)^4 + \frac{1}{\rho \delta} + \frac{1}{\rho} \left(\log \left(\frac{r}{\delta} \right) \right)^2 \right), \quad (103)$$

$$\eta = \left(\frac{\epsilon}{\left(\frac{\sqrt{r} B_{a2} B_b B_{x1}}{(mp)^{\frac{1}{4}}} + m \mathfrak{b} \right) \left(\frac{\sqrt{\log m B_{x1} (mp)^{\frac{1}{4}}}}{\sqrt{B_b B_G}} + 1 \right)} \right), \quad T = \left(\frac{1}{\eta B_{x1} (mp)^{\frac{1}{4}}} \right), \quad (104)$$

we have there exists $t \geq [T]$ with

$$\Pr[\text{sign}(f_{(\cdot; \cdot; \cdot)}(\mathbf{x})) \neq y] \leq L_{\mathcal{D}}(f_{(\cdot; \cdot; \cdot)}) - \text{OPT}_{d; r; B_F; S_{\rho}; B_G} + r B_{a1} B_{x1} \sqrt{2\gamma} + \epsilon. \quad (105)$$

Proof of Theorem D.9. By $m = \left(\frac{1}{\sqrt{\rho}} + \frac{1}{\rho} (\log(\frac{\rho}{\delta}))^2\right)$ we have $2re^{-\sqrt{m\rho}} + \frac{1}{m^2} \leq \delta$. For any $B \geq (0, B_b)$, when $\sigma_a = \left(\frac{\rho}{-\ell'(0)(B_{x1} \sqrt{B_G B_b}})\right)$, by Theorem D.5, Lemma D.4, Lemma D.8, with probability at least $1 - \delta$ over the initialization, we have

$$\frac{1}{T} \sum_{t=1}^T L_{\mathcal{D}}(f^{(t)}) \quad (106)$$

$$\frac{1}{T} \sum_{t=1}^T j(L_{\mathcal{D}}(f_{(\mathbf{a}; \mathbf{W}^{(t); \mathbf{b})})} - L_{\mathcal{D}}(f_{(\mathbf{a}; \mathbf{W}^{(1); \mathbf{b})})) + L_{\mathcal{D}}(f_{(\mathbf{a}; \mathbf{W}^{(1); \mathbf{b})})) \quad (107)$$

$$+ \frac{\mathbf{k}\mathbf{a}k_2^2}{2\eta T} + (2\mathbf{k}\mathbf{a}^{(1)}k_2 \frac{\rho}{m} + 4\eta m) \max_{i \in [4m]} \left| \frac{\partial L_{\mathcal{D}}(f^{(T)})}{\partial \mathbf{a}_i^{(T)}} \right| \quad (108)$$

$$\text{OPT}_{d; r; B_F; S_P; B_G} + rB_{a1} \left(\frac{B_{x1}^2 B_b}{\rho m \rho B_G B} + B_{x1} \sqrt{2\gamma} + B \right) \quad (109)$$

$$+ B_{x1} \mathbf{k}\mathbf{a}k_2 \sqrt{\mathbf{k}\mathbf{a}k_0} \max_{i \in [4m]} \mathbf{k}\mathbf{w}_i^{(T)} \mathbf{w}_i^{(1)} k_2 \quad (110)$$

$$+ \frac{\mathbf{k}\mathbf{a}k_2^2}{2\eta T} + 4mB_{x1} (\mathbf{k}\mathbf{a}^{(1)}k_{\infty} + \eta) \left(\max_{i \in [4m]} \mathbf{k}\mathbf{w}_i^{(T)} \mathbf{w}_i^{(1)} k_2 + \max_{i \in [4m]} \mathbf{k}\mathbf{w}_i^{(1)} k_2 + \frac{\mathbf{b}}{B_{x1}} \right). \quad (111)$$

By Lemma D.4, Lemma D.6, Lemma D.7, when $\eta^{(1)} = \left(\frac{\min\{O(\cdot); O(\frac{\mathbf{b}}{\rho})\}}{-\ell'(0)(B_{x1} \sqrt{d+\mathbf{b}})}\right)$, we have

$$\mathbf{k}\mathbf{a}k_0 = O\left(r(mp)^{\frac{1}{2}}\right), \quad \mathbf{k}\mathbf{a}k_2 = O\left(\frac{B_{a2}B_b}{\mathbf{b}(mp)^{\frac{1}{4}}}\right) \quad (112)$$

$$\mathbf{k}\mathbf{a}^{(1)}k_{\infty} = O\left(\eta^{(1)} \ell'(0)(B_{x1} \sigma_w \frac{\rho}{d} + \mathbf{b})\right) \quad (113)$$

$$= \min\{O(\eta), O(\eta\mathbf{b})\} \quad (114)$$

$$\max_{i \in [4m]} \mathbf{k}\mathbf{w}_i^{(1)} k_2 = O\left(\frac{\mathbf{b} \frac{\rho}{\log m} B_{x1}}{B_G B}\right) \quad (115)$$

$$\max_{i \in [4m]} \mathbf{k}\mathbf{w}_i^{(T)} \mathbf{w}_i^{(1)} k_2 = O\left(T\eta B_{x1} \mathbf{k}\mathbf{a}^{(1)}k_{\infty} + T\eta^2 B_{x1}^2 \max_{i \in [4m]} \mathbf{k}\mathbf{w}_i^{(1)} k_2 + T\eta^2 B_{x1} \mathbf{b}\right) \quad (116)$$

$$= O\left(T\eta^2 B_{x1}^2 \left(\max_{i \in [4m]} \mathbf{k}\mathbf{w}_i^{(1)} k_2 + \frac{\mathbf{b}}{B_{x1}}\right)\right). \quad (117)$$

Set $B = \frac{B_{x1}}{(mp)^{\frac{1}{4}}} \sqrt{\frac{B_b}{B_G}}$, we have $\sigma_a = \left(\frac{\mathbf{b}(mp)^{\frac{1}{4}}}{-\ell'(0)(B_{x1} \sqrt{B_G B_b})}\right)$ which satisfy the requirements.

Then,

$$\frac{1}{T} \sum_{t=1}^T L_{\mathcal{D}}(f^{(t)}) \quad (118)$$

$$\text{OPT}_{d; r; B_F; S_P; B_G} + rB_{a1} \left(\frac{2B_{x1}}{(mp)^{\frac{1}{4}}} \sqrt{\frac{B_b}{B_G}} + B_{x1} \sqrt{2\gamma} \right) \quad (119)$$

$$+ \left(\frac{\rho}{r} B_{a2} B_b T \eta^2 B_{x1}^2 \frac{B_{x1}}{\mathbf{b}} + m\eta B_{x1} \right) O\left(\frac{\mathbf{b} \frac{\rho}{\log m} B_{x1}}{B_G B} + \frac{\mathbf{b}}{B_{x1}}\right) + O\left(\frac{B_{a2}^2 B_b^2}{\eta T \mathbf{b}^2 (mp)^{\frac{1}{2}}}\right)$$

$$\text{OPT}_{d; r; B_F; S_P; B_G} + rB_{a1} \left(\frac{2B_{x1}}{(mp)^{\frac{1}{4}}} \sqrt{\frac{B_b}{B_G}} + B_{x1} \sqrt{2\gamma} \right) \quad (120)$$

$$+ \eta \left(\frac{\rho}{r} B_{a2} B_b T \eta B_{x1}^2 + m\mathbf{b} \right) O\left(\frac{\rho \log m B_{x1} (mp)^{\frac{1}{4}}}{B_b B_G} + 1\right) + O\left(\frac{B_{a2}^2 B_b^2}{\eta T \mathbf{b}^2 (mp)^{\frac{1}{2}}}\right). \quad (121)$$

Furthermore, for any $\epsilon \in (0, 1)$, set

$$\tilde{b} = \left(\frac{B_G^{\frac{1}{4}} B_{a2} B_b^{\frac{3}{4}}}{r B_{a1}} \right), \quad m = \left(\frac{1}{p \epsilon^4} \left(r B_{a1} B_{x1} \sqrt{\frac{B_b}{B_G}} \right)^4 + \frac{1}{\delta} + \frac{1}{p} \left(\log \left(\frac{r}{\delta} \right) \right)^2 \right), \quad (122)$$

$$\eta = \left(\frac{\epsilon}{\left(\frac{\sqrt{r} B_{a2} B_b B_{x1}}{(mp)^{\frac{1}{4}}} + m \tilde{b} \right) \left(\frac{\sqrt{\log m} B_{x1} (mp)^{\frac{1}{4}}}{\sqrt{B_b B_G}} + 1 \right)} \right), \quad T = \left(\frac{1}{\eta B_{x1} (mp)^{\frac{1}{4}}} \right), \quad (123)$$

we have

$$\frac{1}{T} \sum_{t=1}^T L_{\mathcal{D}}(f^{(t)}) \leq \text{OPT}_{d;r;B_F;S_{p:};B_G} + r B_{a1} \left(\frac{2 B_{x1}}{(mp)^{\frac{1}{4}}} \sqrt{\frac{B_b}{B_G}} + B_{x1} \sqrt{2\gamma} \right) + \frac{\epsilon}{2} \quad (124)$$

$$+ O \left(\frac{B_{x1} B_{a2}^2 B_b^2}{\tilde{b}^2 (mp)^{\frac{1}{4}}} \right) \quad (125)$$

$$\text{OPT}_{d;r;B_F;S_{p:};B_G} + r B_{a1} B_{x1} \sqrt{2\gamma} + \epsilon. \quad (126)$$

We finish the proof as the 0-1 classification error is bounded by the loss function, e.g., $|\text{sign}(f(\mathbf{x})) \neq y| \leq \frac{\ell(f(\mathbf{x}))}{\ell(0)}$, where $\ell(0) = 1$. \square

D.3 More Discussion about Setting

Range of σ_w . In practice, the value of σ_w cannot be arbitrary, because its choice will have an effect on the Gradient Feature set $S_{p:};B_G$. On the other hand, $d = \log m$ is a natural assumption, otherwise, the two-layer neural networks may fall in the NTK regime.

Parameter Choice. We use $\lambda = 1/\eta$ in the first step so that the neural network will totally forget its initialization, leading to the feature emergence here. This is a common setting for analysis convenience in previous work, e.g., [32, 33, 105]. We can extend this to other choices (e.g., small initialization and large step size for the first few steps), as long as after the gradient update, the gradient dominates the neuron weights. We use $\lambda = 0$ afterward as the regularization effect is weak in our analysis. We can extend our analysis to λ being a small value.

Early Stopping. Our analysis divides network learning into two stages: the feature learning stage, and then classifier learning over the good features. The feature learning stage is simplified to one gradient step for the convenience of analysis, while in practice feature learning can happen in multiple steps. The current framework focuses on the gradient features in the early gradient steps, while feature learning can also happen in later steps, in particular for more complicated data. It is an interesting direction to extend the analysis to a longer training horizon.

Role of s . The s encodes the sign of the bias term, which is important. Recall that we do not update the bias term for simplicity. Let's consider a simple toy example. Assume we have $f_1(x) = a_1 \sigma(w_1^\top x + 1)$, $f_2(x) = a_2 \sigma(w_2^\top x - 1)$ and $f_3(x) = a_3 \sigma(w_3^\top x + 2)$, where σ is ReLU activation function which is a homogeneous function.

1. The sign of the bias term is important. We can see that we always have $a_1 \sigma(w_1^\top x + 1) \neq a_2 \sigma(w_2^\top x - 1)$ for any a_1, w_1, a_2, w_2 . This means that $f_1(x)$ and $f_2(x)$ are intrinsically different and have different active patterns. Thus, we need to handle the sign of the bias term carefully.
2. The scaling of the bias is absorbed. On the other hand, we can see that $a_1 \sigma(w_1^\top x + 1) = a_3 \sigma(w_3^\top x + 2)$ when $a_1 = 2a_3, 2w_1 = w_3$. It means that the scale of the bias term is less important, which can be absorbed into other terms.

Thus, we only need to handle bias with different signs carefully.

Gradient Feature Distribution. We may define a gradient feature distribution rather than a gradient feature set. However, we find that the technical tools used in this continuous setting are pretty different from the discrete version.

Activation Functions. We can change the ReLU activation function to a sublinear activation function, e.g. leaky ReLU, sigmoid, to get a similar conclusion. First, we need to introduce a corresponding gradient feature set, and then we can make it by following the same analysis pipeline. For simplicity, we present ReLU only.

D.4 Gradient Feature Learning Framework under Empirical Risk with Sample Complexity

In this section, we consider training with empirical risk. Intuitively, the proof is straightforward from the proof for population loss. We can simply replace the population loss with the empirical loss, which will introduce an error term in the gradient analysis. We use concentration inequality to control the error term and show that the error term depends inverse-polynomially on the sample size n .

Definition D.10 (Empirical Simplified Gradient Vector). Recall $Z = f(\mathbf{x}^{(l)}, y^{(l)})g_{l \in [n]}$, for any $\mathbf{w} \in \mathbb{R}^d$, $b \in \mathbb{R}$, an Empirical Simplified Gradient Vector is defined as

$$\tilde{G}(\mathbf{w}, b) := \frac{1}{n} \sum_{l \in [n]} [y^{(l)} \mathbf{x}^{(l)}] \mathbb{1}[\mathbf{w}^\top \mathbf{x}^{(l)} > b]. \quad (127)$$

Definition D.11 (Empirical Gradient Feature). Recall $Z = f(\mathbf{x}^{(l)}, y^{(l)})g_{l \in [n]}$, let $\mathbf{w} \in \mathbb{R}^d$, $b \in \mathbb{R}$ be random variables drawn from some distribution \mathcal{W}, \mathcal{B} . An Empirical Gradient Feature set with parameters p, γ, B_G is defined as:

$$\tilde{S}_{p, \gamma, B_G}(\mathcal{W}, \mathcal{B}) := \left\{ (D, s) \mid \Pr_{\mathbf{w}, b} \left[\tilde{G}(\mathbf{w}, b) \in \mathcal{C}_D \text{ and } k\tilde{G}(\mathbf{w}, b)k_2 \leq B_G \text{ and } s = \frac{b}{\gamma} \right] \geq p \right\}.$$

When clear from context, write it as $\tilde{S}_{p, \gamma, B_G}$.

Considering training by Algorithm 1, we have the following results.

Theorem 3.12 (Main Result). Assume Assumption 3.1. For any $\epsilon, \delta \in (0, 1)$, if $m \geq e^d$ and

$$\begin{aligned} m &= \left(\frac{1}{p\epsilon^4} \left(rB_{a1}B_{x1} \sqrt{\frac{B_b}{B_G}} \right)^4 + \frac{1}{\delta} + \frac{1}{p} \left(\log \left(\frac{r}{\delta} \right) \right)^2 \right), \\ T &= \left(\frac{1}{\epsilon} \left(\frac{\rho r B_{a2} B_b B_{x1}}{(mp)^{\frac{1}{4}}} + m\mathfrak{b} \right) \left(\frac{\rho \log m}{B_b B_G} + \frac{1}{B_{x1} (mp)^{\frac{1}{4}}} \right) \right), \\ \frac{n}{\log n} &\geq \left(\frac{m^3 p B_x^2 B_{a2}^4 B_b}{\epsilon^2 r^2 B_{a1}^2 B_G} + \frac{(mp)^{\frac{1}{2}} B_{x2}}{B_b B_G} + \frac{B_x^2}{B_{x2}} + \frac{1}{p} + \left(\frac{1}{B_G^2} + \frac{1}{B_{x1}^2} \right) \frac{B_{x2}}{j\ell'(0)j^2} + \frac{Tm}{\delta} \right), \end{aligned}$$

then with initialization (8) and proper hyper-parameter values, we have with probability $1 - \delta$ over the initialization and training samples, there exists $t \geq T$ in Algorithm 1 with:

$$\begin{aligned} \Pr[\text{sign}(f^{(t)}(\mathbf{x})) \neq y] &\leq L_{\mathcal{D}}(f^{(t)}) \\ &\leq \text{OPT}_{d, r, B_F, S_{p, \gamma, B_G}} + rB_{a1}B_{x1} \sqrt{2\gamma} + O\left(\frac{\rho B_{x2} \log n}{B_G j\ell'(0)jn^{\frac{1}{2}}} \right) + \epsilon. \end{aligned}$$

See the full statement and proof in Theorem D.17. Below, we show some lemmas used in the analysis under empirical loss.

Lemma D.12 (Empirical Gradient Concentration Bound). When $\frac{n}{\log n} > \frac{B_x^2}{B_{x2}}$, with probability at least $1 - O\left(\frac{1}{n}\right)$ over training samples, for all $i \in [4m]$, we have

$$\left\| \frac{\partial \tilde{L}_{\mathcal{Z}}(f)}{\partial \mathbf{w}_i} - \frac{\partial L_{\mathcal{D}}(f)}{\partial \mathbf{w}_i} \right\|_2 \leq O\left(\frac{j\mathbf{a}_i j^{\rho} B_{x2} \log n}{n^{\frac{1}{2}}} \right), \quad (128)$$

$$\left| \frac{\partial \tilde{L}_{\mathcal{Z}}(f)}{\partial \mathbf{a}_i} - \frac{\partial L_{\mathcal{D}}(f)}{\partial \mathbf{a}_i} \right| \leq O\left(\frac{k\mathbf{w}_i k_2^{\rho} B_{x2} \log n}{n^{\frac{1}{2}}} \right), \quad (129)$$

$$\left| \tilde{L}_{\mathcal{Z}}(f) - L_{\mathcal{D}}(f) \right| \leq O\left(\frac{\left(k\mathbf{a}_0 k_0 k\mathbf{a}_\infty (\max_{i \in [4m]} k\mathbf{w}_i k_2 B_x + \mathfrak{b}) + 1 \right)^{\rho} \log n}{n^{\frac{1}{2}}} \right). \quad (130)$$

Proof of Lemma D.12. First, we define,

$$\mathbf{z}^{(l)} = \ell'(y^{(l)} f(\mathbf{x}^{(l)})) y^{(l)} \left[\sigma' \left(\langle \mathbf{w}_i, \mathbf{x}^{(l)} \rangle - \mathbf{b}_i \right) \mathbf{x}^{(l)} \right] \quad (131)$$

$$\mathbb{E}_{(\mathbf{x};y)} [\ell'(y f(\mathbf{x})) y [\sigma'(\langle \mathbf{w}_i, \mathbf{x} \rangle - \mathbf{b}_i)] \mathbf{x}]. \quad (132)$$

As $j\ell'(z)j \leq 1, jyj \leq 1, j\sigma'(z)j \leq 1$, we have $\mathbf{z}^{(l)}$ is zero-mean random vector with $\|\mathbf{z}^{(l)}\|_2 \leq 2B_x$ as well as $\mathbb{E} [\|\mathbf{z}^{(l)}\|_2^2] \leq B_{x2}$. Then by Vector Bernstein Inequality, Lemma 18 in [66], for $0 < z < \frac{B_{x2}}{B_x}$ we have

$$\Pr \left(\left\| \frac{\partial \tilde{L}_{\mathcal{Z}}(f)}{\partial \mathbf{w}_i} - \frac{\partial L_{\mathcal{D}}(f)}{\partial \mathbf{w}_i} \right\|_2 \geq j\mathbf{a}_i j z \right) = \Pr \left(\left\| \frac{1}{n} \sum_{l \in [n]} \mathbf{z}^{(l)} \right\|_2 \geq z \right) \quad (133)$$

$$\exp \left(-n \frac{z^2}{8B_{x2}} + \frac{1}{4} \right). \quad (134)$$

Thus, let $z = n^{-\frac{1}{2}} \frac{B_{x2}}{B_x} \log n$, with probability at least $1 - O\left(\frac{1}{n}\right)$, we have

$$\left\| \frac{\partial \tilde{L}_{\mathcal{Z}}(f)}{\partial \mathbf{w}_i} - \frac{\partial L_{\mathcal{D}}(f)}{\partial \mathbf{w}_i} \right\|_2 \leq O \left(\frac{j\mathbf{a}_i j \frac{B_{x2}}{B_x} \log n}{n^{\frac{1}{2}}} \right). \quad (135)$$

On the other hand, by Bernstein Inequality, for $z > 0$ we have

$$\Pr \left(\left| \frac{\partial \tilde{L}_{\mathcal{Z}}(f)}{\partial \mathbf{a}_i} - \frac{\partial L_{\mathcal{D}}(f)}{\partial \mathbf{a}_i} \right| > z k \mathbf{w}_i k_2 \right) \quad (136)$$

$$= \Pr \left(\left| \frac{1}{n} \sum_{l \in [n]} \left(\ell'(y^{(l)} f(\mathbf{x}^{(l)})) y^{(l)} \left[\sigma' \left(\langle \mathbf{w}_i, \mathbf{x}^{(l)} \rangle - \mathbf{b}_i \right) \right] \right) \right| > z k \mathbf{w}_i k_2 \right) \quad (137)$$

$$\mathbb{E}_{(\mathbf{x};y)} [\ell'(y f(\mathbf{x})) y [\sigma'(\langle \mathbf{w}_i, \mathbf{x} \rangle - \mathbf{b}_i)]] \right| > z k \mathbf{w}_i k_2 \right) \quad (138)$$

$$2 \exp \left(-\frac{\frac{1}{2} n z^2}{B_{x2} + \frac{1}{3} B_x z} \right). \quad (139)$$

Thus, when $\frac{n}{\log n} > \frac{B_x^2}{B_{x2}}$, let $z = n^{-\frac{1}{2}} \frac{B_{x2}}{B_x} \log n$, with probability at least $1 - O\left(\frac{1}{n}\right)$, we have

$$\left| \frac{\partial \tilde{L}_{\mathcal{Z}}(f)}{\partial \mathbf{a}_i} - \frac{\partial L_{\mathcal{D}}(f)}{\partial \mathbf{a}_i} \right| \leq O \left(\frac{k \mathbf{w}_i k_2 \frac{B_{x2}}{B_x} \log n}{n^{\frac{1}{2}}} \right). \quad (140)$$

Finally, we have

$$\left| \tilde{L}_{\mathcal{Z}}(f) - L_{\mathcal{D}}(f) \right| \quad (141)$$

$$= \left| \frac{1}{n} \sum_{l=1}^n \left(\ell \left(y^{(l)} \mathbf{a}^\top \left[\sigma(\mathbf{W}^\top \mathbf{x}^{(l)} - \mathbf{b}) \right] \right) - \mathbb{E}_{(\mathbf{x};y) \sim \mathcal{D}} \left[\ell \left(y \mathbf{a}^\top \left[\sigma(\mathbf{W}^\top \mathbf{x} - \mathbf{b}) \right] \right) \right] \right) \right|. \quad (142)$$

By Assumption 3.1, we have $\ell \left(y^{(l)} \mathbf{a}^\top \left[\sigma(\mathbf{W}^\top \mathbf{x}^{(l)} - \mathbf{b}) \right] \right) - \mathbb{E}_{(\mathbf{x};y) \sim \mathcal{D}} \left[\ell \left(y \mathbf{a}^\top \left[\sigma(\mathbf{W}^\top \mathbf{x} - \mathbf{b}) \right] \right) \right]$ is a zero-mean random variable, with bound $2k\mathbf{a}k_0 k\mathbf{a}k_\infty (\max_{i \in [4m]} k\mathbf{w}_i k_2 B_x + \mathfrak{b}) + 2$. By Hoeffding's inequality, for all $z > 0$, we have

$$\Pr \left(\left| \tilde{L}_{\mathcal{Z}}(f) - L_{\mathcal{D}}(f) \right| \geq z \right) \leq 2 \exp \left(-\frac{z^2 n}{(k\mathbf{a}k_0 k\mathbf{a}k_\infty (\max_{i \in [4m]} k\mathbf{w}_i k_2 B_x + \mathfrak{b}) + 1)^2} \right).$$

Thus, with probability at least $1 - O\left(\frac{1}{n}\right)$, we have

$$\left| \tilde{L}_{\mathcal{Z}}(f) - L_{\mathcal{D}}(f) \right| \leq O \left(\frac{(k\mathbf{a}k_0 k\mathbf{a}k_\infty (\max_{i \in [4m]} k\mathbf{w}_i k_2 B_x + \mathfrak{b}) + 1)^{\frac{D}{2}} \log n}{n^{\frac{1}{2}}} \right). \quad (143)$$

□

The gradients allow for obtaining a set of neurons approximating the ‘‘ground-truth’’ network with comparable loss.

Lemma D.13 (Existence of Good Networks under Empirical Risk). *Suppose $\frac{n}{\log n} > \left(\frac{B_x^2}{B_{x2}} + \frac{1}{\rho} + \frac{B_{x2}}{B_G^2 |\Gamma^0(0)|^2}\right)$. Let $\lambda^{(1)} = \frac{1}{(1)}$. For any $B \geq (0, B_b)$, let $\sigma_a = \left(\frac{b}{|\Gamma^0(0)|^{(1)} B_G B}\right)$ and $\delta = 2re^{-\frac{\rho mp}{2}}$. Then, with probability at least $1 - \delta$ over the initialization and training samples, there exists \mathbf{a}_i 's such that $f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}(\mathbf{x}) = \sum_{i=1}^{4m} \mathbf{a}_i \sigma \left(\langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle - \mathbf{b}_i\right)$ satisfies*

$$L_{\mathcal{D}}(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}) \quad (144)$$

$$rB_{a1} \left(\frac{2B_{x1}^2 B_b}{\rho mp B_G B} + B_{x1} \sqrt{2\gamma + O\left(\frac{\rho B_{x2} \log n}{B_G j \ell'(0) j n^{\frac{1}{2}}}\right)} + B \right) + \text{OPT}_{d; r; B_F; S_{p^-}; B_G}, \quad (145)$$

$$\text{and } k\mathbf{a}k_0 = O\left(r(mp)^{\frac{1}{2}}\right), k\mathbf{a}k_2 = O\left(\frac{B_{a2} B_b}{b(mp)^{\frac{1}{4}}}\right), k\mathbf{a}k_{\infty} = O\left(\frac{B_{a1} B_b}{b(mp)^{\frac{1}{2}}}\right).$$

Proof of Lemma D.13. Denote $\rho = O\left(\frac{1}{n}\right)$ and $\beta = O\left(\frac{\sqrt{B_{x2} \log n}}{n^{\frac{1}{2}}}\right)$. Note that by symmetric initialization, we have $\ell'(y f_{(0)}(\mathbf{x})) = j\ell'(0)j$ for any $\mathbf{x} \geq X$, so that, by Lemma D.12, we have $\left\| \tilde{G}(\mathbf{w}_i^{(0)}, \mathbf{b}_i) - G(\mathbf{w}_i^{(0)}, \mathbf{b}_i) \right\|_2 \leq \frac{\rho}{|\Gamma^0(0)|}$ with probability at least $1 - \rho$. Thus, by union bound, we can see that $S_{p^-}; B_G \subseteq \tilde{S}_{p^-}; \frac{B_G}{|\Gamma^0(0)|}; B_G - \frac{\rho}{|\Gamma^0(0)|}$. Consequently, we have $\text{OPT}_{d; r; B_F; \tilde{S}_{p^-}; \frac{B_G}{|\Gamma^0(0)|}; B_G - \frac{\rho}{|\Gamma^0(0)|}} \leq \text{OPT}_{d; r; B_F; S_{p^-}; B_G}$. Exactly follow the proof in

Lemma D.4 by replacing $S_{p^-}; B_G$ to $\tilde{S}_{p^-}; \frac{B_G}{|\Gamma^0(0)|}; B_G - \frac{\rho}{|\Gamma^0(0)|}$. Then, we finish the proof by $\rho \leq \frac{\rho}{2}, \frac{\rho}{|\Gamma^0(0)|} \leq (1 - \frac{\rho}{2}) B_G$. \square

We will use Theorem D.5 to prove that gradient descent learns a good classifier (Theorem D.17). Theorem 3.12 is simply a direct corollary of Theorem D.17. To apply the theorem we first present a few lemmas bounding the change in the network during steps.

Lemma D.14 (Bound of $(0), (1)$ under Empirical Risk). *Assume the same conditions as in Lemma D.13, and $d \geq \log m$, with probability at least $1 - \delta = \frac{1}{m^2} = O\left(\frac{m}{n}\right)$ over the initialization and training samples, $k\mathbf{a}^{(0)}k_{\infty} = O\left(\frac{b\sqrt{\log m}}{|\Gamma^0(0)|^{(1)} B_G B}\right)$, and for all $i \geq [4m]$, we have $k\mathbf{w}_i^{(0)}k_2 = O\left(\sigma_w \frac{\rho}{d}\right)$. Finally, $k\mathbf{a}^{(1)}k_{\infty} = O\left(\eta^{(1)} j\ell'(0)j (B_{x1} \sigma_w \frac{\rho}{d} + b) + \eta^{(1)} \frac{w\sqrt{dB_{x2} \log n}}{n^{\frac{1}{2}}}\right)$, and for all $i \geq [4m]$, $k\mathbf{w}_i^{(1)}k_2 = O\left(\frac{b\sqrt{\log m} B_{x1}}{B_G B} + \frac{b\sqrt{\log m} B_{x2} \log n}{|\Gamma^0(0)| B_G B n^{\frac{1}{2}}}\right)$.*

Proof of Lemma D.14. The proof exactly follows the proof of Lemma D.6 with Lemma D.12. \square

Lemma D.15 (Bound of (t) under Empirical Risk). *Assume the same conditions as in Lemma D.14, and let $\eta = \eta^{(t)}$ for all $t \geq f2, 3, \dots, Tg$, $0 < T\eta B_{x1} = o(1)$, and $0 = \lambda = \lambda^{(t)}$ for all $t \geq f2, 3, \dots, Tg$. With probability at least $1 - O\left(\frac{Tm}{n}\right)$ over training samples, for all $i \geq [4m]$, for all $t \geq f2, 3, \dots, Tg$, we have*

$$j\mathbf{a}_i^{(t)}j = O\left(j\mathbf{a}_i^{(1)}j + k\mathbf{w}_i^{(1)}k_2 + \frac{b}{\left(B_{x1} + \frac{\sqrt{B_{x2} \log n}}{n^{\frac{1}{2}}}\right)} + \eta b\right) \quad (146)$$

$$k\mathbf{w}_i^{(t)}k_2 = k\mathbf{w}_i^{(1)}k_2 + O\left(t\eta \left(B_{x1} + \frac{\rho B_{x2} \log n}{n^{\frac{1}{2}}}\right) j\mathbf{a}_i^{(1)}j + t\eta^2 \left(B_{x1} + \frac{\rho B_{x2} \log n}{n^{\frac{1}{2}}}\right)^2 k\mathbf{w}_i^{(1)}k_2 + t\eta^2 \left(B_{x1} + \frac{\rho B_{x2} \log n}{n^{\frac{1}{2}}}\right) b\right). \quad (147)$$

Proof of Lemma D.15. The proof exactly follows the proof of Lemma D.7 with Lemma D.12. Note that, we have

$$j\mathbf{a}_i^{(t)j} \left| \mathbf{a}_i^{(t-1)} \right| + \eta(B_{x1}k\mathbf{w}_i^{(t-1)}k_2 + \mathfrak{b}) + \eta \frac{k\mathbf{w}_i^{(t-1)}k_2 \rho_{B_{x2} \log n}}{n^{\frac{1}{2}}} \quad (148)$$

$$\left| \mathbf{a}_i^{(t-1)} \right| + \eta \left(B_{x1} + \frac{\rho_{B_{x2} \log n}}{n^{\frac{1}{2}}} \right) k\mathbf{w}_i^{(t-1)}k_2 + \eta Z_i, \quad (149)$$

where we denote $Z_i = \left(B_{x1} + \frac{\sqrt{B_{x2} \log n}}{n^{\frac{1}{2}}} \right) k\mathbf{w}_i^{(1)}k_2 + \mathfrak{b}$. Similarly, we have

$$k\mathbf{w}_i^{(t)}k_2 - \mathbf{w}_i^{(1)}k_2 - k\mathbf{w}_i^{(t-1)}k_2 - \mathbf{w}_i^{(1)}k_2 + \eta \left(B_{x1} + \frac{\rho_{B_{x2} \log n}}{n^{\frac{1}{2}}} \right) j\mathbf{a}_i^{(t-1)j}. \quad (150)$$

We finish the proof by following the same arguments in the proof of Lemma D.7 and union bound. \square

Lemma D.16 (Bound of Loss Gap and Gradient under Empirical Risk). *Assume the same conditions as in Lemma D.15. With probability at least $1 - O\left(\frac{T}{n}\right)$, for all $t \geq [T]$, we have*

$$\left| \tilde{L}_{\mathcal{Z}^{(t)}}(f_{(\mathbf{a}, \mathbf{W}^{(t)}, \mathfrak{b})}) - L_{\mathcal{D}}(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathfrak{b})}) \right| \quad (151)$$

$$O \left(\frac{\left(k\mathbf{a}k_0 k\mathbf{a}k_{\infty} (\max_{i \in [4m]} k\mathbf{w}_i^{(t)}k_2 B_x + \mathfrak{b}) + 1 \right) \rho_{\log n}}{n^{\frac{1}{2}}} \right) \quad (152)$$

$$+ B_{x1} k\mathbf{a}k_2 \sqrt{k\mathbf{a}k_0} \max_{i \in [4m]} k\mathbf{w}_i^{(t)}k_2 - \mathbf{w}_i^{(1)}k_2. \quad (153)$$

With probability at least $1 - O\left(\frac{T}{n}\right)$, for all $t \geq [T]$, $i \geq [4m]$ we have

$$\left| \frac{\partial \tilde{L}_{\mathcal{Z}^{(t)}}(f_{(t)})}{\partial \mathbf{a}_i^{(t)}} \right| \leq B_{x1} (k\mathbf{w}_i^{(t)}k_2 - \mathbf{w}_i^{(1)}k_2 + k\mathbf{w}_i^{(1)}k_2) + \mathfrak{b} + O \left(\frac{k\mathbf{w}_i^{(t)}k_2 \rho_{B_{x2} \log n}}{n^{\frac{1}{2}}} \right). \quad (154)$$

Proof of Lemma D.16. By Lemma D.8 and Lemma D.12, with probability at least $1 - O\left(\frac{T}{n}\right)$, for all $t \geq [T]$, we have

$$\left| \tilde{L}_{\mathcal{Z}^{(t)}}(f_{(\mathbf{a}, \mathbf{W}^{(t)}, \mathfrak{b})}) - L_{\mathcal{D}}(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathfrak{b})}) \right| \quad (155)$$

$$\left| \tilde{L}_{\mathcal{Z}^{(t)}}(f_{(\mathbf{a}, \mathbf{W}^{(t)}, \mathfrak{b})}) - L_{\mathcal{D}}(f_{(\mathbf{a}, \mathbf{W}^{(t)}, \mathfrak{b})}) \right| + \left| L_{\mathcal{D}}(f_{(\mathbf{a}, \mathbf{W}^{(t)}, \mathfrak{b})}) - L_{\mathcal{D}}(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathfrak{b})}) \right| \quad (156)$$

$$O \left(\frac{\left(k\mathbf{a}k_0 k\mathbf{a}k_{\infty} (\max_{i \in [4m]} k\mathbf{w}_i^{(t)}k_2 B_x + \mathfrak{b}) + 1 \right) \rho_{\log n}}{n^{\frac{1}{2}}} \right) \quad (157)$$

$$+ B_{x1} k\mathbf{a}k_2 \sqrt{k\mathbf{a}k_0} \max_{i \in [4m]} k\mathbf{w}_i^{(t)}k_2 - \mathbf{w}_i^{(1)}k_2. \quad (158)$$

By Lemma D.8 and Lemma D.12, with probability at least $1 - O\left(\frac{T}{n}\right)$, for all $t \geq [T]$, $i \geq [4m]$ we have

$$\left| \frac{\partial \tilde{L}_{\mathcal{Z}^{(t)}}(f_{(t)})}{\partial \mathbf{a}_i^{(t)}} \right| \leq B_{x1} (k\mathbf{w}_i^{(t)}k_2 - \mathbf{w}_i^{(1)}k_2 + k\mathbf{w}_i^{(1)}k_2) + \mathfrak{b} + O \left(\frac{k\mathbf{w}_i^{(t)}k_2 \rho_{B_{x2} \log n}}{n^{\frac{1}{2}}} \right). \quad (159)$$

\square

We are now ready to prove the main theorem.

Theorem D.17 (Online Convex Optimization under Empirical Risk. Full Statement of Theorem 3.12). *Consider training by Algorithm I, and any $\delta \geq (0, 1)$. Assume $d \leq \log m$. Set*

$$\sigma_w > 0, \quad \mathfrak{b} > 0, \quad \eta^{(t)} = \eta, \quad \lambda^{(t)} = 0 \text{ for all } t \geq 2, 3, \dots, T, \\ \eta^{(1)} = \left(\frac{\min\{fO(\eta), O(\eta\mathfrak{b})g\}}{\ell'(0)(B_{x1}\sigma_w - \mathfrak{d} + \mathfrak{b})} \right)^{\frac{1}{2}}, \quad \lambda^{(1)} = \frac{1}{\eta^{(1)}}, \quad \sigma_a = \left(\frac{\mathfrak{b}(mp)^{\frac{1}{4}}}{\ell'(0)\eta^{(1)}B_{x1} \rho_{B_G B_b}} \right).$$

Let $0 < T\eta B_{X1} = o(1)$, $m = \left(\frac{1}{\sqrt{\rho}} + \frac{1}{p} (\log(\rho))^2\right)$ and $\frac{n}{\log n} > \left(\frac{B_x^2}{B_{X2}} + \frac{1}{p} + \left(\frac{1}{B_G^2} + \frac{1}{B_{X1}^2}\right) \frac{B_{X2}}{|\ell'(0)|^2} + \frac{Tm}{\rho}\right)$. With probability at least $1 - \delta$ over the initialization and training samples, there exists $t \geq 2[T]$ such that

$$L_{\mathcal{D}}(f^{(t)}) \tag{160}$$

$$\text{OPT}_{d;r;B_F;S_{p^*};B_G} + rB_{a1} \left(\frac{2^{\rho} B_{X1}}{(mp)^{\frac{1}{4}}} \sqrt{\frac{B_b}{B_G}} + B_{X1} \sqrt{2\gamma} + O\left(\frac{\rho B_{X2} \log n}{B_G j\ell'(0) j n^{\frac{1}{2}}}\right) \right) \tag{161}$$

$$+ \eta \left(\frac{\rho}{r B_{a2} B_b T \eta B_{X1}^2} + m \bar{b} \right) O\left(\frac{\rho \log m B_{X1} (mp)^{\frac{1}{4}}}{\rho B_b B_G} + 1\right) + O\left(\frac{B_{a2}^2 B_b^2}{\eta T \bar{b}^2 (mp)^{\frac{1}{2}}}\right) \tag{162}$$

$$+ \frac{\rho \log n}{n^{\frac{1}{2}}} O\left(\left(\frac{r B_{a1} B_b}{\bar{b}} + m \left(\frac{\rho \log m (mp)^{\frac{1}{4}}}{\rho B_b B_G} + \frac{\bar{b}}{B_{X1}}\right)\right)\right) \tag{163}$$

$$\left(\left(\frac{\rho \log m (mp)^{\frac{1}{4}}}{\rho B_b B_G} + T \eta^2 B_{X1} \bar{b}\right) B_x + \bar{b}\right) + 2 \tag{164}$$

$$+ \frac{\rho \log n}{n^{\frac{1}{2}}} O\left(m \eta \left(\frac{\rho \log m (mp)^{\frac{1}{4}}}{\rho B_b B_G} + T \eta^2 B_{X1} \bar{b}\right) \sqrt{B_{X2}}\right). \tag{165}$$

Furthermore, for any $\epsilon \in (0, 1)$, set

$$\bar{b} = \left(\frac{B_G^{\frac{1}{4}} B_{a2} B_b^{\frac{3}{4}}}{r B_{a1}}\right), \quad m = \left(\frac{1}{p \epsilon^4} \left(r B_{a1} B_{X1} \sqrt{\frac{B_b}{B_G}}\right)^4 + \frac{1}{\rho} + \frac{1}{p} (\log(\frac{r}{\delta}))^2\right),$$

$$\eta = \left(\frac{\epsilon}{\left(\frac{\sqrt{r} B_{a2} B_b B_{X1}}{(mp)^{\frac{1}{4}}} + m \bar{b}\right) \left(\frac{\sqrt{\log m} B_{X1} (mp)^{\frac{1}{4}}}{\sqrt{B_b B_G}} + 1\right)}\right), \quad T = \left(\frac{1}{\eta B_{X1} (mp)^{\frac{1}{4}}}\right),$$

$$\frac{n}{\log n} = \left(\frac{m^3 p B_x^2 B_{a2}^4 B_b (\log m)^2}{\epsilon^2 r^2 B_{a1}^2 B_G} + \frac{(mp)^{\frac{1}{2}} B_{X2} \log m}{B_b B_G} + \frac{B_x^2}{B_{X2}} + \frac{1}{p} + \left(\frac{1}{B_G^2} + \frac{1}{B_{X1}^2}\right) \frac{B_{X2}}{j\ell'(0)^2} + \frac{Tm}{\delta}\right),$$

we have there exists $t \geq 2[T]$ with

$$\Pr[\text{sign}(f^{(t)}(\mathbf{x})) \neq y] \leq L_{\mathcal{D}}(f^{(t)}) \tag{166}$$

$$\text{OPT}_{d;r;B_F;S_{p^*};B_G} + rB_{a1} B_{X1} \sqrt{2\gamma} + O\left(\frac{\rho B_{X2} \log n}{B_G j\ell'(0) j n^{\frac{1}{2}}}\right) + \epsilon. \tag{167}$$

Proof of Theorem D.17. We follow the proof in Theorem D.9. By $m = \left(\frac{1}{\sqrt{\rho}} + \frac{1}{p} (\log(\rho))^2\right)$ and $\frac{n}{\log n} > \left(\frac{B_x^2}{B_{X2}} + \frac{1}{p} + \frac{B_{X2}}{B_G^2 |\ell'(0)|^2} + \frac{Tm}{\rho}\right)$, we have $2re^{-\frac{\rho}{m^2}} + \frac{1}{m^2} + O\left(\frac{Tm}{n}\right) \leq \delta$. For any $B \geq (0, B_b)$, when $\sigma_a = \left(\frac{b}{|\ell'(0)| B_G B}\right)$, by Theorem D.5, Lemma D.12, Lemma D.13,

Lemma D.16, with probability at least $1 - \delta$ over the initialization and training samples, we have

$$\frac{1}{T} \sum_{t=1}^T L_{\mathcal{D}}(f^{(t)}) \quad (168)$$

$$\frac{1}{T} \sum_{t=1}^T j L_{\mathcal{D}}(f^{(t)}) - \tilde{L}_{\mathcal{Z}^{(t)}}(f^{(t)}) j + \frac{1}{T} \sum_{t=1}^T \tilde{L}_{\mathcal{Z}^{(t)}}(f^{(t)}) \quad (169)$$

$$\frac{1}{T} \sum_{t=1}^T j L_{\mathcal{D}}(f^{(t)}) - \tilde{L}_{\mathcal{Z}^{(t)}}(f^{(t)}) j + \frac{1}{T} \sum_{t=1}^T \left| \tilde{L}_{\mathcal{Z}^{(t)}}(f_{(\mathbf{a}, \mathbf{W}^{(t)}, \mathbf{b})}) - L_{\mathcal{D}}(f_{(\mathbf{a}, \mathbf{W}^{(t)}, \mathbf{b})}) \right| \quad (170)$$

$$+ L_{\mathcal{D}}(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}) + \frac{\mathbf{k}\mathbf{a}k_2^2}{2\eta T} + (2\mathbf{k}\mathbf{a}^{(1)}k_2 \rho_{\bar{m}} + 4\eta m) \max_{t \in [T], i \in [4m]} \left| \frac{\partial \tilde{L}_{\mathcal{Z}^{(t)}}(f^{(t)})}{\partial \mathbf{a}_i^{(t)}} \right| \quad (171)$$

$$B_{x1} \mathbf{k}\mathbf{a}k_2 \sqrt{\mathbf{k}\mathbf{a}k_0} \max_{i \in [4m]} \mathbf{k}\mathbf{w}_i^{(T)} \quad \mathbf{w}_i^{(1)} k_2 \quad (172)$$

$$+ O \left(\frac{(\mathbf{k}\mathbf{a}k_0 \mathbf{k}\mathbf{a}k_{\infty} + m \mathbf{k}\mathbf{a}^{(T)} k_{\infty}) (\max_{i \in [4m]} \mathbf{k}\mathbf{w}_i^{(T)} k_2 B_x + \mathfrak{b}) + 2}{n^{\frac{1}{2}}} \rho_{\log n} \right) \quad (173)$$

$$+ \text{OPT}_{d;r;B_F;S_P;B_G} + r B_{a1} \left(\frac{2B_{x1}^2 B_b}{\rho_{\bar{m}} B_G B} + B_{x1} \sqrt{2\gamma + O \left(\frac{\rho_{\bar{m}}}{B_{x2} \log n} \right)} + B \right) \quad (174)$$

$$+ \frac{\mathbf{k}\mathbf{a}k_2^2}{2\eta T} + 4m B_{x1} (\mathbf{k}\mathbf{a}^{(1)} k_{\infty} + \eta) \quad (175)$$

$$\left(\max_{i \in [4m]} \mathbf{k}\mathbf{w}_i^{(T)} \quad \mathbf{w}_i^{(1)} k_2 + \max_{i \in [4m]} \mathbf{k}\mathbf{w}_i^{(1)} k_2 + \frac{\mathfrak{b}}{B_{x1}} + O \left(\frac{\max_{i \in [4m]} \mathbf{k}\mathbf{w}_i^{(T)} k_2 \rho_{\bar{m}} \log n}{B_{x1} n^{\frac{1}{2}}} \right) \right).$$

Set $B = \frac{B_{x1}}{(mp)^{\frac{1}{4}}} \sqrt{\frac{2B_b}{B_G}}$, we have $\sigma_a = \left(\frac{b(mp)^{\frac{1}{4}}}{-\rho(0) \binom{1}{1} B_{x1} \sqrt{B_G B_b}} \right)$ which satisfy the requirements. By Lemma D.13, Lemma D.14, Lemma D.15, $\frac{n}{\log n} > \left(\frac{B_{x2}}{B_{x1}^{|\rho(0)|^2}} \right)$, when $\eta^{(1)} = \left(\frac{\min\{O(\cdot), O(\mathfrak{b})\}}{-\rho(0) (B_{x1} \quad w \sqrt{d+\mathfrak{b}})} \right)$, we have

$$\mathbf{k}\mathbf{a}k_0 = O \left(r (mp)^{\frac{1}{2}} \right), \quad \mathbf{k}\mathbf{a}k_2 = O \left(\frac{B_{a2} B_b}{\mathfrak{b} (mp)^{\frac{1}{4}}} \right), \quad \mathbf{k}\mathbf{a}k_{\infty} = O \left(\frac{B_{a1} B_b}{\mathfrak{b} (mp)^{\frac{1}{2}}} \right) \quad (176)$$

$$\mathbf{k}\mathbf{a}^{(1)} k_{\infty} = O \left(\eta^{(1)} j \ell'(0) j (B_{x1} \sigma_w \rho_{\bar{d}} + \mathfrak{b}) + \eta^{(1)} \frac{\sigma_w \rho_{\bar{d}} B_{x2} \log n}{n^{\frac{1}{2}}} \right) \quad (177)$$

$$= \min f O(\eta), O(\eta \mathfrak{b}) g \quad (178)$$

$$\mathbf{k}\mathbf{a}^{(T)} k_{\infty} = O \left(\mathbf{k}\mathbf{a}^{(1)} k_{\infty} + \max_{i \in [4m]} \mathbf{k}\mathbf{w}_i^{(1)} k_2 + \frac{\mathfrak{b}}{\left(B_{x1} + \frac{\sqrt{B_{x2} \log n}}{n^{\frac{1}{2}}} \right)} + \eta \mathfrak{b} \right) \quad (179)$$

$$O \left(\max_{i \in [4m]} \mathbf{k}\mathbf{w}_i^{(1)} k_2 + \frac{\mathfrak{b}}{B_{x1}} \right) \quad (180)$$

$$\max_{i \in [4m]} k\mathbf{w}_i^{(1)} k_2 = O\left(\frac{\mathfrak{b}^{\rho} \overline{\log m} B_{x1}}{B_G B} + \frac{\mathfrak{b}^{\rho} \overline{\log m}^{\rho} B_{x2} \log n}{j\ell'(0)jB_G B n^{\frac{1}{2}}}\right) \quad (181)$$

$$= O\left(\frac{\mathfrak{b}^{\rho} \overline{\log m} B_{x1}}{B_G B}\right) \quad (182)$$

$$= O\left(\frac{\mathfrak{b}^{\rho} \overline{\log m} (mp)^{\frac{1}{4}}}{B_b B_G}\right) \quad (183)$$

$$\max_{i \in [4m]} k\mathbf{w}_i^{(T)} \quad \mathbf{w}_i^{(1)} k_2 = O\left(T\eta \left(B_{x1} + \frac{\rho B_{x2} \log n}{n^{\frac{1}{2}}}\right) j\mathbf{a}_i^{(1)} j\right) \quad (184)$$

$$+ T\eta^2 \left(B_{x1} + \frac{\rho B_{x2} \log n}{n^{\frac{1}{2}}}\right)^2 k\mathbf{w}_i^{(1)} k_2 \quad (185)$$

$$+ T\eta^2 \left(B_{x1} + \frac{\rho B_{x2} \log n}{n^{\frac{1}{2}}}\right) \mathfrak{b} \quad (186)$$

$$= O\left(T\eta^2 B_{x1}^2 \left(\max_{i \in [4m]} k\mathbf{w}_i^{(1)} k_2 + \frac{\mathfrak{b}}{B_{x1}}\right)\right). \quad (187)$$

Then, following the proof in Theorem D.9, we have

$$\frac{1}{T} \sum_{t=1}^T L_D(f^{(t)}) \quad (188)$$

$$\text{OPT}_{d,r;B_F;S_{p^*};B_G} + rB_{a1} \left(\frac{2^{\rho-2} B_{x1}}{(mp)^{\frac{1}{4}}} \sqrt{\frac{B_b}{B_G}} + B_{x1} \sqrt{2\gamma + O\left(\frac{\rho B_{x2} \log n}{B_G j\ell'(0)j n^{\frac{1}{2}}}\right)} \right) \quad (189)$$

$$+ \eta \left(\rho \bar{r} B_{a2} B_b T \eta B_{x1}^2 + m\mathfrak{b} \right) O\left(\frac{\rho \overline{\log m} B_{x1} (mp)^{\frac{1}{4}}}{B_b B_G} + 1\right) + O\left(\frac{B_{a2}^2 B_b^2}{\eta T \mathfrak{b}^2 (mp)^{\frac{1}{2}}}\right) \quad (190)$$

$$+ O\left(\frac{\left((k\mathbf{a}_0 k_{\infty} + m k\mathbf{a}^{(T)} k_{\infty}) (\max_{i \in [4m]} k\mathbf{w}_i^{(T)} k_2 B_x + \mathfrak{b}) + 2\right)^{\rho} \overline{\log n}}{n^{\frac{1}{2}}}\right) \quad (191)$$

$$+ O\left(\frac{m\eta \max_{i \in [4m]} k\mathbf{w}_i^{(T)} k_2 \rho \overline{B_{x2} \log n}}{n^{\frac{1}{2}}}\right) \quad (192)$$

$$\text{OPT}_{d,r;B_F;S_{p^*};B_G} + rB_{a1} \left(\frac{2^{\rho-2} B_{x1}}{(mp)^{\frac{1}{4}}} \sqrt{\frac{B_b}{B_G}} + B_{x1} \sqrt{2\gamma + O\left(\frac{\rho B_{x2} \log n}{B_G j\ell'(0)j n^{\frac{1}{2}}}\right)} \right) \quad (193)$$

$$+ \eta \left(\rho \bar{r} B_{a2} B_b T \eta B_{x1}^2 + m\mathfrak{b} \right) O\left(\frac{\rho \overline{\log m} B_{x1} (mp)^{\frac{1}{4}}}{B_b B_G} + 1\right) + O\left(\frac{B_{a2}^2 B_b^2}{\eta T \mathfrak{b}^2 (mp)^{\frac{1}{2}}}\right) \quad (194)$$

$$+ \frac{\rho \overline{\log n}}{n^{\frac{1}{2}}} O\left(\left(\frac{rB_{a1} B_b}{\mathfrak{b}} + m \left(\frac{\mathfrak{b}^{\rho} \overline{\log m} (mp)^{\frac{1}{4}}}{B_b B_G} + \frac{\mathfrak{b}}{B_{x1}}\right)\right)\right) \quad (195)$$

$$\left(\left(\frac{\mathfrak{b}^{\rho} \overline{\log m} (mp)^{\frac{1}{4}}}{B_b B_G} + T\eta^2 B_{x1} \mathfrak{b}\right) B_x + \mathfrak{b}\right) + 2 \quad (196)$$

$$+ \frac{\rho \overline{\log n}}{n^{\frac{1}{2}}} O\left(m\eta \left(\frac{\mathfrak{b}^{\rho} \overline{\log m} (mp)^{\frac{1}{4}}}{B_b B_G} + T\eta^2 B_{x1} \mathfrak{b}\right) \sqrt{B_{x2}}\right). \quad (197)$$

Furthermore, for any $\epsilon \in (0, 1)$, set

$$\begin{aligned} \mathfrak{b} &= \left(\frac{B_G^{\frac{1}{2}} B_{a2} B_b^{\frac{3}{4}}}{r B_{a1}} \right), \quad m = \left(\frac{1}{p \epsilon^4} \left(r B_{a1} B_{x1} \sqrt{\frac{B_b}{B_G}} \right)^4 + \frac{1}{\delta} + \frac{1}{p} \left(\log \left(\frac{r}{\delta} \right) \right)^2 \right), \\ \eta &= \left(\frac{\epsilon}{\left(\frac{\sqrt{r} B_{a2} B_b B_{x1}}{(mp)^{\frac{1}{4}}} + m \mathfrak{b} \right) \left(\frac{\sqrt{\log m} B_{x1} (mp)^{\frac{1}{4}}}{\sqrt{B_b B_G}} + 1 \right)} \right), \quad T = \left(\frac{1}{\eta B_{x1} (mp)^{\frac{1}{4}}} \right), \\ \frac{n}{\log n} &= \left(\frac{m^3 p B_x^2 B_{a2}^4 B_b (\log m)^2}{\epsilon^2 r^2 B_{a1}^2 B_G} + \frac{(mp)^{\frac{1}{2}} B_{x2} \log m}{B_b B_G} + \frac{B_x^2}{B_{x2}} + \frac{1}{p} + \left(\frac{1}{B_G^2} + \frac{1}{B_{x1}^2} \right) \frac{B_{x2}}{j \ell'(0) j^2} + \frac{T m}{\delta} \right), \end{aligned}$$

and note that B_G , B_{x1} , B_x and $\sqrt{B_{x2}}$, B_x naturally, we have

$$\frac{1}{T} \sum_{t=1}^T L_{\mathcal{D}}(f^{(t)}) \tag{198}$$

$$\text{OPT}_{d;r;B_F;S_P;B_G} + r B_{a1} \left(\frac{2^{\frac{\rho}{2}} B_{x1}}{(mp)^{\frac{1}{4}}} \sqrt{\frac{B_b}{B_G}} + B_{x1} \sqrt{2\gamma + O\left(\frac{\rho B_{x2} \log n}{B_G j \ell'(0) j n^{\frac{1}{2}}}\right)} \right) \tag{199}$$

$$+ \frac{\epsilon}{2} + O\left(\frac{B_{x1} B_{a2}^2 B_b^2}{\mathfrak{b}^2 (mp)^{\frac{1}{4}}}\right) + \frac{\rho \log n}{n^{\frac{1}{2}}} O\left(\frac{m B_x B_{a2}^{\frac{\rho}{2}} B_b (mp)^{\frac{1}{2}} \log m}{r B_{a1} B_G}\right) \tag{200}$$

$$+ \frac{\rho \log n}{n^{\frac{1}{2}}} O\left(\frac{\epsilon^{\frac{\rho}{2}} B_{x2} \log m (mp)^{\frac{1}{4}}}{B_b B_G}\right) \tag{201}$$

$$\text{OPT}_{d;r;B_F;S_P;B_G} + r B_{a1} B_{x1} \sqrt{2\gamma + O\left(\frac{\rho B_{x2} \log n}{B_G j \ell'(0) j n^{\frac{1}{2}}}\right)} + \epsilon. \tag{202}$$

We finish the proof as the 0-1 classification error is bounded by the loss function, e.g., $|\text{sign}(f(\mathbf{x})) \neq y| \leq \frac{\ell(f(\mathbf{x}))}{\ell(0)}$, where $\ell(0) = 1$. □

E Applications in Special Cases

We present the case study of linear data in Appendix E.1, mixtures of Gaussians in Appendix E.2 and Appendix E.3, parity functions in Appendix E.4, Appendix E.5 and Appendix E.6, and multiple-index models in Appendix E.7.

In special case applications, we consider binary classification with hinge loss, e.g., $\ell(z) = \max\{z, 0\}$. Let $X = \mathbb{R}^d$ be the input space, and $Y = \{-1, 1\}$ be the label space.

Remark E.1 (Hinge Loss and Logistic Loss). *Both hinge loss and logistic loss can be used in special cases and general cases. For convenience, we use hinge loss in special cases, where we can directly get the ground-truth NN close form of the optimal solution which has zero loss. For logistic loss, there is no zero-loss solution. We can still show that the OPT value has an exponentially small upper bound at the cost of more computation.*

E.1 Linear Data

Data Distributions. Suppose two labels are equiprobable, i.e., $\mathbb{E}[y = -1] = \mathbb{E}[y = +1] = \frac{1}{2}$. The input data are linearly separable and there is a ground truth direction \mathbf{w}^* , where $\langle \mathbf{w}^*, \mathbf{x} \rangle = \mathfrak{b}$, such that $y \langle \mathbf{w}^*, \mathbf{x} \rangle > 0$. We also assume $\mathbb{E}[y P_{\mathbf{w}^\perp} \mathbf{x}] = 0$, where $P_{\mathbf{w}^\perp}$ is the projection operator on the complementary space of the ground truth, i.e., the components of input data being orthogonal with the ground truth are independent of the label y . We define the input data signal level as $\rho := \mathbb{E}[y \langle \mathbf{w}^*, \mathbf{x} \rangle] > 0$ and the margin as $\beta := \min_{(x,y)} y \langle \mathbf{w}^*, \mathbf{x} \rangle > 0$.

We call this data distribution D_{Linear} .

Lemma E.2 (Linear Data: Gradient Feature Set). Let $\tilde{b} = d B_{x1} \sigma_w$, where τ is any number large enough to satisfy $d^{-2-\frac{1}{4}} > \left(\frac{\sqrt{B_{x2}}}{d}\right)$. For D_{linear} setting, we have $(\mathbf{w}^*, 1) \in S_{p_i; B_G}$ where

$$p = \frac{1}{2}, \quad \gamma = \left(\frac{\rho \overline{B_{x2}}}{\rho d^{-2-\frac{1}{4}}}\right), \quad B_G = \rho \left(\frac{\rho \overline{B_{x2}}}{d^{-2-\frac{1}{4}}}\right). \quad (203)$$

Proof of Lemma E.2. By data distribution, we have

$$\mathbb{E}_{(\mathbf{x}, y)}[y\mathbf{x}] = \rho \mathbf{w}^*. \quad (204)$$

Define $S_{Sure} : \tilde{f}_i \in [m] : k\mathbf{w}_i^{(0)} k_2 \leq 2 \frac{\rho \overline{B_{x2}}}{d \sigma_w}$. For all $i \in [m]$, we have

$$\Pr[i \in S_{Sure}] = \Pr[k\mathbf{w}_i^{(0)} k_2 \leq 2 \frac{\rho \overline{B_{x2}}}{d \sigma_w}] = \frac{1}{2}. \quad (205)$$

For all $i \in S_{Sure}$, by Markov's inequality and considering neuron $i + m$, we have

$$\Pr_{\mathbf{x}} \left[\langle \mathbf{w}_{i+m}^{(0)}, \mathbf{x} \rangle - \mathbf{b}_{i+m} < 0 \right] = \Pr_{\mathbf{x}} \left[\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle + \mathbf{b}_i < 0 \right] \quad (206)$$

$$\Pr_{\mathbf{x}} \left[k\mathbf{w}_i^{(0)} k_2 k\mathbf{x} k_2 - \mathbf{b}_i \right] \quad (207)$$

$$\Pr_{\mathbf{x}} \left[k\mathbf{x} k_2 \leq \frac{d^{-\frac{1}{2}} B_{x1}}{2} \right] \quad (208)$$

$$\left(\frac{1}{d^{-\frac{1}{2}}} \right). \quad (209)$$

For all $i \in S_{Sure}$, by Hölder's inequality, we have

$$\left\| \mathbb{E}_{(\mathbf{x}, y)} \left[y \left(1 - \sigma' \left[\langle \mathbf{w}_{i+m}^{(0)}, \mathbf{x} \rangle - \mathbf{b}_{i+m} \right] \right) \mathbf{x} \right] \right\|_2 \quad (210)$$

$$= \left\| \mathbb{E}_{(\mathbf{x}, y)} \left[y \left(1 - \sigma' \left[\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle + \mathbf{b}_i \right] \right) \mathbf{x} \right] \right\|_2 \quad (211)$$

$$\sqrt{\mathbb{E}[k\mathbf{x} k_2^2] \mathbb{E} \left[\left(1 - \sigma' \left[\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle + \mathbf{b}_i \right] \right)^2 \right]} \quad (212)$$

$$\left(\frac{\rho \overline{B_{x2}}}{d^{-2-\frac{1}{4}}} \right). \quad (213)$$

We have

$$1 - \frac{\left| \langle G(\mathbf{w}_{i+m}^{(0)}, \mathbf{b}_{i+m}), \mathbf{w}^* \rangle \right|}{kG(\mathbf{w}_{i+m}^{(0)}, \mathbf{b}_{i+m}) k_2} = 1 - \frac{\left| \langle G(\mathbf{w}_i^{(0)}, \mathbf{b}_i), \mathbf{w}^* \rangle \right|}{kG(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2} \quad (214)$$

$$1 - \frac{\rho \left(\frac{\sqrt{B_{x2}}}{d^{-2-\frac{1}{4}}} \right)}{\rho + \left(\frac{\sqrt{B_{x2}}}{d^{-2-\frac{1}{4}}} \right)} \quad (215)$$

$$= \left(\frac{\rho \overline{B_{x2}}}{\rho d^{-2-\frac{1}{4}}} \right) = \gamma. \quad (216)$$

We finish the proof by $\frac{\mathbf{b}_{i+m}}{|\mathbf{b}_{i+m}|} = 1$. \square

Lemma E.3 (Linear Data: Existence of Good Networks). Assume the same conditions as in Lemma E.2. Define

$$f^*(\mathbf{x}) = \frac{1}{\beta} \sigma(h\mathbf{w}^*, \mathbf{x}) - \frac{1}{\beta} \sigma(h - \mathbf{w}^*, \mathbf{x}). \quad (217)$$

For D_{linear} setting, we have $f^* \in F_{d; r; B_F; S_{p_i; B_G}}$, where $r = 2, B_F = (B_{a1}, B_{a2}, B_b) = \left(1, \frac{\sqrt{2}}{B_{x1}}, \frac{1}{B_{x1}}\right)$, $p = \frac{1}{2}$, $\gamma = \left(\frac{\sqrt{B_{x2}}}{d^{-2-\frac{1}{4}}}\right)$, $B_G = \rho \left(\frac{\sqrt{B_{x2}}}{d^{-2-\frac{1}{4}}}\right)$. We also have $\text{OPT}_{d; r; B_F; S_{p_i; B_G}} = 0$.

Proof of Lemma E.3. By Lemma E.2 and Lemma E.3, we have $f^* \in \mathcal{F}_{d;r;B_F;S_{p_i};B_G}$. We also have

$$\text{OPT}_{d;r;B_F;S_{p_i};B_G} \leq L_{\mathcal{D}_{\text{linear}}}(f^*) \quad (218)$$

$$= \mathbb{E}_{(\mathbf{x};y) \sim \mathcal{D}_{\text{linear}}} L_{(\mathbf{x};y)}(f^*) \quad (219)$$

$$= 0. \quad (220)$$

□

Theorem E.4 (Linear Data: Main Result). For $\mathcal{D}_{\text{linear}}$ setting, for any $\delta \in (0, 1)$ and for any $\epsilon \in (0, 1)$ when

$$m = \text{poly}\left(\frac{1}{\delta}, \frac{1}{\epsilon}, \frac{1}{\beta}, \frac{1}{\rho}\right) e^d, \quad T = \text{poly}(m, B_{X1}), \quad n = \text{poly}\left(m, B_X, \frac{1}{\delta}, \frac{1}{\epsilon}, \frac{1}{\beta}, \frac{1}{\rho}\right), \quad (221)$$

trained by Algorithm 1 with hinge loss, with probability at least $1 - \delta$ over the initialization, with proper hyper-parameters, there exists $t \geq T$ such that

$$\Pr[\text{sign}(f_{(t)}(\mathbf{x})) \neq y] \leq \epsilon. \quad (222)$$

Proof of Theorem E.4. Let $\tilde{b} = d B_{X1} \sigma_w$, where τ is a number large enough to satisfy $d^{\tau-1/4} > \left(\frac{\sqrt{B_{X2}}}{\sqrt{d}}\right)$ and $O\left(\frac{B_{X1} B_{X2}^{1/4}}{\sqrt{d}^{\tau-1/8}}\right) \leq \frac{\epsilon}{2}$. By Lemma E.3, we have $f^* \in \mathcal{F}_{d;r;B_F;S_{p_i};B_G}$, where

$$r = 2, B_F = (B_{a1}, B_{a2}, B_b) = \left(\frac{1}{\beta}, \frac{\sqrt{2}}{B_{X1}}, \frac{1}{B_{X1}}\right), p = \frac{1}{2}, \gamma = \left(\frac{\sqrt{B_{X2}}}{d^{\tau-1/4}}\right), B_G = \rho \left(\frac{\sqrt{B_{X2}}}{d^{\tau-1/4}}\right).$$

We also have $\text{OPT}_{d;r;B_F;S_{p_i};B_G} = 0$.

Adjust σ_w such that $\tilde{b} = d B_{X1} \sigma_w = \left(\frac{B_G^{1/2} B_{a2} B_b^{3/4}}{\sqrt{r} B_{a1}}\right)$. Injecting above parameters into Theorem 3.12, we have with probability at least $1 - \delta$ over the initialization, with proper hyper-parameters, there exists $t \geq T$ such that

$$\Pr[\text{sign}(f_{(t)}(\mathbf{x})) \neq y] \leq O\left(\frac{B_{X1} B_{X2}^{1/4}}{\beta^{\tau} \rho d^{\tau-1/8}}\right) + O\left(\frac{B_{X1} B_{X2}^{1/4} (\log n)^{1/4}}{\beta^{\tau} \rho n^{1/4}}\right) + \epsilon/2 \leq \epsilon. \quad (223)$$

□

E.2 Mixture of Gaussians

We recap the problem setup in Section 4.1 for readers' convenience.

E.2.1 Problem Setup

Data Distributions. We follow the notations from [99]. The data are from a mixture of r high-dimensional Gaussians, and each Gaussian is assigned to one of two possible labels in $\mathcal{Y} = \{f, g\}$. Let $S(y) \subseteq [r]$ denote the set of indices of the Gaussians associated with the label y . The data distribution is then:

$$q(\mathbf{x}, y) = q(y)q(\mathbf{x}|y), \quad q(\mathbf{x}|y) = \sum_{j \in S(y)} p_j N_j(\mathbf{x}), \quad (224)$$

where $N_j(\mathbf{x})$ is a multivariate normal distribution with mean μ_j and covariance Σ_j , and p_j are chosen such that $q(\mathbf{x}, y)$ is correctly normalized.

We call this data distribution $\mathcal{D}_{\text{mixture}}$.

We will make some assumptions about the Gaussians, for which we first introduce some notations. For all $j \in [r]$, let $y_{(j)} \in \{f, g\}$ be the label for $N_j(\mathbf{x})$.

$$D_j := \frac{\mu_j}{k_{\mu_j} k_2}, \quad \mu_j := \mu_j / \rho d, \quad B_1 := \min_{j \in [r]} k_{\mu_j} k_2, \quad B_2 := \max_{j \in [r]} k_{\mu_j} k_2, \quad p_B := \min_{j \in [r]} p_j.$$

Assumption E.5 (Mixture of Gaussians. Recap of Assumption 4.1). Let $\delta \leq \tau \leq d$ be a parameter that will control our final error guarantee. Assume

- Equiprobable labels: $q(-1) = q(+1) = 1/2$.
- For all $j \geq [r]$, $\sigma_j = \sigma_j I_{d \times d}$. Let $\sigma_B := \max_{j \in [r]} \sigma_j$ and $\sigma_{B^+} := \max_{B \geq 2} \sigma_B$.
- $r \geq 2d$, $p_B = \frac{1}{2d}$, $\left(\frac{1}{d} + \sqrt{\frac{B^+ \log d}{d}} \right) \leq B \leq \frac{1}{2} d$.
- The Gaussians are well-separated: for all $i \neq j \geq [r]$, we have $\frac{1}{2} \geq \langle \mathbf{D}_i, \mathbf{D}_j \rangle \geq \theta$, where $0 < \theta \leq \min \left\{ \frac{1}{2r}, \frac{B^+}{B} \sqrt{\frac{\log d}{d}} \right\}$.

Below, we define a sufficient condition that randomly initialized weights will fall in nice gradients set after the first gradient step update.

Definition E.6 (Mixture of Gaussians: Subset of Nice Gradients Set). Recall $\mathbf{w}_i^{(0)}$ is the weight for the i -th neuron at initialization. For all $j \geq [r]$, let $S_{D_j; \text{Sure}} \subseteq [m]$ be those neurons that satisfy

- $\langle \mathbf{w}_i^{(0)}, \mu_j \rangle \geq C_{\text{Sure},1} \mathbf{b}_i$,
- $\langle \mathbf{w}_i^{(0)}, \mu_{j^0} \rangle \leq C_{\text{Sure},2} \mathbf{b}_i$, for all $j' \neq j, j' \geq [r]$.
- $\|\mathbf{w}_i^{(0)}\|_2 \leq \left(\frac{p_B}{d} \sigma_w \right)$.

E.2.2 Mixture of Gaussians: Feature Learning

We show the important Lemma E.7 first and defer other Lemmas after it.

Lemma E.7 (Mixture of Gaussians: Gradient Feature Set. Part statement of Lemma 4.3). Let $C_{\text{Sure},1} = \frac{3}{2}$, $C_{\text{Sure},2} = \frac{1}{2}$, $\mathfrak{b} = C_b \frac{p_B}{\tau d \log d \sigma_w \sigma_{B^+}}$, where C_b is a large enough universal constant. For D_{mixture} setting, we have $(D_j, +1) \geq S_{p; :B_G}$ for all $j \geq [r]$, where

$$p = \left(\frac{B \leq \frac{1}{2} d}{\frac{p_B}{\tau \log d \sigma_{B^+}} d^{(9C_b^2 - B^+ - (2B^2)_1)}} \right), \quad \gamma = \frac{1}{d^{0.9 - 1.5}}, \quad (225)$$

$$B_G = p_B B \leq \frac{1}{d} = O\left(\frac{\sigma_{B^+}}{d^{0.9}}\right). \quad (226)$$

Proof of Lemma E.7. For all $j \geq [r]$, by Lemma E.10, for all $i \geq S_{D_j; \text{Sure}}$,

$$\frac{1}{\|\mathbf{G}(\mathbf{w}_i^{(0)}, \mathbf{b}_i)\|_2} \left| \langle \mathbf{G}(\mathbf{w}_i^{(0)}, \mathbf{b}_i), D_j \rangle \right| \quad (227)$$

$$\frac{1}{\sqrt{\left| \langle \mathbf{G}(\mathbf{w}_i^{(0)}, \mathbf{b}_i), D_j \rangle \right|^2 + \max_{D_j^\perp: D_j^\perp = 0, \|D_j^\perp\|_2 = 1} \left| \langle \mathbf{G}(\mathbf{w}_i^{(0)}, \mathbf{b}_i), D_j^\perp \rangle \right|^2}} \left| \langle \mathbf{G}(\mathbf{w}_i^{(0)}, \mathbf{b}_i), D_j \rangle \right| \quad (228)$$

$$\frac{1}{\left| \langle \mathbf{G}(\mathbf{w}_i^{(0)}, \mathbf{b}_i), D_j \rangle \right| + \max_{D_j^\perp: D_j^\perp = 0, \|D_j^\perp\|_2 = 1} \left| \langle \mathbf{G}(\mathbf{w}_i^{(0)}, \mathbf{b}_i), D_j^\perp \rangle \right|}} \left| \langle \mathbf{G}(\mathbf{w}_i^{(0)}, \mathbf{b}_i), D_j \rangle \right| \quad (229)$$

$$\frac{1}{1 + \frac{B \leq \frac{1}{2} d \cdot O\left(\frac{1}{d^{\frac{1}{2}}}\right) + B^+ \cdot O\left(\frac{1}{d^{0.9}}\right)}{p_B B \leq \frac{1}{d} \sqrt{d(1 - O\left(\frac{1}{d}\right))} - B \leq \frac{1}{2} d \cdot O\left(\frac{1}{d^{\frac{1}{2}}}\right) - B^+ \cdot O\left(\frac{1}{d^{0.9}}\right)}} \quad (230)$$

$$\frac{\frac{\sigma_{B^+} \cdot O\left(\frac{1}{d^{0.9}}\right)}{p_B B \leq \frac{1}{d}}}{\sigma_{B^+} \cdot O\left(\frac{1}{d^{0.9}}\right)} \quad (231)$$

$$< \frac{1}{d^{0.9 - 1.5}} = \gamma, \quad (232)$$

where the last inequality follows $B \leq \frac{1}{2} d \leq \left(\sigma_{B^+} \sqrt{\frac{\log d}{d}} \right)$.

Thus, we have $G(\mathbf{w}_i^{(0)}, \mathbf{b}_i) \geq C_{D_j}$; and $\left| \langle G(\mathbf{w}_i^{(0)}, \mathbf{b}_i), D_j \rangle \right| \leq kG(\mathbf{w}_i^{(0)}, \mathbf{b}_i)k_2 \leq B_{x1}, \frac{b_i}{\|\mathbf{b}_i\|} = +1$. Thus, by Lemma E.8, we have

$$\Pr_{\mathbf{w}; b} \left[G(\mathbf{w}, b) \geq C_{D_j}; \text{ and } kG(\mathbf{w}, b)k_2 \leq B_G \text{ and } \frac{b}{\|b\|} = +1 \right] \quad (233)$$

$$\Pr [i \geq S_{D_j}; \text{Sure}] \quad (234)$$

$$p. \quad (235)$$

Thus, $(D_j, +1) \geq S_{p_i}; B_G$. We finish the proof. \square

Below are Lemmas used in the proof of Lemma E.7. In Lemma E.8, we calculate p used in $S_{p_i}; B_G$.

Lemma E.8 (Mixture of Gaussians: Geometry at Initialization. Lemma B.2 in [7]). *Assume the same conditions as in Lemma E.7, recall for all $i \geq [m]$, $\mathbf{w}_i^{(0)} \sim \mathcal{N}(0, \sigma_w^2 I_{d \times d})$, over the random initialization, we have for all $i \geq [m], j \geq [r]$,*

$$\Pr [i \geq S_{D_j}; \text{Sure}] \leq \left(\frac{B_{x1}}{p \frac{1}{\tau} \log d \sigma_{B+}} \frac{1}{d^{(9C_b^2 - B_{x+}^2 - (2B^2 - 1))}} \right). \quad (236)$$

Proof of Lemma E.8. Recall for all $l \geq [r]$, $\mu_l = \mu_l / \rho_{\bar{d}}$.

WLOG, let $j = r$. For all $l \geq [r - 1]$. We define $Z_1 = \{l \geq [r - 1] : \langle \mathbf{w}_l, D_r \rangle \geq \theta g\}$ and $Z_2 = \{l \geq [r - 1] : -1 < \langle \mathbf{w}_l, D_r \rangle < \theta g\}$. WLOG, let $Z_1 = [r_1]$, $Z_2 = [r_1 + 1, \dots, r_2]$, where $0 \leq r_1 \leq r_2 \leq r - 1$. We define the following events

$$\zeta_l = \left\{ \langle \mathbf{w}_l^{(0)}, \mu_l \rangle \geq C_{\text{Sure}; 2} \mathbf{b}_l \right\}, \hat{\zeta}_l = \left\{ \left| \langle \mathbf{w}_l^{(0)}, \mu_l \rangle \right| \geq C_{\text{Sure}; 2} \mathbf{b}_l \right\}. \quad (237)$$

We define space $A = \text{span}(\mu_1, \dots, \mu_{r_1})$ and $\underline{\mu}_l = P_{A^\perp} \mu_l$, where P_{A^\perp} is the projection operator on the complementary space of A . For $l \geq Z_2$, we also define $\underline{\mu}_l = \mu_l - \frac{\langle \mu_l, \mu_r \rangle}{\|\mu_r\|_2^2} \mu_r$, and the event

$$\zeta_l = \left\{ \langle \mathbf{w}_l^{(0)}, \underline{\mu}_l \rangle \geq C_{\text{Sure}; 2} \mathbf{b}_l \right\}, \hat{\zeta}_l = \left\{ \left| \langle \mathbf{w}_l^{(0)}, \underline{\mu}_l \rangle \right| \geq C_{\text{Sure}; 2} \mathbf{b}_l \right\}. \quad (238)$$

For $l \geq Z_2$, we have $\mu_l = \underline{\mu}_l + \rho \mu_r$, where $\rho \geq 0$. So $\langle \mathbf{w}_l, \mu_l \rangle = \langle \mathbf{w}_l, \underline{\mu}_l \rangle + \rho \langle \mathbf{w}_l, \mu_r \rangle \geq \langle \mathbf{w}_l, \underline{\mu}_l \rangle$ when $\langle \mathbf{w}_l, \mu_r \rangle \geq 0$. As a result, we have

$$\zeta_l \setminus \left\{ \langle \mathbf{w}_l^{(0)}, \mu_r \rangle \geq C_{\text{Sure}; 1} \mathbf{b}_l \right\} \subseteq \hat{\zeta}_l \setminus \left\{ \langle \mathbf{w}_l^{(0)}, \mu_r \rangle \geq C_{\text{Sure}; 1} \mathbf{b}_l \right\}. \quad (239)$$

By Assumption 4.1, we have

$$\frac{1}{2} \leq 1 - r\theta \leq 1 - r_1\theta \leq \frac{k\underline{\mu}_r k_2}{k\mu_r k_2} \leq 1. \quad (240)$$

We also have,

$$\Pr \left[\langle \mathbf{w}_i^{(0)}, \mu_r \rangle \geq C_{\text{Sure}; 1} \mathbf{b}_i, \zeta_1, \dots, \zeta_{r-1} \right] \quad (241)$$

$$= \Pr \left[\langle \mathbf{w}_i^{(0)}, \mu_r \rangle \geq C_{\text{Sure}; 1} \mathbf{b}_i, \zeta_1, \dots, \zeta_{r_2} \right] \quad (242)$$

$$\Pr \left[\langle \mathbf{w}_i^{(0)}, \mu_r \rangle \geq C_{\text{Sure}; 1} \mathbf{b}_i, \zeta_1, \dots, \zeta_{r_1}, \zeta_{r_1+1}, \dots, \zeta_{r_2} \right] \quad (243)$$

$$\Pr \left[\langle \mathbf{w}_i^{(0)}, \mu_r \rangle \geq C_{\text{Sure}; 1} \mathbf{b}_i, \hat{\zeta}_1, \dots, \hat{\zeta}_{r_1}, \hat{\zeta}_{r_1+1}, \dots, \hat{\zeta}_{r_2} \right] \quad (244)$$

$$= \Pr \left[\underbrace{\langle \mathbf{w}_i^{(0)}, \mu_r \rangle \geq C_{\text{Sure}; 1} \mathbf{b}_i}_{p_r} \left| \hat{\zeta}_1, \dots, \hat{\zeta}_{r_1}, \hat{\zeta}_{r_1+1}, \dots, \hat{\zeta}_{r_2} \right. \right] \Pr \left[\hat{\zeta}_1, \dots, \hat{\zeta}_{r_1}, \hat{\zeta}_{r_1+1}, \dots, \hat{\zeta}_{r_2} \right].$$

$1 \geq [r_2] p_l$

For the first condition in Definition E.6, we have,

$$p_r = \Pr \left[\left\langle \mathbf{w}_i^{(0)}, \mu_r \right\rangle C_{\text{Sure};1} \mathbf{b}_i \middle| \hat{\zeta}_1, \dots, \hat{\zeta}_{r_1}, \hat{\zeta}_{r_1+1}, \dots, \hat{\zeta}_{r_2} \right] \quad (245)$$

$$= \Pr \left[\left\langle \mathbf{w}_i^{(0)}, \hat{\mu}_r + \mu_r \right\rangle C_{\text{Sure};1} \mathbf{b}_i \middle| \hat{\zeta}_1, \dots, \hat{\zeta}_{r_1} \right] \quad (246)$$

$$\Pr \left[\left\langle \mathbf{w}_i^{(0)}, \hat{\mu}_r + \mu_r \right\rangle C_{\text{Sure};1} \mathbf{b}_i, \left\langle \mathbf{w}_i^{(0)}, \mu_r \right\rangle = 0 \middle| \hat{\zeta}_1, \dots, \hat{\zeta}_{r_1} \right] \quad (247)$$

$$= \Pr \left[\left\langle \mathbf{w}_i^{(0)}, \hat{\mu}_r + \mu_r \right\rangle C_{\text{Sure};1} \mathbf{b}_i \middle| \left\langle \mathbf{w}_i^{(0)}, \mu_r \right\rangle = 0, \hat{\zeta}_1, \dots, \hat{\zeta}_{r_1} \right] \quad (248)$$

$$\Pr \left[\left\langle \mathbf{w}_i^{(0)}, \mu_r \right\rangle = 0 \middle| \hat{\zeta}_1, \dots, \hat{\zeta}_{r_1} \right] \quad (249)$$

$$= \frac{1}{2} \Pr \left[\left\langle \mathbf{w}_i^{(0)}, \hat{\mu}_r + \mu_r \right\rangle C_{\text{Sure};1} \mathbf{b}_i \middle| \left\langle \mathbf{w}_i^{(0)}, \mu_r \right\rangle = 0, \hat{\zeta}_1, \dots, \hat{\zeta}_{r_1} \right] \quad (250)$$

$$\frac{1}{2} \Pr \left[\left\langle \mathbf{w}_i^{(0)}, \hat{\mu}_r \right\rangle C_{\text{Sure};1} \mathbf{b}_i \middle| \left\langle \mathbf{w}_i^{(0)}, \mu_r \right\rangle = 0, \hat{\zeta}_1, \dots, \hat{\zeta}_{r_1} \right] \quad (251)$$

$$= \frac{1}{2} \Pr \left[\left\langle \mathbf{w}_i^{(0)}, \hat{\mu}_r \right\rangle C_{\text{Sure};1} \mathbf{b}_i \right] \quad (252)$$

$$\left(\frac{k_{\mu_r} k_2}{\rho \tau \log d \sigma_{B_+}} d^{(9C_b^2 - B_+^2 = (2\|\cdot\|_2^2))} \right), \quad (253)$$

where the last equality following that $\hat{\mu}_r$ is orthogonal with μ_1, \dots, μ_{r_1} and the property of the standard Gaussian vector, and the last inequality follows Lemma F.6.

For the second condition in Definition E.6, by Lemma F.6, we have,

$$p_1 = \Pr \left[\hat{\zeta}_1 \right] = 1 \left(\frac{k_{\mu_1} k_2}{\rho \tau \log d \sigma_{B_+}} d^{(C_b^2 - B_+^2 = (8\|\cdot\|_2^2))} \right) \quad (254)$$

$$p_2 = \Pr \left[\hat{\zeta}_2 \middle| \hat{\zeta}_1 \right] = \Pr \left[\hat{\zeta}_2 \right] = 1 \left(\frac{k_{\mu_2} k_2}{\rho \tau \log d \sigma_{B_+}} d^{(C_b^2 - B_+^2 = (8\|\cdot\|_2^2))} \right) \quad (255)$$

$$\vdots \quad (256)$$

$$p_{r-1} = \Pr \left[\hat{\zeta}_{r_2} \middle| \hat{\zeta}_1, \dots, \hat{\zeta}_{r_1}, \hat{\zeta}_{r_1+1}, \dots, \hat{\zeta}_{r_2} \right] = \Pr \left[\hat{\zeta}_{r_2} \right] = \Pr \left[\hat{\zeta}_{r_2} \right] \quad (257)$$

$$1 \left(\frac{k_{\mu_{r-1}} k_2}{\rho \tau \log d \sigma_{B_+}} d^{(C_b^2 - B_+^2 = (8\|\cdot\|_2^2))} \right). \quad (258)$$

On the other hand, if X is a $\chi^2(k)$ random variable. Then we have

$$\Pr(X \leq k + 2\sqrt{kx} + 2x) \leq e^{-x}. \quad (259)$$

Therefore, by assumption $B_1 \leq \left(\sigma_{B_+} \sqrt{\frac{\log d}{d}} \right)$, we have

$$\Pr \left(\frac{1}{\sigma_w^2} \left\| \mathbf{w}_i^{(0)} \right\|_2^2 \leq d + 2\sqrt{(9C_b^2 \tau \sigma_{B_+}^2 / (2B_1^2) + 2) d \log d} \right) \quad (260)$$

$$+ 2(9C_b^2 \tau \sigma_{B_+}^2 / (2B_1^2) + 2) \log d \quad (261)$$

$$O \left(\frac{1}{d^2 d^{(9C_b^2 - B_+^2 = (2B_1^2))}} \right). \quad (262)$$

Recall $B_1 = \min_{j \in [r]} k_{\mu_j} k_2$, $B_2 = \max_{j \in [r]} k_{\mu_j} k_2$. Thus, by union bound, we have

$$\Pr [i \geq S_{D_j; \text{Sure}}] \quad (263)$$

$$\leq O\left(\frac{1}{d^2 d^{(9C_b^2 B_+^2 = (2B^2_1))}}\right) \quad (264)$$

$$\left(\frac{B_1}{\rho \frac{1}{\tau \log d} \sigma_{B_+} d^{(9C_b^2 B_+^2 = (2B^2_1))}} \left(1 - \frac{rB_2}{\rho \frac{1}{\tau \log d} \sigma_{B_+} d^{(C_b^2 B_+^2 = (8B^2_2))}}\right)\right) \quad (265)$$

$$\leq O\left(\frac{1}{d^2 d^{(9C_b^2 B_+^2 = (2B^2_1))}}\right) \quad (266)$$

$$\left(\frac{B_1}{\rho \frac{1}{\tau \log d} \sigma_{B_+} d^{(9C_b^2 B_+^2 = (2B^2_1))}}\right). \quad (267)$$

□

In Lemma E.9, we compute the activation pattern for the neurons in $S_{D_j; \text{Sure}}$.

Lemma E.9 (Mixture of Gaussians: Activation Pattern). *Assume the same conditions as in Lemma E.7, for all $j \geq [r]$, $i \geq S_{D_j; \text{Sure}}$, we have*

(1) When $\mathbf{x} \sim \mathcal{N}_j(\mu_j, \sigma_j I_{d \times d})$, the activation probability satisfies,

$$\Pr_{\mathbf{x} \sim \mathcal{N}_j(\mu_j; \sigma_j I_{d \times d})} [\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \geq \mathbf{b}_i] \geq 1 - O\left(\frac{1}{d}\right). \quad (268)$$

(2) For all $j' \neq j, j' \geq [r]$, when $\mathbf{x} \sim \mathcal{N}_{j'}(\mu_{j'}, \sigma_{j'} I_{d \times d})$, the activation probability satisfies,

$$\Pr_{\mathbf{x} \sim \mathcal{N}_{j'}(\mu_{j'}; \sigma_{j'} I_{d \times d})} [\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \geq \mathbf{b}_i] \leq O\left(\frac{1}{d}\right). \quad (269)$$

Proof of Lemma E.9. In the proof, we need $\tilde{b} = C_b \frac{\rho}{\tau d \log d} \sigma_w \sigma_{B_+}$, where C_b is a large enough universal constant. For the first statement, when $\mathbf{x} \sim \mathcal{N}_j(\mu_j, \sigma_j I_{d \times d})$, by $C_{\text{Sure};1} \geq \frac{3}{2}$, we have

$$\Pr_{\mathbf{x} \sim \mathcal{N}_j(\mu_j; \sigma_j I_{d \times d})} [\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \geq \mathbf{b}_i] \geq \Pr_{\mathbf{x} \sim \mathcal{N}(0; \sigma_j I_{d \times d})} [\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \geq (1 - C_{\text{Sure};1}) \mathbf{b}_i] \quad (270)$$

$$\geq \Pr_{\mathbf{x} \sim \mathcal{N}(0; \sigma_j I_{d \times d})} [\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \geq \frac{\mathbf{b}_i}{2}] \quad (271)$$

$$= 1 - \Pr_{\mathbf{x} \sim \mathcal{N}(0; \sigma_j I_{d \times d})} [\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \leq -\frac{\mathbf{b}_i}{2}] \quad (272)$$

$$\leq 1 - \exp\left(-\frac{\mathbf{b}_i^2}{(d\sigma_w^2\sigma_j^2)}\right) \quad (273)$$

$$\geq 1 - O\left(\frac{1}{d}\right), \quad (274)$$

where the third inequality follows the Chernoff bound and symmetricity of the Gaussian vector.

For the second statement, we prove similarly by $0 < C_{\text{Sure};2} \leq \frac{1}{2}$. □

Then, Lemma E.10 gives gradients of neurons in $S_{D_j; \text{Sure}}$. It shows that these gradients are highly aligned with D_j .

Lemma E.10 (Mixture of Gaussians: Feature Emergence). *Assume the same conditions as in Lemma E.7, for all $j \geq [r]$, $i \geq S_{D_j; \text{Sure}}$, we have*

$$\left\langle \mathbb{E}_{(\mathbf{x}; y)} \left[y \sigma' \left(\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \geq \mathbf{b}_i \right) \mathbf{x} \right], y_{(j)} D_j \right\rangle \quad (275)$$

$$\geq \rho B_1 \frac{\rho}{d} \left(1 - O\left(\frac{1}{d}\right)\right) \geq B_2 O\left(\frac{1}{d^{-\frac{1}{2}}}\right) \geq \sigma_{B_+} O\left(\frac{1}{d^{0.9}}\right). \quad (276)$$

For any unit vector D_j^\perp which is orthogonal with D_j , we have

$$\left| \left\langle \mathbb{E}_{(\mathbf{x};y)} \left[y \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \mathbf{x} \right], D_j^\perp \right\rangle \right| \leq B_2 O \left(\frac{1}{d^{-\frac{1}{2}}} \right) + \sigma_{B+} O \left(\frac{1}{d^{0.9}} \right). \quad (277)$$

Proof of Lemma E.10. For all $j \geq [r]$, $i \geq S_{D_j; \text{Sure}}$, we have

$$\mathbb{E}_{(\mathbf{x};y)} \left[y \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \mathbf{x} \right] \quad (278)$$

$$= \sum_{l \in [r]} p_l \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{x})} \left[y \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \mathbf{x} \right] \quad (279)$$

$$= \sum_{l \in [r]} p_l y(l) \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \frac{1}{d} \mathbf{I}_d)} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_l \right\rangle - \mathbf{b}_i \right) (\mathbf{x} + \mu_l) \right]. \quad (280)$$

Thus, by Lemma F.7 and Lemma E.9,

$$\left\langle \mathbb{E}_{(\mathbf{x};y)} \left[y \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \mathbf{x} \right], y(j) D_j \right\rangle \quad (281)$$

$$= p_j \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \frac{1}{d} \mathbf{I}_d)} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_j \right\rangle - \mathbf{b}_i \right) (\mathbf{x} + \mu_j)^\top D_j \right] \quad (282)$$

$$+ \sum_{l \in [r]; l \neq j} p_l y(l) y(j) \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \frac{1}{d} \mathbf{I}_d)} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_l \right\rangle - \mathbf{b}_i \right) (\mathbf{x} + \mu_l)^\top D_j \right] \quad (283)$$

$$p_j \mu_j^\top D_j \left(1 - O \left(\frac{1}{d} \right) \right) - \sum_{l \in [r]; l \neq j} p_l \mu_l^\top D_j j O \left(\frac{1}{d} \right) \quad (284)$$

$$p_j \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \frac{1}{d} \mathbf{I}_d)} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_j \right\rangle - \mathbf{b}_i \right) \mathbf{x}^\top D_j \right] \right| \quad (285)$$

$$\sum_{l \in [r]; l \neq j} p_l \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \frac{1}{d} \mathbf{I}_d)} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_l \right\rangle - \mathbf{b}_i \right) \mathbf{x}^\top D_j \right] \right| \quad (286)$$

$$p_j B_1 \frac{1}{d} \left(1 - O \left(\frac{1}{d} \right) \right) - B_2 O \left(\frac{1}{d^{-\frac{1}{2}}} \right) \quad (287)$$

$$p_j \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \frac{1}{d} \mathbf{I}_d)} \left[\left(1 - \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_j \right\rangle - \mathbf{b}_i \right) \right) \mathbf{x}^\top D_j \right] \right| \quad (288)$$

$$\sum_{l \in [r]; l \neq j} p_l \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \frac{1}{d} \mathbf{I}_d)} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_l \right\rangle - \mathbf{b}_i \right) \mathbf{x}^\top D_j \right] \right| \quad (289)$$

$$= p_j B_1 \frac{1}{d} \left(1 - O \left(\frac{1}{d} \right) \right) - B_2 O \left(\frac{1}{d^{-\frac{1}{2}}} \right) \quad (290)$$

$$p_j \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \frac{1}{d} \mathbf{I}_d)} \left[\left(1 - \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_j \right\rangle - \mathbf{b}_i \right) \right) \mathbf{x}^\top D_j \right] \right| \quad (291)$$

$$\sum_{l \in [r]; l \neq j} p_l \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \frac{1}{d} \mathbf{I}_d)} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_l \right\rangle - \mathbf{b}_i \right) \mathbf{x}^\top D_j \right] \right| \quad (292)$$

$$p_j B_1 \frac{1}{d} \left(1 - O \left(\frac{1}{d} \right) \right) - B_2 O \left(\frac{1}{d^{-\frac{1}{2}}} \right) - \sigma_{B+} O \left(\frac{1}{d^{0.9}} \right). \quad (293)$$

For any unit vector D_j^\perp which is orthogonal with D_j , similarly, we have

$$\left| \left\langle \mathbb{E}_{(\mathbf{x}, y)} \left[y \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \mathbf{x} \right], D_j^\perp \right\rangle \right| \quad (294)$$

$$p_j \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \cdot, I)} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_j \right\rangle - \mathbf{b}_i \right) \mathbf{x}^\top D_j^\perp \right] \right| \quad (295)$$

$$+ \sum_{l \in [r]: l \neq j} p_l \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \cdot, I)} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_l \right\rangle - \mathbf{b}_i \right) (\mathbf{x} + \mu_l)^\top D_j^\perp \right] \right| \quad (296)$$

$$B_{2O} \left(\frac{1}{d^{-\frac{1}{2}}} \right) + p_j \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \cdot, I)} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_j \right\rangle - \mathbf{b}_i \right) \mathbf{x}^\top D_j^\perp \right] \right| \quad (297)$$

$$+ \sum_{l \in [r]: l \neq j} p_l \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \cdot, I)} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_l \right\rangle - \mathbf{b}_i \right) \mathbf{x}^\top D_j^\perp \right] \right| \quad (298)$$

$$B_{2O} \left(\frac{1}{d^{-\frac{1}{2}}} \right) + \sigma_{B+} O \left(\frac{1}{d^{0.9}} \right). \quad (299)$$

□

E.2.3 Mixture of Gaussians: Final Guarantee

Lemma E.11 (Mixture of Gaussians: Existence of Good Networks. Part statement of Lemma 4.3). *Assume the same conditions as in Lemma E.7. Define*

$$f^*(\mathbf{x}) = \sum_{j=1}^r \frac{y^{(j)}}{\tau \log d \sigma_{B+}} \left[\sigma \left(h D_j, \mathbf{x} i - 2\sqrt{\tau \log d} \sigma_{B+} \right) \right]. \quad (300)$$

For D_{mixture} setting, we have $f^* \geq 2 F_{d; r; B_F; S_{p_i}; B_G}$, where $B_F = (B_{a1}, B_{a2}, B_b) = \left(\frac{1}{\sqrt{\log d} B_+}, \frac{\sqrt{r}}{\sqrt{\log d} B_+}, 2^{\rho} \tau \log d \sigma_{B+} \right)$, $p = \left(\frac{B_1}{\sqrt{\log d} B_+}, \beta_{(9C_b^2 B_+^2 - (2B^2, 1))} \right)$, $\gamma = \frac{1}{d^{0.9} 1.5}$, $B_G = p B B_1 \bar{d} = O \left(\frac{B_+}{d^{0.9}} \right)$ and $B_{X1} = (B_2 + \sigma_{B+}) \bar{d}$, $B_{X2} = (B_2 + \sigma_{B+})^2 d$. We also have $\text{OPT}_{d; r; B_F; S_{p_i}; B_G} \geq \frac{3}{d} + \frac{4}{d^{0.9} 1 \sqrt{\log d}}$.

Proof of Lemma E.11. We can check $B_{X1} = (B_2 + \sigma_{B+}) \bar{d}$, $B_{X2} = (B_2 + \sigma_{B+})^2 d$ by direct calculation. By Lemma E.7, we have $f^* \geq 2 F_{d; r; B_F; S_{p_i}; B_G}$.

For any $j \geq [r]$, by $B_1 \left(\sigma_{B+} \sqrt{\frac{\log d}{d}} \right) \leq 4\sigma_{B+} \sqrt{\frac{\log d}{d}}$, we have

$$\Pr_{\mathbf{x} \sim \mathcal{N}_j \left(\frac{j}{d} \right)} \left[h D_j, \mathbf{x} i - 2\sqrt{\tau \log d} \sigma_{B+} \leq \sqrt{\tau \log d} \sigma_{B+} \right] \quad (301)$$

$$= \Pr_{\mathbf{x} \sim \mathcal{N}_j(0; \cdot, I_{d \times d})} \left[h D_j, \mathbf{x} i + k \mu_j k_2 - 2\sqrt{\tau \log d} \sigma_{B+} \leq \sqrt{\tau \log d} \sigma_{B+} \right] \quad (302)$$

$$\Pr_{\mathbf{x} \sim \mathcal{N}_j(0; \cdot, I_{d \times d})} \left[h D_j, \mathbf{x} i + \frac{\rho}{d} B_1 - 2\sqrt{\tau \log d} \sigma_{B+} \leq \sqrt{\tau \log d} \sigma_{B+} \right] \quad (303)$$

$$\Pr_{\mathbf{x} \sim \mathcal{N}_j(0; \cdot, I_{d \times d})} \left[h D_j, \mathbf{x} i \leq \sqrt{\tau \log d} \sigma_{B+} \right] \quad (304)$$

$$1 - \frac{1}{d}, \quad (305)$$

where the last inequality follows Chernoff bound.

For any $l \notin j, l \geq [r]$, by $\theta = \frac{B_+}{B_-} \sqrt{\frac{\log d}{d}}$, we have

$$\Pr_{\mathbf{x} \sim \mathcal{N}_j(\cdot; j, l, d)} \left[hD_l, \mathbf{x} \geq 2\sqrt{\tau \log d} \sigma_{B_+} \right] \quad (306)$$

$$\Pr_{\mathbf{x} \sim \mathcal{N}_j(0; j, l, d)} \left[hD_l, \mathbf{x} + \theta B_- \frac{\rho_-}{d} \geq 2\sqrt{\tau \log d} \sigma_{B_+} \right] \quad (307)$$

$$\Pr_{\mathbf{x} \sim \mathcal{N}_j(0; j, l, d)} \left[hD_l, \mathbf{x} \geq \sqrt{\tau \log d} \sigma_{B_+} \right] \quad (308)$$

$$\frac{1}{d}. \quad (309)$$

Thus, we have

$$\Pr_{(\mathbf{x}, y) \sim \mathcal{D}_{mixture}} [yf^*(\mathbf{x}) > 1] \quad (310)$$

$$\sum_{j \in [r]} p_j \left(\Pr_{\mathbf{x} \sim \mathcal{N}_j(\cdot; j, l, d)} \left[hD_j, \mathbf{x} \geq 2\sqrt{\tau \log d} \sigma_{B_+} + \sqrt{\tau \log d} \sigma_{B_+} \right] \right) \quad (311)$$

$$\sum_{j \in [r]} p_j \left(\sum_{l \neq j, l \in [r]} \Pr_{\mathbf{x} \sim \mathcal{N}_j(\cdot; j, l, d)} \left[hD_l, \mathbf{x} \geq 2\sqrt{\tau \log d} \sigma_{B_+} < 0 \right] \right) \quad (312)$$

$$1 - \frac{2}{d}. \quad (313)$$

We also have

$$\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_{mixture}} [1 - |yf^*(\mathbf{x})|] \quad (314)$$

$$\begin{aligned} & \sum_{j \in [r]} p_j \left(\Pr_{\mathbf{x} \sim \mathcal{N}_j(\cdot; j, l, d)} \left[hD_j, \mathbf{x} \geq 2\sqrt{\tau \log d} \sigma_{B_+} < \sqrt{\tau \log d} \sigma_{B_+} \right] \frac{y_{(j)}^2 \rho_- \tau \log d \sigma_{B_+}}{\rho_- \tau \log d \sigma_{B_+}} \right) \\ & + \sum_{j \in [r]} p_j \left(\sum_{l \neq j, l \in [r]} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_j(\cdot; j, l, d)} \left[\sigma' \left[hD_l, \mathbf{x} \geq 2\sqrt{\tau \log d} \sigma_{B_+} > 0 \right] \frac{hD_l, \mathbf{x} \geq 2\sqrt{\tau \log d} \sigma_{B_+}}{\rho_- \tau \log d \sigma_{B_+}} \right] \right) \\ & \frac{1}{d} + \sum_{j \in [r]} p_j \left(\sum_{l \neq j, l \in [r]} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_j(0; j, l, d)} \left[\sigma' \left[hD_l, \mathbf{x} > \sqrt{\tau \log d} \sigma_{B_+} \right] \frac{hD_l, \mathbf{x} \geq \sqrt{\tau \log d} \sigma_{B_+}}{\rho_- \tau \log d \sigma_{B_+}} \right] \right) \end{aligned} \quad (315)$$

$$\frac{1}{d} + \frac{1}{\rho_- \tau \log d} \sum_{j \in [r]} p_j \left(\sum_{l \neq j, l \in [r]} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_j(0; j, l, d)} \left[\sigma' \left[hD_l, \mathbf{x} > \sqrt{\tau \log d} \right] hD_l, \mathbf{x} \right] \right) \quad (316)$$

$$\frac{1}{d} + \frac{4}{d^{0.9} \rho_- \tau \log d}, \quad (316)$$

where the second last inequality follows Lemma F.7 and $r \leq 2d$. Thus, we have

$$\text{OPT}_{d; r; B_F; S_p; B_G} = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_{mixture}} [\ell(yf^*(\mathbf{x}))] \quad (317)$$

$$= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_{mixture}} [1 - |yf^*(\mathbf{x})|] \quad (318)$$

$$\begin{aligned} & \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_{mixture}} [1 - |yf^*(\mathbf{x})|] + \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_{mixture}} [1 - |yf^*(\mathbf{x})|] \\ & \frac{3}{d} + \frac{4}{d^{0.9} \rho_- \tau \log d}. \end{aligned} \quad (319)$$

□

Theorem 4.4 (Mixtures of Gaussians: Main Result). *Assume Assumption 4.1. For any $\epsilon, \delta \geq (0, 1)$, when Algorithm 1 uses hinge loss with*

$$m = \text{poly} \left(\frac{1}{\delta}, \frac{1}{\epsilon}, d \left(B_+^2 = B_-^2 \right), r, \frac{1}{\rho_B} \right) \quad e^d, \quad T = \text{poly}(m), \quad n = \text{poly}(m)$$

and proper hyper-parameters, then with probability at least $1 - \delta$, there exists $t \geq [T]$ such that

$$\Pr[\text{sign}(f_{(t)}(\mathbf{x})) \neq y] \leq \frac{\rho_{\bar{2}r}}{d^{0.4-0.8}} + \epsilon.$$

Proof of Theorem 4.4. Let $\tilde{b} = C_b \sqrt{\rho_{\bar{2}r} \tau d \log d} \sigma_w \sigma_{B_+}$, where C_b is a large enough universal constant. By Lemma E.11, we have $f^* \in \mathcal{F}_{d;r;B_F;S_{p_j};B_G}$, where $B_F = (B_{a1}, B_{a2}, B_b) = \left(\frac{1}{\sqrt{\log d} B_+}, \frac{\sqrt{r}}{\sqrt{\log d} B_+}, 2 \sqrt{\rho_{\bar{2}r} \tau d \log d} \sigma_{B_+} \right)$, $p = \left(\frac{B_{a1}}{\sqrt{\log d} B_+}, \frac{B_{a2}}{\sqrt{\log d} B_+}, \frac{B_b}{\sqrt{\log d} B_+} \right)$, $\gamma = \frac{1}{d^{0.9-1.5}}$, $B_G = p_B B_{a1} \sqrt{d} = O\left(\frac{B_+}{d^{0.9}}\right)$ and $B_{X1} = (B_{X2} + \sigma_{B_+}) \sqrt{d}$, $B_{X2} = (B_{X2} + \sigma_{B_+})^2 d$. We also have $\text{OPT}_{d;r;B_F;S_{p_j};B_G} \leq \frac{3}{d} + \frac{4}{d^{0.9-1} \sqrt{\log d}}$.

Adjust σ_w such that $\tilde{b} = C_b \sqrt{\rho_{\bar{2}r} \tau d \log d} \sigma_w \sigma_{B_+} = \left(\frac{B_G^{1/4} B_{a2} B_b^{3/4}}{\sqrt{r} B_{a1}} \right)$. Injecting above parameters into Theorem 3.12, we have with probability at least $1 - \delta$ over the initialization, with proper hyper-parameters, there exists $t \geq [T]$ such that

$$\Pr[\text{sign}(f_{(t)}(\mathbf{x})) \neq y] \tag{320}$$

$$\begin{aligned} & \frac{3}{d} + \frac{4}{d^{0.9-1} \sqrt{\log d}} + \frac{\rho_{\bar{2}r} B_{X2}}{d^{(0.9-1.5)2} \sqrt{\log d} \sigma_{B_+}} + O\left(\frac{r B_{a1} B_{X1} B_{X2}^{1/4} (\log n)^{1/4}}{B_G n^{1/4}}\right) + \epsilon/2 \\ & \leq \frac{\rho_{\bar{2}r}}{d^{0.4-0.8}} + \epsilon. \end{aligned} \tag{321}$$

□

E.3 Mixture of Gaussians - XOR

We consider a special Mixture of Gaussians distribution studied in [99]. Consider the same data distribution in Appendix E.2.1 and Definition E.6 with the following assumptions.

Assumption E.12 (Mixture of Gaussians in [99]). *Assume four Gaussians cluster with XOR-like pattern, for any $\tau > 0$,*

- $r = 4$ and $p_1 = p_2 = p_3 = p_4 = \frac{1}{4}$.
- $\mu_1 = \mu_2, \mu_3 = \mu_4$ and $k\mu_1 k_2 = k\mu_3 k_2 = k\mu_4 k_2 = \frac{\rho_{\bar{d}}}{d}$ and $\langle \mu_1, \mu_3 \rangle = 0$.
- For all $j \geq [4]$, $\mu_j = \sigma_B I_{d \times d}$ and $1 - \sigma_B = \sqrt{\frac{d}{\log \log d}}$.
- $y_{(1)} = y_{(2)} = 1$ and $y_{(3)} = y_{(4)} = -1$.

We denote this data distribution as $D_{\text{mixture-xor}}$ setting.

E.3.1 Mixture of Gaussians - XOR: Feature Learning

Lemma E.13 (Mixture of Gaussians in [99]: Gradient Feature Set). *Let $C_{\text{Sure},1} = \frac{6}{5}$, $C_{\text{Sure},2} = \frac{\sqrt{2}}{\sqrt{\log \log d}}$, $\tilde{b} = \sqrt{\rho_{\bar{2}r} \tau d \log \log d} \sigma_w \sigma_B$ and d is large enough. For $D_{\text{mixture-xor}}$ setting, we have $(D_j, +1) \in S_{p_j; B_G}$ for all $j \geq [4]$, where*

$$p = \left(\frac{1}{\sqrt{\rho_{\bar{2}r} \tau d \log \log d} \sigma_B (\log d)^{\frac{18}{25} \frac{2}{B}}}, \gamma = \frac{\sigma_B}{d} \right), \tag{322}$$

$$B_G = \frac{\rho_{\bar{d}}}{4} \left(1 + O\left(\frac{1}{(\log d)^{50}}\right) \right) \sigma_B O\left(\frac{1}{(\log d)^{0.018}}\right). \tag{323}$$

Proof of Lemma E.13. For all $j \geq [r]$, by Lemma E.16, for all $i \in S_{D_j} : \text{Sure}$,

$$1 \frac{\left| \langle G(\mathbf{w}_i^{(0)}, \mathbf{b}_i), D_j \rangle \right|}{kG(\mathbf{w}_i^{(0)}, \mathbf{b}_i)k_2} \quad (324)$$

$$1 \frac{\left| \langle G(\mathbf{w}_i^{(0)}, \mathbf{b}_i), D_j \rangle \right|}{\sqrt{\left| \langle G(\mathbf{w}_i^{(0)}, \mathbf{b}_i), D_j \rangle \right|^2 + \max_{D_j^\perp : \|D_j^\perp\|_2=1} \left| \langle G(\mathbf{w}_i^{(0)}, \mathbf{b}_i), D_j^\perp \rangle \right|^2}} \quad (325)$$

$$1 \frac{\left| \langle G(\mathbf{w}_i^{(0)}, \mathbf{b}_i), D_j \rangle \right|}{\left| \langle G(\mathbf{w}_i^{(0)}, \mathbf{b}_i), D_j \rangle \right| + \max_{D_j^\perp : \|D_j^\perp\|_2=1} \left| \langle G(\mathbf{w}_i^{(0)}, \mathbf{b}_i), D_j^\perp \rangle \right|} \quad (326)$$

$$1 \frac{1}{1 + \frac{B O\left(\frac{1}{(\log d)^{0.018}}\right)}{\frac{1}{4} \sqrt{d} \left(1 - O\left(\frac{1}{(\log d)^{50}}\right)\right) - B O\left(\frac{1}{(\log d)^{0.018}}\right)}} \quad (327)$$

$$\frac{\sigma_B O\left(\frac{1}{(\log d)^{0.018}}\right)}{\frac{1}{4} \sqrt{d} \left(1 - O\left(\frac{1}{(\log d)^{50}}\right)\right) - B O\left(\frac{1}{(\log d)^{0.018}}\right)} \quad (328)$$

$$< \frac{\sigma_B}{d} = \gamma. \quad (329)$$

Thus, we have $G(\mathbf{w}_i^{(0)}, \mathbf{b}_i) \geq C_{D_j}$ and $\left| \langle G(\mathbf{w}_i^{(0)}, \mathbf{b}_i), D_j \rangle \right| \geq kG(\mathbf{w}_i^{(0)}, \mathbf{b}_i)k_2 \geq B_{x1} \frac{b_i}{|\mathbf{b}_i|} = +1$. Thus, by Lemma E.14, we have

$$\Pr_{\mathbf{w}, b} \left[G(\mathbf{w}, b) \geq C_{D_j}; \text{ and } kG(\mathbf{w}, b)k_2 \geq B_G \text{ and } \frac{b}{|b|} = +1 \right] \quad (330)$$

$$\Pr [i \in S_{D_j} : \text{Sure}] \quad (331)$$

$$p. \quad (332)$$

Thus, $(D_j + 1) \geq S_{p_i} : B_G$. We finish the proof. \square

Lemma E.14 (Mixture of Gaussians in [99]: Geometry at Initialization). *Assume the same conditions as in Lemma E.13. Recall for all $i \in [m]$, $\mathbf{w}_i^{(0)} \sim \mathcal{N}(0, \sigma_w^2 I_{d \times d})$, over the random initialization, we have for all $i \in [m]$, $j \in [4]$,*

$$\Pr [i \in S_{D_j} : \text{Sure}] \geq \left(\frac{1}{\frac{\rho}{\tau} \log \log d \sigma_B \left(\log d\right)^{\frac{18}{25} \frac{2}{B}}} \right). \quad (333)$$

Proof of Lemma E.14. WLOG, let $j = 1$. By Assumption E.12, for the first condition in Definition E.6, we have,

$$\Pr \left[\left| \langle \mathbf{w}_i^{(0)}, \mu_1 \rangle \right| \geq C_{\text{Sure},1} b_i \right] \geq \left(\frac{1}{\frac{\rho}{\tau} \log \log d \sigma_B \left(\log d\right)^{\frac{18}{25} \frac{2}{B}}} \right), \quad (334)$$

where the last inequality follows Lemma F.6.

For the second condition in Definition E.6, by Lemma F.6, we have,

$$\Pr \left[\left| \langle \mathbf{w}_i^{(0)}, \mu_2 \rangle \right| \geq C_{\text{Sure},2} b_i \right] \geq \frac{1}{2} \frac{1}{\pi \sigma_B} \frac{1}{e^{\frac{2}{B}}}, \quad (335)$$

On the other hand, if X is a $\chi^2(k)$ random variable. Then we have

$$\Pr(X \leq k + 2 \frac{\rho}{kx + 2x}) \leq e^{-x}. \quad (336)$$

Therefore, we have

$$\Pr \left(\frac{1}{\sigma_w^2} \left\| \mathbf{w}_i^{(0)} \right\|_2^2 \geq d + 2\sqrt{\left(\frac{18\tau\sigma_B^2}{25} + 2 \right) d \log \log d} + 2 \left(\frac{18\tau\sigma_B^2}{25} + 2 \right) \log \log d \right) \quad (337)$$

$$O \left(\frac{1}{(\log d)^2 (\log d)^{\frac{18}{25} \frac{2}{B}}} \right). \quad (338)$$

Thus, by union bound, we have

$$\Pr [i \in S_{D_j; \text{Sure}}] \leq \left(\frac{1}{\frac{\rho}{\tau \log \log d} \sigma_B (\log d)^{\frac{18}{25} \frac{2}{B}}} \right). \quad (339)$$

□

Lemma E.15 (Mixture of Gaussians in [99]: Activation Pattern). *Assume the same conditions as in Lemma E.13, for all $j \geq 4, i \in S_{D_j; \text{Sure}}$, we have*

(1) When $\mathbf{x} \sim \mathcal{N}_j(\mu_j, \sigma_B I_{d \times d})$, the activation probability satisfies,

$$\Pr_{\mathbf{x} \sim \mathcal{N}_j(\mu_j; \sigma_B I_{d \times d})} \left[\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \geq \mathbf{b}_i \right] \geq 1 - \frac{1}{(\log d)^{\frac{6}{5}}}. \quad (340)$$

(2) For all $j' \neq j, j' \geq 4$, when $\mathbf{x} \sim \mathcal{N}_{j'}(\mu_{j'}, \sigma_B I_{d \times d})$, the activation probability satisfies,

$$\Pr_{\mathbf{x} \sim \mathcal{N}_{j'}(\mu_{j'}; \sigma_B I_{d \times d})} \left[\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \geq \mathbf{b}_i \right] = O \left(\frac{1}{(\log d)^{\frac{6}{5}}} \right). \quad (341)$$

Proof of Lemma E.15. In the proof, we need $\tilde{b} = \frac{\rho}{\tau d \log \log d} \sigma_w \sigma_B$. For the first statement, when $\mathbf{x} \sim \mathcal{N}_j(\mu_j, \sigma_B I_{d \times d})$, by $C_{\text{Sure};1} = \frac{6}{5}$, we have

$$\Pr_{\mathbf{x} \sim \mathcal{N}_j(\mu_j; \sigma_B I_{d \times d})} \left[\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \geq \mathbf{b}_i \right] \geq \Pr_{\mathbf{x} \sim \mathcal{N}(0; \sigma_B I_{d \times d})} \left[\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \geq (1 - C_{\text{Sure};1}) \mathbf{b}_i \right] \quad (342)$$

$$\geq \Pr_{\mathbf{x} \sim \mathcal{N}(0; \sigma_B I_{d \times d})} \left[\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \geq \frac{\mathbf{b}_i}{5} \right] \quad (343)$$

$$= 1 - \Pr_{\mathbf{x} \sim \mathcal{N}(0; \sigma_B I_{d \times d})} \left[\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \leq \frac{\mathbf{b}_i}{5} \right] \quad (344)$$

$$= 1 - \exp \left(- \frac{\mathbf{b}_i^2}{50 d \sigma_w^2 \sigma_B^2} \right) \quad (345)$$

$$\geq 1 - \frac{1}{(\log d)^{\frac{6}{5}}}, \quad (346)$$

where the third inequality follows the Chernoff bound and symmetricity of the Gaussian vector.

For the second statement, we prove similarly by $0 < C_{\text{Sure};2} = \frac{\sqrt{2}}{\sqrt{\log \log d}}$. □

Then, Lemma E.16 gives gradients of neurons in $S_{D_j; \text{Sure}}$. It shows that these gradients are highly aligned with D_j .

Lemma E.16 (Mixture of Gaussians in [99]: Feature Emergence). *Assume the same conditions as in Lemma E.13, for all $j \geq 4, i \in S_{D_j; \text{Sure}}$, we have*

$$\left\langle \mathbb{E}_{(\mathbf{x}, y)} \left[y \sigma' \left(\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle - \mathbf{b}_i \right) \mathbf{x} \right], y_{(j)} D_j \right\rangle \quad (347)$$

$$\geq \frac{1}{4} \frac{\rho}{d} \left(1 - O \left(\frac{1}{(\log d)^{\frac{6}{5}}} \right) \right) \sigma_B O \left(\frac{1}{(\log d)^{0.018}} \right). \quad (348)$$

For any unit vector D_j^\perp which is orthogonal with D_j , we have

$$\left| \left\langle \mathbb{E}_{(\mathbf{x}, y)} \left[y \sigma' \left(\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle - \mathbf{b}_i \right) \mathbf{x} \right], D_j^\perp \right\rangle \right| \leq \sigma_B O \left(\frac{1}{(\log d)^{0.018}} \right). \quad (349)$$

Proof of Lemma E.16. For all $j \geq 2$, $i \geq 2$, S_{D_j} , *Sure*, we have

$$\mathbb{E}_{(\mathbf{x}; y)} \left[y \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \mathbf{x} \right] \quad (350)$$

$$= \sum_{l \in [4]} \frac{1}{4} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}_l(\mathbf{x})} \left[y \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \mathbf{x} \right] \quad (351)$$

$$= \sum_{l \in [4]} \frac{1}{4} y_{(l)} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \mathbf{I}_{l,d-d})} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_l \right\rangle - \mathbf{b}_i \right) (\mathbf{x} + \mu_l) \right]. \quad (352)$$

Thus, by Lemma F.7 and Lemma E.15,

$$\left\langle \mathbb{E}_{(\mathbf{x}; y)} \left[y \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \mathbf{x} \right], y_{(j)} D_j \right\rangle \quad (353)$$

$$= \frac{1}{4} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \mathbf{B} \mathbf{I}_{d-d})} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_j \right\rangle - \mathbf{b}_i \right) (\mathbf{x} + \mu_j)^\top D_j \right] \quad (354)$$

$$+ \sum_{l \in [4]; l \neq j} \frac{1}{4} y_{(l)} y_{(j)} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \mathbf{I}_{l,d-d})} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_l \right\rangle - \mathbf{b}_i \right) (\mathbf{x} + \mu_l)^\top D_j \right] \quad (355)$$

$$\frac{1}{4} \mu_j^\top D_j \left(1 - O \left(\frac{1}{(\log d)^{50}} \right) \right) \sum_{l \in [4]; l \neq j} \frac{1}{4} \mu_l^\top D_j O \left(\frac{1}{d^2} \right) \quad (356)$$

$$\frac{1}{4} \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \mathbf{B} \mathbf{I})} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_j \right\rangle - \mathbf{b}_i \right) \mathbf{x}^\top D_j \right] \right| \quad (357)$$

$$\sum_{l \in [4]; l \neq j} \frac{1}{4} \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \mathbf{I}_l)} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_l \right\rangle - \mathbf{b}_i \right) \mathbf{x}^\top D_j \right] \right| \quad (358)$$

$$\frac{1}{4} \rho_{\bar{d}} \left(1 - O \left(\frac{1}{(\log d)^{50}} \right) \right) \frac{1}{4} \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \mathbf{B} \mathbf{I})} \left[\left(1 - \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_j \right\rangle - \mathbf{b}_i \right) \right) \mathbf{x}^\top D_j \right] \right| \quad (359)$$

$$\sum_{l \in [4]; l \neq j} \frac{1}{4} \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \mathbf{I}_l)} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_l \right\rangle - \mathbf{b}_i \right) \mathbf{x}^\top D_j \right] \right| \quad (359)$$

$$= \frac{1}{4} \rho_{\bar{d}} \left(1 - O \left(\frac{1}{(\log d)^{50}} \right) \right) \frac{1}{4} \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \mathbf{B} \mathbf{I})} \left[\left(1 - \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_j \right\rangle - \mathbf{b}_i \right) \right) \mathbf{x}^\top D_j \right] \right| \quad (360)$$

$$\sum_{l \in [4]; l \neq j} \frac{1}{4} \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \mathbf{I}_l)} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_l \right\rangle - \mathbf{b}_i \right) \mathbf{x}^\top D_j \right] \right| \quad (360)$$

$$\frac{1}{4} \rho_{\bar{d}} \left(1 - O \left(\frac{1}{(\log d)^{50}} \right) \right) \sigma_B O \left(\frac{1}{(\log d)^{0.018}} \right). \quad (361)$$

For any unit vector D_j^\perp which is orthogonal with D_j , similarly, we have

$$\left| \left\langle \mathbb{E}_{(\mathbf{x}; y)} \left[y \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \mathbf{x} \right], D_j^\perp \right\rangle \right| \quad (362)$$

$$\frac{1}{4} \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \mathbf{B} \mathbf{I})} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_j \right\rangle - \mathbf{b}_i \right) \mathbf{x}^\top D_j^\perp \right] \right| \quad (363)$$

$$+ \sum_{l \in [4]; l \neq j} \frac{1}{4} \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \mathbf{I}_l)} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_l \right\rangle - \mathbf{b}_i \right) (\mathbf{x} + \mu_l)^\top D_j^\perp \right] \right| \quad (364)$$

$$\frac{1}{4} \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \mathbf{B} \mathbf{I})} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_j \right\rangle - \mathbf{b}_i \right) \mathbf{x}^\top D_j^\perp \right] \right| \quad (365)$$

$$+ \sum_{l \in [4]; l \neq j} \frac{1}{4} \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0; \mathbf{I}_l)} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} + \mu_l \right\rangle - \mathbf{b}_i \right) \mathbf{x}^\top D_j^\perp \right] \right| \quad (366)$$

$$\sigma_B O \left(\frac{1}{(\log d)^{0.018}} \right). \quad (367)$$

□

E.3.2 Mixture of Gaussians - XOR: Final Guarantee

Lemma E.17 (Mixture of Gaussians in [99]: Existence of Good Networks). *Assume the same conditions as in Lemma E.13 and let $\tau = 1$ and when $0 < \tau = O\left(\frac{d}{B \log d}\right)$. Define*

$$f^*(\mathbf{x}) = \sum_{j=1}^4 \frac{y^{(j)}}{\tau \log d \sigma_B} \left[\sigma \left(h D_j, \mathbf{x} \mid 2\sqrt{\tau \log d \sigma_B} \right) \right]. \quad (368)$$

For $D_{\text{mixture-xor}}$ setting, we have $f^* \in \mathcal{F}_{d;r;B_F;S_{p_i};B_G}$, where $B_F = (B_{a1}, B_{a2}, B_b) = \left(\frac{1}{\sqrt{-\log d} \sigma_B}, \frac{2}{\sqrt{-\log d} \sigma_B}, 2^{\rho} \frac{\rho}{\tau \log d \sigma_B}\right)$, $p = \left(\frac{1}{B \cdot (\log d)^{\frac{2}{B}}}\right)$, $\gamma = \frac{B}{\sqrt{d}}$, $r = 4$, $B_G = \frac{1}{5} \rho \bar{d}$ and $B_{x1} = (1 + \sigma_B) \rho \bar{d}$, $B_{x2} = (1 + \sigma_B)^2 d$. We also have $\text{OPT}_{d;r;B_F;S_{p_i};B_G} \leq \frac{3}{d} + \frac{4}{d^{0.9-1} \sqrt{-\log d}}$.

Proof of Lemma E.17. We finish the proof by following the proof of Lemma E.11 \square

Theorem E.18 (Mixture of Gaussians in [99]: Main Result). *For $D_{\text{mixture-xor}}$ setting with Assumption E.12, when d is large enough, for any $\delta \in (0, 1)$ and for any $\epsilon \in (0, 1)$ when*

$$m = \left(\sigma_B (\log d)^{\frac{2}{B}} \left(\left(\log \left(\frac{1}{\delta} \right) \right)^2 + \frac{1 + \sigma_B}{\epsilon^4} \right) + \frac{1}{\delta} \right) e^d, \quad (369)$$

$$T = \text{poly}(\sigma_B, 1/\epsilon, 1/\delta, \log d), \quad (370)$$

$$n = \left\lceil \left(\frac{m^3 (1 + \sigma_B^2)}{\epsilon^2 \max\{\sigma_B (\log d)^{\frac{2}{B}}, 1\}} + \sigma_B (\log d)^{\frac{2}{B}} + \frac{Tm}{\delta} \right) \right\rceil, \quad (371)$$

trained by Algorithm 1 with hinge loss, with probability at least $1 - \delta$ over the initialization and training samples, with proper hyper-parameters, there exists $t \geq [T]$ such that

$$\Pr[\text{sign}(f^{(t)}(\mathbf{x})) \neq y] \leq O\left(\left(1 + \sigma_B^{\frac{3}{B}}\right) \left(\frac{1}{d^{\frac{1}{4}}} + \frac{(\log n)^{\frac{1}{4}}}{n^{\frac{1}{4}}}\right)\right) + \epsilon. \quad (372)$$

Proof of Theorem E.18. Let $\bar{b} = \rho \bar{d} \log \log d \sigma_w \sigma_B$. By Lemma E.17, let $\tau = 1$ and when $\tau = O\left(\frac{d}{B \log d}\right)$, we have $f^* \in \mathcal{F}_{d;r;B_F;S_{p_i};B_G}$, where $B_F = (B_{a1}, B_{a2}, B_b) = \left(\frac{1}{\sqrt{-\log d} \sigma_B}, \frac{2}{\sqrt{-\log d} \sigma_B}, 2^{\rho} \frac{\rho}{\tau \log d \sigma_B}\right)$, $p = \left(\frac{1}{B \cdot (\log d)^{\frac{2}{B}}}\right)$, $\gamma = \frac{B}{\sqrt{d}}$, $r = 4$, $B_G = \frac{1}{5} \rho \bar{d}$ and $B_{x1} = (1 + \sigma_B) \rho \bar{d}$, $B_{x2} = (1 + \sigma_B)^2 d$. We also have $\text{OPT}_{d;r;B_F;S_{p_i};B_G} \leq \frac{3}{d} + \frac{4}{d^{0.9-1} \sqrt{-\log d}}$.

Adjust σ_w such that $\bar{b} = \rho \bar{d} \log \log d \sigma_w \sigma_B = \left(\frac{B_G^{\frac{1}{4}} B_{a2} B_b^{\frac{3}{4}}}{\sqrt{F} B_{a1}}\right)$. Injecting above parameters into Theorem 3.12, we have with probability at least $1 - \delta$ over the initialization, with proper hyper-parameters, there exists $t \geq [T]$ such that

$$\Pr[\text{sign}(f^{(t)}(\mathbf{x})) \neq y] \leq O\left(\left(1 + \sigma_B^{\frac{3}{B}}\right) \left(\frac{1}{d^{\frac{1}{4}}} + \frac{(\log n)^{\frac{1}{4}}}{n^{\frac{1}{4}}}\right)\right) + \epsilon. \quad (373)$$

\square

E.4 Parity Functions

We recap the problem setup in Section 4.2 for readers' convenience.

E.4.1 Problem Setup

Data Distributions. Suppose $\mathbf{M} \in \mathbb{R}^{d \times D}$ is an unknown dictionary with D columns that can be regarded as patterns. For simplicity, assume $d = D$ and \mathbf{M} is orthonormal. Let $\phi \in \mathbb{R}^d$ be a hidden representation vector. Let $A \subseteq [D]$ be a subset of size rk corresponding to the class relevant patterns

and r is an odd number. Then the input is generated by $\mathbf{M}\phi$, and some function on ϕ_A generates the label. WLOG, let $A = \{1, \dots, rk\}$, $A^\perp = \{rk+1, \dots, d\}$. Also, we split A such that for all $j \geq [r]$, $A_j = \{(j-1)k+1, \dots, jk\}$. Then the input \mathbf{x} and the class label y are given by:

$$\mathbf{x} = \mathbf{M}\phi, \quad y = g^*(\phi_A) = \text{sign} \left(\sum_{j=1}^r \text{XOR}(\phi_{A_j}) \right), \quad (374)$$

where g^* is the ground-truth labeling function mapping from \mathbb{R}^{rk} to $Y = \{-1, 1\}$, ϕ_A is the sub-vector of ϕ with indices in A , and $\text{XOR}(\phi_{A_j}) = \prod_{l \in A_j} \phi_l$ is the parity function.

We still need to specify the distribution of ϕ , which determines the structure of the input distribution:

$$X := (1 - 2rp_A)X_U + \sum_{j \in [r]} p_A(X_{j,+} + X_{j,-}). \quad (375)$$

For all corresponding ϕ_{A_j} in X , we have $\phi_l \in \{-1, 1\}$, independently:

$$\phi_l = \begin{cases} +1, & \text{w.p. } p_o \\ 1, & \text{w.p. } p_o \\ 0, & \text{w.p. } 1 - 2p_o \end{cases}$$

where p_o controls the signal noise ratio: if p_o is large, then there are many nonzero entries in A^\perp which are noise interfering with the learning of the ground-truth labeling function on A .

For corresponding ϕ_A , any $j \geq [r]$, we have

- In $X_{j,+}$, $\phi_{A_j} = [+1, +1, \dots, +1]^\top$ and $\phi_{A \setminus A_j}$ only have zero elements.
- In $X_{j,-}$, $\phi_{A_j} = [-1, -1, \dots, -1]^\top$ and $\phi_{A \setminus A_j}$ only have zero elements.
- In X_U , we have ϕ_A draw from $\{-1, 1\}^{rk}$ uniformly.

We call this data distribution D_{parity} .

Assumption E.19 (Parity Functions. Recap of Assumption 4.5). *Let $\delta = \tau \log d$ be a parameter that will control our final error guarantee. Assume k is an odd number and:*

$$k \geq (\tau \log d), \quad d \geq rk + (\tau r \log d), \quad p_o = O\left(\frac{rk}{d - rk}\right), \quad p_A = \frac{1}{d}. \quad (376)$$

Remark E.20. *The assumptions require k, d , and p_A to be sufficiently large so as to provide enough large signals for learning. When $p_o = \frac{rk}{d - rk}$ means that the signal noise ratio is constant: the expected norm of ϕ_A and that of ϕ_{A^\perp} are comparable.*

To apply our framework, again we only need to compute the parameters in the Gradient Feature set and the corresponding optimal approximation loss. To this end, we first define the gradient features: For all $j \geq [r]$, let

$$D_j = \frac{\sum_{l \in A_j} \mathbf{M}_l}{k \sum_{l \in A_j} \mathbf{M}_l^2}. \quad (377)$$

Remark E.21. *Our data distribution is symmetric, which means for any $\phi \in \mathbb{R}^d$:*

- $y = g^*(\phi_A)$ and $x = \mathbf{M}(\phi)$,
- $\mathbb{P}(\phi) = \mathbb{P}(-\phi)$,
- $\mathbb{E}_{(\mathbf{x}, y)}[y\mathbf{x}] = \mathbf{0}$.

Below, we define a sufficient condition that randomly initialized weights will fall in nice gradients set after the first gradient step update.

Definition E.22 (Parity Functions: Subset of Nice Gradients Set). *Recall $\mathbf{w}_i^{(0)}$ is the weight for the i -th neuron at initialization. For all $j \geq [r]$, let $S_{D_j; \text{Sure}} \subseteq [m]$ be those neurons that satisfy*

$$\bullet \left\langle \mathbf{w}_i^{(0)}, D_j \right\rangle \geq \frac{C_{\text{Sure}; 1}}{\sqrt{k}} \mathbf{b}_i,$$

- $\left| \left\langle \mathbf{w}_i^{(0)}, D_j \right\rangle \right| \leq \frac{C_{\text{Sure},2}}{\sqrt{k}} \mathbf{b}_i$, for all $j' \notin j, j' \geq [r]$,
- $\left\| P_A \mathbf{w}_i^{(0)} \right\|_2 \leq \left(\frac{\rho_{\bar{k}}}{r k \sigma_w} \right)$,
- $\left\| P_{A^c} \mathbf{w}_i^{(0)} \right\|_2 \leq \left(\frac{\rho_{\bar{k}}}{d - r k \sigma_w} \right)$,

where P_A, P_{A^c} are the projection operator on the space \mathbf{M}_A and \mathbf{M}_{A^c} .

E.4.2 Parity Functions: Feature Learning

We show the important Lemma E.23 first and defer other Lemmas after it.

Lemma E.23 (Parity Functions: Gradient Feature Set. Part statement of Lemma 4.7). *Let $C_{\text{Sure},1} = \frac{3}{2}, C_{\text{Sure},2} = \frac{1}{2}, \mathbf{b} = C_b \frac{\rho_{\bar{k}}}{\tau r k \log d \sigma_w}$, where C_b is a large enough universal constant. For D_{parity} setting, we have $(D_j, +1), (D_j, -1) \geq S_{p_i; B_G}$ for all $j \geq [r]$, where*

$$p = \left(\frac{1}{\rho_{\bar{k}} \tau r \log d} \frac{1}{d^{(9C_b^2 r - 8)}} \right), \quad \gamma = \frac{1}{d - 2}, \quad B_G = \frac{\rho_{\bar{k}}}{k p_A} \quad O\left(\frac{\rho_{\bar{k}}}{d}\right). \quad (378)$$

Proof of Lemma E.23. Note that for all $l \geq [d]$, we have $\mathbf{M}_l^\top \mathbf{x} = \phi_l$. For all $j \geq [r]$, by Lemma E.26, for all $i \geq S_{D_j; \text{Sure}}$, when $\gamma = \frac{1}{d - 2}$,

$$\left| \left\langle G(\mathbf{w}_i^{(0)}, \mathbf{b}_i), D_j \right\rangle \right| \leq (1 - \gamma) k G(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \quad (379)$$

$$= \left| \left\langle G(\mathbf{w}_i^{(0)}, \mathbf{b}_i), \frac{\sum_{l \in A_j} \mathbf{M}_l}{\rho_{\bar{k}}} \right\rangle \right| \leq (1 - \gamma) k G(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \quad (380)$$

$$\leq \frac{\rho_{\bar{k}}}{k p_A} \quad O\left(\frac{\rho_{\bar{k}}}{d}\right) \quad \left(1 - \frac{1}{d - 2}\right) \sqrt{k p_A^2 + \sum_{l \in [d]} O\left(\frac{1}{d}\right)^2} \quad (381)$$

$$\leq \frac{\rho_{\bar{k}}}{k p_A} \quad O\left(\frac{\rho_{\bar{k}}}{d}\right) \quad \left(1 - \frac{1}{d - 2}\right) \left(\frac{\rho_{\bar{k}}}{k p_A} + O\left(\frac{1}{d - \frac{1}{2}}\right)\right) \quad (382)$$

$$\leq \frac{\rho_{\bar{k}}}{k p_A} \quad O\left(\frac{\rho_{\bar{k}}}{d}\right) \quad O\left(\frac{1}{d - \frac{1}{2}}\right) \quad (383)$$

$$> 0. \quad (384)$$

Thus, we have $G(\mathbf{w}_i^{(0)}, \mathbf{b}_i) \geq C_{D_j}$; and $\frac{\rho_{\bar{k}}}{k p_A} \quad O\left(\frac{\sqrt{k}}{d}\right) \leq k G(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \leq \frac{\rho_{\bar{k}}}{k p_A} + O\left(\frac{1}{d - \frac{1}{2}}\right)$, $\frac{\mathbf{b}_i}{|\mathbf{b}_i|} = +1$. Thus, by Lemma E.24, we have

$$\Pr_{\mathbf{w}, \mathbf{b}} \left[G(\mathbf{w}, \mathbf{b}) \geq C_{D_j}; \text{ and } k G(\mathbf{w}, \mathbf{b}) k_2 \leq B_G \text{ and } \frac{\mathbf{b}}{|\mathbf{b}|} = +1 \right] \quad (385)$$

$$\Pr [i \geq S_{D_j; \text{Sure}}] \quad (386)$$

$$p. \quad (387)$$

Thus, $(D_j, +1) \geq S_{p_i; B_G}$. Since $\mathbb{E}_{(\mathbf{x}; \mathbf{y})}[\mathbf{y} \mathbf{x}] = \mathbf{0}$, by Lemma F.2 and considering $i \geq [2m] \cap [m]$, we have $(D_j, -1) \geq S_{p_i; B_G}$. We finish the proof. \square

Below are Lemmas used in the proof of Lemma E.23. In Lemma E.24, we calculate p used in $S_{p_i; B_G}$.

Lemma E.24 (Parity Functions: Geometry at Initialization. Lemma B.2 in [7]). *Assume the same conditions as in Lemma E.23, recall for all $i \geq [m]$, $\mathbf{w}_i^{(0)} \sim N(0, \sigma_w^2 I_{d \times d})$, over the random initialization, we have for all $i \geq [m], j \geq [r]$,*

$$\Pr [i \geq S_{D_j; \text{Sure}}] \leq \left(\frac{1}{\rho_{\bar{k}} \tau r \log d} \frac{1}{d^{(9C_b^2 r - 8)}} \right). \quad (388)$$

Proof of Lemma E.24. For every $i \in [m], j, j' \in [r], j \neq j'$, by Lemma F.6,

$$p_1 = \Pr \left[\left\langle \mathbf{w}_i^{(0)}, D_j \right\rangle \frac{C_{\text{Sure},1}}{\rho \frac{r}{k}} \mathbf{b}_i \right] = \left(\frac{1}{\rho \frac{r}{\tau r \log d} d^{(9C_b^2 r-8)}} \right) \quad (389)$$

$$p_2 = \Pr \left[\left| \left\langle \mathbf{w}_i^{(0)}, D_{j'} \right\rangle \right| \frac{C_{\text{Sure},2}}{\rho \frac{r}{k}} \mathbf{b}_i \right] = \left(\frac{1}{\rho \frac{r}{\tau r \log d} d^{(C_b^2 r-8)}} \right). \quad (390)$$

On the other hand, if X is a $\chi^2(k)$ random variable, by Lemma F.5, we have

$$\Pr(X \leq k + 2\sqrt{\rho \frac{r}{k} x} + 2x) \leq e^{-x}. \quad (391)$$

Therefore, by assumption $rk \geq (\tau r \log d), d \geq rk \geq (\tau r \log d)$, we have

$$\Pr \left(\frac{1}{\sigma_w^2} \left\| P_A \mathbf{w}_i^{(0)} \right\|_2^2 \leq rk + 2\sqrt{(9C_b^2 \tau r / 8 + 2)rk \log d} + 2(9C_b^2 \tau r / 8 + 2) \log d \right) \quad (392)$$

$$O \left(\frac{1}{d^2 d^{(9C_b^2 r-8)}} \right), \quad (393)$$

$$\Pr \left(\frac{1}{\sigma_w^2} \left\| P_A \mathbf{w}_i^{(0)} \right\|_2^2 \leq (d - rk) + 2\sqrt{(9C_b^2 \tau r / 8 + 2)(d - rk) \log d} + 2(9C_b^2 \tau r / 8 + 2) \log d \right)$$

$$O \left(\frac{1}{d^2 d^{(9C_b^2 r-8)}} \right). \quad (394)$$

Thus, by union bound, and D_1, \dots, D_r being orthogonal with each other, we have

$$\Pr [i \in S_{D_j, \text{Sure}}] \leq p_1 (1 - p_2)^{r-1} = O \left(\frac{1}{d^2 d^{(9C_b^2 r-8)}} \right) \quad (395)$$

$$= \left(\frac{1}{\rho \frac{r}{\tau r \log d} d^{(9C_b^2 r-8)}} \left(1 - \frac{r}{\rho \frac{r}{\tau r \log d} d^{(C_b^2 r-8)}} \right) \right) \quad (396)$$

$$O \left(\frac{1}{d^2 d^{(9C_b^2 r-8)}} \right) \quad (397)$$

$$= \left(\frac{1}{\rho \frac{r}{\tau r \log d} d^{(9C_b^2 r-8)}} \right). \quad (398)$$

□

In Lemma E.25, we compute the activation pattern for the neurons in $S_{D_j, \text{Sure}}$.

Lemma E.25 (Parity Functions: Activation Pattern). *Assume the same conditions as in Lemma E.23, for all $j \in [r], i \in S_{D_j, \text{Sure}}$, we have*

(1) When $\mathbf{x} \sim \mathcal{X}$, we have

$$\Pr_{\mathbf{x} \sim \mathcal{X}} \left[\left| \sum_{l \in A^?} \left\langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \phi_l \right\rangle \right| \geq t \right] \leq \exp \left(- \frac{t^2}{(rk\sigma_w^2)} \right). \quad (399)$$

(2) When $\mathbf{x} \sim \mathcal{X}_U$, we have

$$\Pr_{\mathbf{x} \sim \mathcal{X}_U} \left[\left| \sum_{l \in A} \left\langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \phi_l \right\rangle \right| \geq t \right] \leq \exp \left(- \frac{t^2}{(rk\sigma_w^2)} \right). \quad (400)$$

(3) When $\mathbf{x} \sim \mathcal{X}_U$, the activation probability satisfies,

$$\Pr_{\mathbf{x} \sim \mathcal{X}_U} \left[\sum_{l \in [d]} \left\langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \phi_l \right\rangle \mathbf{b}_l = 0 \right] = O \left(\frac{1}{d} \right). \quad (401)$$

(4) When $\mathbf{x} \sim \mathcal{X}_{j;+}$, the activation probability satisfies,

$$\Pr_{\mathbf{x} \sim \mathcal{X}_{j;+}} \left[\sum_{l \in [d]} \langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \phi_l \rangle \cdot \mathbf{b}_l \geq 0 \right] = 1 - O\left(\frac{1}{d}\right). \quad (402)$$

(5) For all $j' \notin j, j' \geq [r], s \geq f+, g$, when $\mathbf{x} \sim \mathcal{X}_{j^0;s}$, or $\mathbf{x} \sim \mathcal{X}_{j;-}$, the activation probability satisfies,

$$\Pr \left[\sum_{l \in [d]} \langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \phi_l \rangle \cdot \mathbf{b}_l \geq 0 \right] = O\left(\frac{1}{d}\right). \quad (403)$$

Proof of Lemma E.25. For the first statement, when $\mathbf{x} \sim \mathcal{X}$, note that $\langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \rangle \phi_l$ is a mean-zero sub-Gaussian random variable with sub-Gaussian norm $\left(\left| \langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \rangle \right| \sqrt{p_0} \right)$.

$$\Pr_{\mathbf{x} \sim \mathcal{X}} \left[\left| \sum_{l \in A^?} \langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \phi_l \rangle \right| \geq t \right] = \Pr_{\mathbf{x} \sim \mathcal{X}} \left[\left| \sum_{l \in A^?} \langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \rangle \phi_l \right| \geq t \right] \quad (404)$$

$$\exp \left(- \frac{t^2}{\sum_{l \in A^?} \left(\langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \rangle^2 p_0 \right)} \right) \quad (405)$$

$$\exp \left(- \frac{t^2}{((d-rk)\sigma_w^2 p_0)} \right) \quad (406)$$

$$\exp \left(- \frac{t^2}{(rk\sigma_w^2)} \right), \quad (407)$$

where the inequality follows general Hoeffding's inequality.

For the second statement, when $\mathbf{x} \sim \mathcal{X}_U$, by Hoeffding's inequality,

$$\Pr_{\mathbf{x} \sim \mathcal{X}_U} \left[\left| \sum_{l \in A} \langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \phi_l \rangle \right| \geq t \right] = \Pr_{\mathbf{x} \sim \mathcal{X}_U} \left[\left| \sum_{l \in A} \langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \rangle \phi_l \right| \geq t \right] \quad (408)$$

$$2 \exp \left(- \frac{t^2}{2 \sum_{l \in A} \langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \rangle^2} \right) \quad (409)$$

$$\exp \left(- \frac{t^2}{(rk\sigma_w^2)} \right). \quad (410)$$

In the proof of the third to the last statement, we need $\tilde{b} = C_b \frac{D}{\tau rk \log d} \sigma_w$, where C_b is a large enough universal constant.

For the third statement, when $\mathbf{x} \sim \mathcal{X}_U$, by union bound and previous statements,

$$\Pr_{\mathbf{x} \sim \mathcal{X}_U} \left[\sum_{l \in [d]} \langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \phi_l \rangle \cdot \mathbf{b}_l \geq 0 \right] \quad (411)$$

$$\Pr_{\mathbf{x} \sim \mathcal{X}_U} \left[\sum_{l \in A} \langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \phi_l \rangle \cdot \frac{\mathbf{b}_l}{2} \right] + \Pr_{\mathbf{x} \sim \mathcal{X}_U} \left[\sum_{l \in A^?} \langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \phi_l \rangle \cdot \frac{\mathbf{b}_l}{2} \right] \quad (412)$$

$$O\left(\frac{1}{d}\right). \quad (413)$$

For the forth statement, when $\mathbf{x} \sim \mathcal{X}_{j;+}$, by $C_{Sure;1} \leq \frac{3}{2}$ and previous statements,

$$\Pr_{\mathbf{x} \sim \mathcal{X}_{j;+}} \left[\sum_{l \in [d]} \langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \phi_l \rangle \mathbf{b}_i \geq 0 \right] \quad (414)$$

$$= \Pr_{\mathbf{x} \sim \mathcal{X}_{j;+}} \left[\sum_{l \in A_j} \langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \phi_l \rangle + \sum_{l \in A \setminus A_j} \langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \phi_l \rangle + \sum_{l \in A^?} \langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \phi_l \rangle \mathbf{b}_i \right] \quad (415)$$

$$\Pr_{\mathbf{x} \sim \mathcal{X}_{j;+}} \left[\sum_{l \in A^?} \langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \phi_l \rangle (1 - C_{Sure;1}) \mathbf{b}_i \right] \quad (416)$$

$$\Pr_{\mathbf{x} \sim \mathcal{X}_{j;+}} \left[\sum_{l \in A^?} \langle \mathbf{w}_i^{(0)}, \mathbf{M}_l \phi_l \rangle \frac{\mathbf{b}_i}{2} \right] \quad (417)$$

$$1 - O\left(\frac{1}{d}\right). \quad (418)$$

For the last statement, we prove similarly by $0 < C_{Sure;2} \leq \frac{1}{2}$. \square

Then, Lemma E.26 gives gradients of neurons in $S_{D_j; Sure}$. It shows that these gradients are highly aligned with D_j .

Lemma E.26 (Parity Functions: Feature Emergence). *Assume the same conditions as in Lemma E.23, for all $j \geq 2[r]$, $i \geq 2 S_{D_j; Sure}$, we have the following holds:*

(1) For all $l \geq A_j$, we have

$$p_A - O\left(\frac{1}{d}\right) \mathbb{E}_{(\mathbf{x}; y)} \left[y \sigma' \left(\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \mathbf{b}_i \right) \phi_l \right] \geq p_A + O\left(\frac{1}{d}\right). \quad (419)$$

(2) For all $l \geq A_{j^0}$, any $j' \neq j$, $j' \geq 2[r]$, we have

$$\left| \mathbb{E}_{(\mathbf{x}; y)} \left[y \sigma' \left(\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \mathbf{b}_i \right) \phi_l \right] \right| \leq O\left(\frac{1}{d}\right). \quad (420)$$

(3) For all $l \geq A^\perp$, we have

$$\left| \mathbb{E}_{(\mathbf{x}; y)} \left[y \sigma' \left(\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \mathbf{b}_i \right) \phi_l \right] \right| \leq O\left(\frac{1}{d}\right). \quad (421)$$

Proof of Lemma E.26. For all $l \geq [d]$, we have

$$\mathbb{E}_{(\mathbf{x}; y)} \left[y \sigma' \left(\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \mathbf{b}_i \right) \phi_l \right] \quad (422)$$

$$= p_A \sum_{l \in [r]} \left(\mathbb{E}_{\mathbf{x} \sim \mathcal{X}_{j;+}} \left[\sigma' \left(\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \mathbf{b}_i \right) \phi_l \right] - \mathbb{E}_{\mathbf{x} \sim \mathcal{X}_j} \left[\sigma' \left(\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \mathbf{b}_i \right) \phi_l \right] \right) \quad (423)$$

$$+ (1 - 2rp_A) \mathbb{E}_{\mathbf{x} \sim \mathcal{X}_U} \left[y \sigma' \left(\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \mathbf{b}_i \right) \phi_l \right]. \quad (424)$$

For the first statement, for all $l \geq A_j$, by Lemma E.25 (3) and (4), we have

$$\mathbb{E}_{(\mathbf{x}; y)} \left[y \sigma' \left(\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \mathbf{b}_i \right) \phi_l \right] \quad (425)$$

$$= p_A \left(\mathbb{E}_{\mathbf{x} \sim \mathcal{X}_{j;+}} \left[\sigma' \left(\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \mathbf{b}_i \right) \right] + \mathbb{E}_{\mathbf{x} \sim \mathcal{X}_j} \left[\sigma' \left(\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \mathbf{b}_i \right) \right] \right) \quad (426)$$

$$+ (1 - 2rp_A) \mathbb{E}_{\mathbf{x} \sim \mathcal{X}_U} \left[y \sigma' \left(\langle \mathbf{w}_i^{(0)}, \mathbf{x} \rangle \mathbf{b}_i \right) \phi_l \right] \quad (427)$$

$$p_A \left(1 - O\left(\frac{1}{d}\right) \right) \geq O\left(\frac{1}{d}\right) \quad (428)$$

$$p_A - O\left(\frac{1}{d}\right), \quad (429)$$

and we also have

$$\mathbb{E}_{(\mathbf{x};y)} \left[y \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \phi_l \right] \quad (430)$$

$$= p_A \left(\mathbb{E}_{\mathbf{x} \sim \mathcal{X}_{j,+}} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \right] + \mathbb{E}_{\mathbf{x} \sim \mathcal{X}_j} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \right] \right) \quad (431)$$

$$+ (1 - 2rp_A) \mathbb{E}_{\mathbf{x} \sim \mathcal{X}_U} \left[y \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \phi_l \right] \quad (432)$$

$$p_A + O\left(\frac{1}{d}\right). \quad (433)$$

Similarly, for the second statement, for all $l \geq A_{j^0}$, any $j' \notin j, j' \geq [r]$, by Lemma E.25 (3) and (5), we have

$$\left| \mathbb{E}_{(\mathbf{x};y)} \left[y \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \phi_l \right] \right| \quad (434)$$

$$\left| p_A \left(\mathbb{E}_{\mathbf{x} \sim \mathcal{X}_{j^0,+}} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \right] + \mathbb{E}_{\mathbf{x} \sim \mathcal{X}_{j^0}} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \right] \right) \right| + O\left(\frac{1}{d}\right) \\ O\left(\frac{1}{d}\right). \quad (435)$$

For the third statement, for all $l \geq A^\perp$, by Lemma E.25 (3), (4), (5), we have

$$\left| \mathbb{E}_{(\mathbf{x};y)} \left[y \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \phi_l \right] \right| \quad (436)$$

$$p_A \sum_{l \in [r]} \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{X}_{l,+}} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \phi_l \right] - \mathbb{E}_{\mathbf{x} \sim \mathcal{X}_l} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \phi_l \right] \right| + O\left(\frac{1}{d}\right)$$

$$p_A \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{X}_{j,+}} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \phi_l \right] - \mathbb{E}_{\mathbf{x} \sim \mathcal{X}_j} \left[\sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \phi_l \right] \right| + O\left(\frac{1}{d}\right)$$

$$p_A \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{X}_{j,+}} \left[\left(1 - \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \right) \phi_l \right] - \mathbb{E}_{\mathbf{x} \sim \mathcal{X}_j} \left[\left(1 - \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \right) \phi_l \right] \right| \\ + p_A \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{X}_{j,+}} [\phi_l] - \mathbb{E}_{\mathbf{x} \sim \mathcal{X}_j} [\phi_l] \right| + O\left(\frac{1}{d}\right) \quad (437)$$

$$= p_A \left| \mathbb{E}_{\mathbf{x} \sim \mathcal{X}_{j,+}} \left[\left(1 - \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \right) \phi_l \right] - \mathbb{E}_{\mathbf{x} \sim \mathcal{X}_j} \left[\left(1 - \sigma' \left(\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right) \right) \phi_l \right] \right| \\ + O\left(\frac{1}{d}\right) \quad (438)$$

$$O\left(\frac{1}{d}\right), \quad (439)$$

where the second inequality follows $2rp_A \leq 1$ and the third inequality follows the triangle inequality. \square

E.4.3 Parity Functions: Final Guarantee

Lemma E.27 (Parity Functions: Existence of Good Networks. Part statement of Lemma 4.7). *Assume the same conditions as in Lemma E.23. Define*

$$f^*(\mathbf{x}) = \sum_{j=1}^r \sum_{i=0}^k \binom{k}{i} \rho_{\frac{k}{k}}^{i+1} \\ \left[\sigma \left(\langle \mathbf{h} D_j, \mathbf{x} \rangle - \frac{2i}{k} \rho_{\frac{k}{k}} \right) - 2\sigma \left(\langle \mathbf{h} D_j, \mathbf{x} \rangle - \frac{2i}{k} \rho_{\frac{k}{k}} \right) + \sigma \left(\langle \mathbf{h} D_j, \mathbf{x} \rangle - \frac{2i}{k} \rho_{\frac{k}{k}} \right) \right]. \quad (440)$$

For D_{parity} setting, we have $f^* \geq F_{d;3r(k+1);B_F;S_{p^*};B_G}$, where $B_F = (B_{a1}, B_{a2}, B_b) = \left(2\rho_{\frac{k}{k}}, 2\sqrt{rk(k+1)}, \frac{k+1}{\sqrt{k}} \right)$, $p = \left(\frac{1}{\sqrt{r \log d} \cdot d^{(9C_b^2 r - 8)}} \right)$, $\gamma = \frac{1}{d^2}$, $B_G = \rho_{\frac{k}{k}} p_A + O\left(\frac{\sqrt{k}}{d}\right)$ and $B_{x1} = \rho_{\frac{k}{k}}/d$, $B_{x2} = d$. We also have $\text{OPT}_{d;3r(k+1);B_F;S_{p^*};B_G} = 0$.

Proof of Lemma E.27. We can check $B_{x_1} = \frac{\rho_-}{d}$, $B_{x_2} = d$ by direct calculation. By Lemma E.23, we have $f^* \geq F_{d;3r(k+1);B_F;S_{p_1};B_G}$. We note that

$$\sigma \left(hD_j, \mathbf{x} \mid \frac{2i-k}{\rho_-} \right) = 2\sigma \left(hD_j, \mathbf{x} \mid \frac{2i-k}{\rho_-} \right) + \sigma \left(hD_j, \mathbf{x} \mid \frac{2i-k+1}{\rho_-} \right) \quad (441)$$

is a bump function for $hD_j, \mathbf{x} \mid \frac{2i-k}{\rho_-}$. We can check that $y f^*(\mathbf{x}) = 1$. Thus, we have

$$\text{OPT}_{d;3r(k+1);B_F;S_{p_1};B_G} = L_{\mathcal{D}_{\text{parity}}}(f^*) \quad (442)$$

$$= \mathbb{E}_{(\mathbf{x};y) \sim \mathcal{D}_{\text{parity}}} L_{(\mathbf{x};y)}(f^*) \quad (443)$$

$$= 0. \quad (444)$$

□

Theorem 4.8 (Parity Functions: Main Result). *Assume Assumption 4.5. For any $\epsilon, \delta \geq (0, 1)$, when Algorithm 1 uses hinge loss with*

$$m = \text{poly} \left(\frac{1}{\delta}, \frac{1}{\epsilon}, d^{(r)}, k, \frac{1}{p_A} \right) e^d, \quad T = \text{poly}(m), \quad n = \text{poly}(m)$$

and proper hyper-parameters, then with probability at least $1 - \delta$, there exists $t \geq [T]$ such that

$$\Pr[\text{sign}(f^{(t)}(\mathbf{x})) \neq y] \leq \frac{3r \frac{\rho_-}{k}}{d^{(-3)=2}} + \epsilon.$$

Proof of Theorem 4.8. Let $\tilde{b} = C_b \frac{\rho_-}{\tau r k \log d} \sigma_w$, where C_b is a large enough universal constant. By Lemma E.27, we have $f^* \geq F_{d;3r(k+1);B_F;S_{p_1};B_G}$, where $B_F = (B_{a_1}, B_{a_2}, B_b) = \left(2^{\frac{\rho_-}{k}}, 2\sqrt{rk(k+1)}, \frac{k+1}{\sqrt{k}} \right)$, $p = \left(\frac{1}{\sqrt{r \log d} \cdot d^{(9C_b^2 r - 8)}} \right)$, $\gamma = \frac{1}{d^{1/2}}$, $B_G = \frac{\rho_-}{k} p_A = O\left(\frac{\sqrt{k}}{d}\right)$ and $B_{x_1} = \frac{\rho_-}{d}$, $B_{x_2} = d$. We also have $\text{OPT}_{d;3r(k+1);B_F;S_{p_1};B_G} = 0$.

Adjust σ_w such that $\tilde{b} = C_b \frac{\rho_-}{\tau r k \log d} \sigma_w = \left(\frac{B_G^{\frac{1}{4}} B_{a_2} B_b^{\frac{3}{4}}}{\sqrt{r B_{a_1}}} \right)$. Injecting above parameters into Theorem 3.12, we have with probability at least $1 - \delta$ over the initialization, with proper hyper-parameters, there exists $t \geq [T]$ such that

$$\Pr[\text{sign}(f^{(t)}(\mathbf{x})) \neq y] \leq \frac{2^{\frac{\rho_-}{2r}} \frac{\rho_-}{k}}{d^{(-3)=2}} + O \left(\frac{r B_{a_1} B_{x_1} B_{x_2}^{\frac{1}{4}} (\log n)^{\frac{1}{4}}}{\rho_- B_G n^{\frac{1}{4}}} \right) + \epsilon/2 \leq \frac{3r \frac{\rho_-}{k}}{d^{(-3)=2}} + \epsilon. \quad \square$$

E.5 Uniform Parity Functions

We consider the sparse parity problem over the uniform data distribution studied in [15]. We use the properties of the problem to prove the key lemma (i.e., the existence of good networks) in our framework and then derive the final guarantee from our theorem of the simple setting (Theorem 3.4). We provide Theorem E.31 as (1) use it as a warm-up and (2) follow the original analysis in [15] to give a comparison. We will provide Theorem E.40 as an alternative version that trains both layers.

Consider the same data distribution in Appendix E.4.1 and Definition E.22 with the following assumptions.

Assumption E.28 (Uniform Parity Functions). *We follow the data distribution in Appendix E.4.1. Let $r = 1$, $p_A = 0$, $p_0 = \frac{1}{2}$, $\mathbf{M} = I_{d \times d}$ and $d = 2k^2$, and k is an even number.*

We denote this data distribution as $\mathcal{D}_{\text{parity-uniform}}$ setting.

To apply our framework, again we only need to compute the parameters in the Gradient Feature set and the corresponding optimal approximation loss. To this end, we first define the gradient features: let

$$D = \frac{\sum_{I \in \mathcal{A}} \mathbf{M}_I}{k \sum_{I \in \mathcal{A}} \mathbf{M}_I k_2}. \quad (445)$$

We follow the initialization and training dynamic in [15].

Initialization and Loss. We use hinge loss and we have unbiased initialization, for all $i \geq [m]$,

$$\mathbf{a}_i^{(0)} = \text{Unif}(f-1g), \mathbf{w}_i^{(0)} = \text{Unif}(f-1g^d), \mathbf{b}_i = \text{Unif}(f-1+1/k, \dots, 1-1/kg). \quad (446)$$

Training Process. We use the following one-step training algorithm for this specific data distribution.

Algorithm 4 Network Training via Gradient Descent [15]. Special case of Algorithm 2

Initialize $(\mathbf{a}^{(0)}, \mathbf{W}^{(0)}, \mathbf{b})$ as in Equation (8) and Equation (446); Sample $Z \sim D_{\text{parity-uniform}}^n$
 $\mathbf{W}^{(1)} = \mathbf{W}^{(0)} - \eta^{(1)}(r_{\mathbf{W}} \tilde{L}_Z(f^{(0)})) + \lambda^{(1)} \mathbf{W}^{(0)}$
 $\mathbf{a}^{(1)} = \mathbf{a}^{(0)} - \eta^{(1)}(r_{\mathbf{a}} \tilde{L}_Z(f^{(0)})) + \lambda^{(1)} \mathbf{a}^{(0)}$
for $t = 2$ **to** T **do**
 $\mathbf{a}^{(t)} = \mathbf{a}^{(t-1)} - \eta^{(t)} r_{\mathbf{a}} \tilde{L}_Z(f^{(t-1)})$
end for

Use the notation in Section 5.3 of [92], for every $S \subseteq [n]$, s.t. $|S| = k$, we define

$$\xi_k := \widehat{\text{Maj}}(S) = \binom{k-1}{\frac{d-1}{2}} 2^{-(d-1)} \binom{d-1}{\frac{d-1}{2}}. \quad (447)$$

Lemma E.29 (Uniform Parity Functions: Existence of Good Networks. Rephrase of Lemma 5 in [15]). *For every $\epsilon, \delta \in (0, 1/2)$, denoting $\tau = \frac{1}{16k} \frac{|k-1|}{2d \log(32k^3 d)}$, let $\eta^{(1)} = \frac{1}{k|k-1|}$, $\lambda^{(1)} = \frac{1}{(1)}$, $m = k \cdot 2^k \log(k/\delta)$, $n = \frac{2}{\delta} \log(4dm/\delta)$ and $d = \lceil k^4 \log(kd/\epsilon) \rceil$, w.p. at least $1 - 2\delta$ over the initialization and the training samples, there exists $\mathbf{a} \in \mathbb{R}^m$ with $\|\mathbf{a}\|_{k\infty} \leq 8k$ and $\|\mathbf{a}\|_{k_2} \leq 8k \cdot \frac{1}{k}$ such that $f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}$ satisfies*

$$L_{\mathcal{D}_{\text{parity-uniform}}}(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}) \leq \epsilon. \quad (448)$$

Additionally, it holds that $k\sigma(\mathbf{W}^{(1)\top} \mathbf{x} - \mathbf{b})_{k\infty} \leq d + 1$.

Remark E.30. In [15], they update the bias term in the first gradient step. However, if we check the proof carefully, we can see that the fixed bias still goes through all their analysis.

E.5.1 Uniform Parity Functions: Final Guarantee

Considering training by Algorithm 4, we have the following results.

Theorem E.31 (Uniform Parity Functions: Main Result). *Fix $\epsilon \in (0, 1/2)$ and let $m = \lceil k \cdot 2^k \log(k/\epsilon) \rceil$, $n = \lceil k^{7-6} d \binom{d}{k-1} \log(kd/\epsilon) \log(dm/\epsilon) + \frac{k^3 m d^2}{2} \rceil$, $d = \lceil k^4 \log(kd/\epsilon) \rceil$. Let $\eta^{(1)} = \frac{1}{k|k-1|}$, $\lambda^{(1)} = \frac{1}{(1)}$, and $\eta = \eta^{(t)} = \frac{1}{(d^2 m)}$, for all $t \in \{2, 3, \dots, T\}$. If $T = \lceil \frac{k^3 m d^2}{\epsilon} \rceil$, then training by Algorithm 4 with hinge loss, w.h.p. over the initialization and the training samples, there exists $t \in [T]$ such that*

$$\Pr[\text{sign}(f^{(t)})(\mathbf{x}) \neq y] = L_{\mathcal{D}_{\text{parity-uniform}}}(f_{(\mathbf{a}^{(t)}, \mathbf{W}^{(1)}, \mathbf{b})}) \leq \epsilon. \quad (449)$$

Proof of Theorem E.31. By Lemma E.29, w.h.p., we have for properly chosen hyper-parameters,

$$\text{OPT}_{\mathbf{W}^{(1)}, \mathbf{b}; B_{a_2}} = L_{\mathcal{D}_{\text{parity-uniform}}}(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}) \leq \frac{\epsilon}{3}. \quad (450)$$

We compute the L -smooth constant of $\tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})})$ to \mathbf{a} .

$$\left\| \nabla_{\mathbf{a}} \tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}_1, \mathbf{W}^{(1)}, \mathbf{b})}) - \nabla_{\mathbf{a}} \tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}_2, \mathbf{W}^{(1)}, \mathbf{b})}) \right\|_2 \quad (451)$$

$$= \left\| \frac{1}{n} \sum_{\mathbf{x} \in \mathcal{Z}} \left[\left(\ell'(y f_{(\mathbf{a}_1, \mathbf{W}^{(1)}, \mathbf{b})}(\mathbf{x})) - \ell'(y f_{(\mathbf{a}_2, \mathbf{W}^{(1)}, \mathbf{b})}(\mathbf{x})) \right) \sigma(\mathbf{W}^{(1)\top} \mathbf{x} - \mathbf{b}) \right] \right\|_2 \quad (452)$$

$$\left\| \frac{1}{n} \sum_{\mathbf{x} \in \mathcal{Z}} \left[\left| f_{(\mathbf{a}_1, \mathbf{W}^{(1)}, \mathbf{b})}(\mathbf{x}) - f_{(\mathbf{a}_2, \mathbf{W}^{(1)}, \mathbf{b})}(\mathbf{x}) \right| \sigma(\mathbf{W}^{(1)\top} \mathbf{x} - \mathbf{b}) \right] \right\|_2 \quad (453)$$

$$\frac{1}{n} \sum_{\mathbf{x} \in \mathcal{Z}} \left[k \mathbf{a}_1 - \mathbf{a}_2 k_2 \left\| \sigma(\mathbf{W}^{(1)\top} \mathbf{x} - \mathbf{b}) \right\|_2^2 \right]. \quad (454)$$

By the Lemma E.29, we have $k \sigma(\mathbf{W}^{(1)\top} \mathbf{x} - \mathbf{b}) k_{\infty} \leq d + 1$. Thus, we have,

$$L = O \left(\frac{1}{n} \sum_{\mathbf{x} \in \mathcal{Z}} \left\| \sigma(\mathbf{W}^{(1)\top} \mathbf{x} - \mathbf{b}) \right\|_2^2 \right) \quad (455)$$

$$O(d^2 m). \quad (456)$$

This means that we can let $\eta = \frac{1}{d^2 m}$ and we will get our convergence result. Note that we have $\mathbf{a}^{(1)} = \mathbf{0}$ and $k \mathbf{a} k_2 = O\left(\frac{1}{k} \frac{1}{k}\right)$. So, if we choose $T = O\left(\frac{k^3}{\epsilon}\right)$, there exists $t \geq [T]$ such that $\tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}^{(t)}, \mathbf{W}^{(1)}, \mathbf{b})}) - \tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}) = O\left(\frac{L \|\mathbf{a}^{(1)} - \mathbf{a}\|_2^2}{T}\right) \leq \epsilon/3$.

We also have $\sqrt{\frac{\|\mathbf{a}\|_2^2 (\|\mathbf{W}^{(1)}\|_F^2 B_G^2 + \|\mathbf{b}\|_2^2)}{n}} \leq \frac{\epsilon}{3}$. Then our theorem gets proved by Theorem 3.4. \square

E.6 Uniform Parity Functions: Alternative Analysis

It is also possible to unify [15] into our general Gradient Feature Learning Framework by mildly modifying the framework in Theorem 3.12. In order to do that, we first need to use a different metric in the definition of gradient features.

E.6.1 Modified General Feature Learning Framework for Uniform Parity Functions

Definition E.32 (Gradient Feature with Infinity Norm). For a unit vector $D \in \mathbb{R}^d$ with $k D k_2 = 1$, and a $\gamma_{\infty} \in (0, 1)$, a direction neighborhood (cone) $\mathcal{C}_{D, \gamma_{\infty}}^{\infty}$ is defined as: $\mathcal{C}_{D, \gamma_{\infty}}^{\infty} := \left\{ \mathbf{w} \mid \left\| \frac{\mathbf{w}}{\|\mathbf{w}\|} - D \right\|_{\infty} < \gamma_{\infty} \right\}$. Let $\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$ be random variables drawn from some distribution W, B . A Gradient Feature set with parameters $p, \gamma_{\infty}, B_G, B_{G_1}$ is defined as:

$$S_{p, \gamma_{\infty}; B_G, B_{G_1}}^{\infty}(W, B) := \left\{ (D, s) \mid \Pr_{\mathbf{w}, b} \left[G(\mathbf{w}, b) \in \mathcal{C}_{D, \gamma_{\infty}}^{\infty}, B_{G_1} \leq k G(\mathbf{w}, b) k_2 \leq B_G, s = \frac{b}{\|b\|} \right] \geq p \right\}.$$

When clear from context, write it as $S_{p, \gamma_{\infty}; B_G, B_{G_1}}^{\infty}$.

Definition E.33 (Optimal Approximation via Gradient Features with Infinity Norm). The Optimal Approximation network and loss using gradient feature induced networks $\mathcal{F}_{d, r; B_F; S_{p, \gamma_{\infty}; B_G, B_{G_1}}^1}$ are defined as:

$$f^* := \operatorname{argmin}_{f \in \mathcal{F}_{d, r; B_F; S_{p, \gamma_{\infty}; B_G, B_{G_1}}^1}} L_{\mathcal{D}}(f), \quad (457)$$

$$\operatorname{OPT}_{d, r; B_F; S_{p, \gamma_{\infty}; B_G, B_{G_1}}^1} := \min_{f \in \mathcal{F}_{d, r; B_F; S_{p, \gamma_{\infty}; B_G, B_{G_1}}^1}} L_{\mathcal{D}}(f). \quad (458)$$

We consider the data distribution in Appendix E.4.1 with Assumption E.28, i.e., $D_{\text{parity-uniform}}$ in Appendix E.5. Note that with this dataset, we have $k \mathbf{x} k_{\infty} \leq B_{x_{\infty}} = 1$. We use the following unbiased initialization:

$$\begin{aligned} \text{for } i \in \{1, \dots, mg\}: \quad & \mathbf{a}_i^{(0)} \sim \mathcal{N}(0, \sigma_a^2), \mathbf{w}_i^{(0)} \sim \mathcal{U}(\mathbb{R}^d), \mathbf{b}_i = \mathbf{b} - 1, \\ \text{for } i \in \{fm + 1, \dots, 2mg\}: \quad & \mathbf{a}_i^{(0)} = \mathbf{a}_{i-m}^{(0)}, \mathbf{w}_i^{(0)} = \mathbf{w}_{i-m}^{(0)}, \mathbf{b}_i = \mathbf{b}_{i-m}, \\ \text{for } i \in \{f2m + 1, \dots, 4mg\}: \quad & \mathbf{a}_i^{(0)} = \mathbf{a}_{i-2m}^{(0)}, \mathbf{w}_i^{(0)} = \mathbf{w}_{i-2m}^{(0)}, \mathbf{b}_i = \mathbf{b}_{i-2m} \end{aligned} \quad (459)$$

Let r_i denote the gradient of the i -th neuron $r_{\mathbf{w}_i} L_{\mathcal{D}}(f^{(0)})$. Denote the subset of neurons with nice gradients approximating feature (D, s) as:

$$G_{(D;s):Nice}^{\infty} := \left\{ i \in [2m] : s = \frac{\mathbf{b}_i}{\|\mathbf{b}_i\|}, \left\| \frac{r_i}{k r_i k} - D \right\|_{\infty} \leq \gamma_{\infty}, \left| \mathbf{a}_i^{(0)} \right|_{B_{G1}} \leq k r_i k_2, \left| \mathbf{a}_i^{(0)} \right|_{B_G} \right\}.$$

Lemma E.34 (Existence of Good Networks. Modified Version of Lemma 3.14 Under Uniform Parity Setting). *Let $\lambda^{(1)} = \frac{1}{(1)}$. For any $B \geq (0, B_b)$, let $\sigma_a = \left(\frac{b}{-\nu(0) (1) B_G B} \right)$ and $\delta = 2re^{-\sqrt{mp}} + \frac{1}{\delta^2}$. Then, with probability at least $1 - \delta$ over the initialization, there exists \mathbf{a}_i 's such that $f_{(\mathbf{a}; \mathbf{W}^{(1); \mathbf{b}})}(\mathbf{x}) = \sum_{i=1}^{4m} \mathbf{a}_i \sigma \left(\langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle - \mathbf{b}_i \right)$ satisfies*

$$L_{\mathcal{D}}(f_{(\mathbf{a}; \mathbf{W}^{(1); \mathbf{b}})}) \leq r B_{a1} \left(\frac{B_{x1} B_{G1} B_b}{\delta \sqrt{mp} B_G B} + \sqrt{2 \log(d) d} \gamma_{\infty} + B \right) + \text{OPT}_{d;r;B_F;S_{p_1}^1; \dots; B_G; B_{G1}},$$

$$\text{and } k \mathbf{a} k_0 = O\left(r(mp)^{\frac{1}{2}}\right), k \mathbf{a} k_2 = O\left(\frac{B_{a2} B_b}{b(mp)^{\frac{1}{4}}}\right), k \mathbf{a} k_{\infty} = O\left(\frac{B_{a1} B_b}{b(mp)^{\frac{1}{2}}}\right).$$

Proof of Lemma E.34. Recall $f^*(\mathbf{x}) = \sum_{j=1}^r \mathbf{a}_j^* \sigma(\langle \mathbf{w}_j^*, \mathbf{x} \rangle - \mathbf{b}_j^*)$, where $f^* \in \mathcal{F}_{d;r;B_F;S_{p_1}^1; \dots; B_G; B_{G1}}$ is defined in Definition E.33 and let $s_j^* = \frac{\mathbf{b}_j}{\|\mathbf{b}_j\|}$. By Lemma D.3, with probability at least $1 - \delta_1$, $\delta_1 = 2re^{-cmp}$, for all $j \in [r]$, we have $j G_{(\mathbf{w}_j; s_j):Nice}^{\infty} \geq \frac{mp}{4}$. Then for all $i \in G_{(\mathbf{w}_j; s_j):Nice}^{\infty} \cap [2m]$, we have $\ell'(0) \eta^{(1)} G(\mathbf{w}_i^{(0)}, \mathbf{b}_i) \frac{\mathbf{b}_j}{b}$ only depend on $\mathbf{w}_i^{(0)}$ and \mathbf{b}_i , which is independent of $\mathbf{a}_i^{(0)}$. Given Definition E.32, we have

$$\ell'(0) \eta^{(1)} k G(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \frac{\mathbf{b}_j^*}{b} \geq \left[\ell'(0) \eta^{(1)} B_{x1} \frac{B_b}{b}, \ell'(0) \eta^{(1)} B_{x1} \frac{B_b}{b} \right]. \quad (460)$$

We split $[r]$ into $\mathcal{I} = \{j \in [r] : \|\mathbf{b}_j^*\| < B/g\}$, $\mathcal{J} = \{j \in [r] : \|\mathbf{b}_j^*\| \geq B/g\}$ and $\mathcal{K} = \{j \in [r] : \|\mathbf{b}_j^*\| \geq B/g\}$. Let $\epsilon_a = \frac{B_{G1} B_b}{\sqrt{mp} B_G B}$. Then we know that for all $j \in \mathcal{I} \cup \mathcal{K}$, for all $i \in G_{(\mathbf{w}_j; s_j):Nice}^{\infty}$ we have

$$\Pr_{\mathbf{a}_i^{(0)} \sim \mathcal{N}(0; \frac{2}{g})} \left[\left| \mathbf{a}_i^{(0)} \ell'(0) \eta^{(1)} k G(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \frac{\mathbf{b}_j^*}{b} - 1 \right| \leq \epsilon_a \right] \quad (461)$$

$$= \Pr_{\mathbf{a}_i^{(0)} \sim \mathcal{N}(0; \frac{2}{g})} \left[1 - \epsilon_a \leq \mathbf{a}_i^{(0)} \ell'(0) \eta^{(1)} k G(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \frac{\mathbf{b}_j^*}{b} \leq 1 + \epsilon_a \right] \quad (462)$$

$$= \Pr_{g \sim \mathcal{N}(0;1)} \left[1 - \epsilon_a \leq g \left(\frac{k G(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \mathbf{b}_j^*}{B_G B} \right) \leq 1 + \epsilon_a \right] \quad (463)$$

$$= \Pr_{g \sim \mathcal{N}(0;1)} \left[(1 - \epsilon_a) \left(\frac{B_G B}{k G(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \mathbf{b}_j^*} \right) \leq g \leq (1 + \epsilon_a) \left(\frac{B_G B}{k G(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \mathbf{b}_j^*} \right) \right]$$

$$= \left(\frac{\epsilon_a B_G B}{k G(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \mathbf{b}_j^*} \right) \quad (464)$$

$$\left(\frac{\epsilon_a B_G B}{B_{G1} B_b} \right) \quad (465)$$

$$= \left(\frac{1}{\delta \sqrt{mp}} \right). \quad (466)$$

Thus, with probability $\left(\frac{1}{\sqrt{mp}} \right)$ over $\mathbf{a}_i^{(0)}$, we have

$$\left| \mathbf{a}_i^{(0)} \ell'(0) \eta^{(1)} k G(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \frac{\mathbf{b}_j^*}{b} - 1 \right| \leq \epsilon_a, \quad \left| \mathbf{a}_i^{(0)} \right| = O\left(\frac{b}{\ell'(0) \eta^{(1)} B_G B}\right). \quad (467)$$

Similarly, for $j \geq 2$, for all $i \geq G_{(\mathbf{w}_j; s_j)}^\infty$:Nice, with probability $\left(\frac{1}{\sqrt{mp}}\right)$ over $\mathbf{a}_i^{(0)}$, we have

$$\left| \mathbf{a}_i^{(0)} \ell'(0) \eta^{(1)} kG(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \frac{B}{\bar{b}} - 1 \right| \leq \epsilon_a, \quad \left| \mathbf{a}_i^{(0)} \right| = O\left(\frac{\bar{b}}{\ell'(0) \eta^{(1)} B_G B}\right). \quad (468)$$

For all $j \geq [r]$, let $J_j = G_{(\mathbf{w}_j; s_j)}^\infty$:Nice be the set of i 's such that condition Equation (467) or Equation (468) are satisfied. By Chernoff bound and union bound, with probability at least $1 - \delta_2$, $\delta_2 = re^{-\sqrt{mp}}$, for all $j \geq [r]$ we have $|J_j| \geq \frac{1}{2} \binom{[r]}{j}$. We have for $\delta_j \geq \frac{1}{2} + [r]$, $\delta_i \geq j$,

$$\begin{aligned} & \left| \frac{j \mathbf{b}_j^*}{\bar{b}} \langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle - \langle \mathbf{w}_j^*, \mathbf{x} \rangle \right| \quad (469) \\ & \left| \left\langle \mathbf{a}_i^{(0)} \ell'(0) \eta^{(1)} kG(\mathbf{w}_i^{(0)}, \mathbf{b}_i) k_2 \frac{j \mathbf{b}_j^*}{\bar{b}} \frac{\mathbf{w}_i^{(1)}}{k \mathbf{w}_i^{(1)} k_2} - \frac{\mathbf{w}_j^*}{k \mathbf{w}_j^* k_2}, \mathbf{x} \right\rangle + \left\langle \frac{\mathbf{w}_i^{(1)}}{k \mathbf{w}_i^{(1)} k_2} - \frac{\mathbf{w}_j^*}{k \mathbf{w}_j^* k_2}, \mathbf{x} \right\rangle \right| \\ & \leq \epsilon_a k \mathbf{x} k_2 + \sqrt{2 \log(d) d} \gamma_\infty. \quad (470) \end{aligned}$$

With probability $1 - \frac{\delta_2}{d^2}$ by Hoeffding's inequality. Similarly, for $\delta_j \geq \frac{1}{2}$, $\delta_i \geq j$,

$$\left| \frac{B}{\bar{b}} \langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle - \langle \mathbf{w}_j^*, \mathbf{x} \rangle \right| \leq \epsilon_a k \mathbf{x} k_2 + \sqrt{2 \log(d) d} \gamma_\infty. \quad (471)$$

If $i \geq j$, $j \geq \frac{1}{2} + [r]$, set $\mathbf{a}_i = \mathbf{a}_j^* \frac{|\mathbf{b}_j|}{|j| \bar{b}}$, if $i \geq j$, $j \geq 2$, set $\mathbf{a}_i = \mathbf{a}_j^* \frac{B}{|j| \bar{b}}$, otherwise set $\mathbf{a}_i = 0$, we have $k \mathbf{a}_0 k_0 = O\left(r(mp)^{\frac{1}{2}}\right)$, $k \mathbf{a}_k k_2 = O\left(\frac{B a_2 B_p}{b(mp)^{\frac{1}{4}}}\right)$, $k \mathbf{a}_\infty k_\infty = O\left(\frac{B a_1 B_p}{b(mp)^{\frac{1}{2}}}\right)$.

Finally, we have

$$L_{\mathcal{D}}(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}) \quad (472)$$

$$= L_{\mathcal{D}}(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}) - L_{\mathcal{D}}(f^*) + L_{\mathcal{D}}(f^*) \quad (473)$$

$$E_{(\mathbf{x}, \mathbf{y})} \left[\left| f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}(\mathbf{x}) - f^*(\mathbf{x}) \right| \right] + L_{\mathcal{D}}(f^*) \quad (474)$$

$$E_{(\mathbf{x}, \mathbf{y})} \left[\left| \sum_{i=1}^m \mathbf{a}_i \sigma(\langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle - \mathbf{b}) + \sum_{i=m+1}^{2m} \mathbf{a}_i \sigma(\langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle + \mathbf{b}) - \sum_{j=1}^r \mathbf{a}_j^* \sigma(\langle \mathbf{w}_j^*, \mathbf{x} \rangle - \mathbf{b}_j^*) \right| \right] + L_{\mathcal{D}}(f^*) \quad (475)$$

$$E_{(\mathbf{x}, \mathbf{y})} \left[\left| \sum_{j \in \mathcal{I}} \sum_{i \in \mathcal{J}} \mathbf{a}_j^* \frac{1}{j} \left| \frac{j \mathbf{b}_j^*}{\mathbf{b}} \sigma(\langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle - \mathbf{b}) - \sigma(\langle \mathbf{w}_j^*, \mathbf{x} \rangle - \mathbf{b}_j^*) \right| \right| \right] \quad (476)$$

$$+ E_{(\mathbf{x}, \mathbf{y})} \left[\left| \sum_{j \in \mathcal{I}} \sum_{i \in \mathcal{J}} \mathbf{a}_j^* \frac{1}{j} \left| \frac{j \mathbf{b}_j^*}{\mathbf{b}} \sigma(\langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle + \mathbf{b}) - \sigma(\langle \mathbf{w}_j^*, \mathbf{x} \rangle - \mathbf{b}_j^*) \right| \right| \right] \quad (477)$$

$$+ E_{(\mathbf{x}, \mathbf{y})} \left[\left| \sum_{j \in \mathcal{I}} \sum_{i \in \mathcal{J}} \mathbf{a}_j^* \frac{1}{j} \left| \frac{B}{\mathbf{b}} \sigma(\langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle - \mathbf{b}) - \sigma(\langle \mathbf{w}_j^*, \mathbf{x} \rangle - \mathbf{b}_j^*) \right| \right| \right] + L_{\mathcal{D}}(f^*) \quad (478)$$

$$E_{(\mathbf{x}, \mathbf{y})} \left[\left| \sum_{j \in \mathcal{I}} \sum_{i \in \mathcal{J}} \mathbf{a}_j^* \frac{1}{j} \left| \frac{j \mathbf{b}_j^*}{\mathbf{b}} \langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle - \langle \mathbf{w}_j^*, \mathbf{x} \rangle \right| \right| \right] \quad (479)$$

$$+ E_{(\mathbf{x}, \mathbf{y})} \left[\left| \sum_{j \in \mathcal{I}} \sum_{i \in \mathcal{J}} \mathbf{a}_j^* \frac{1}{j} \left| \frac{j \mathbf{b}_j^*}{\mathbf{b}} \langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle + \langle \mathbf{w}_j^*, \mathbf{x} \rangle \right| \right| \right] \quad (480)$$

$$+ E_{(\mathbf{x}, \mathbf{y})} \left[\left| \sum_{j \in \mathcal{I}} \sum_{i \in \mathcal{J}} \mathbf{a}_j^* \frac{1}{j} \left| \frac{B}{\mathbf{b}} \langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle + B - \langle \mathbf{w}_j^*, \mathbf{x} \rangle \right| \right| \right] + L_{\mathcal{D}}(f^*) \quad (481)$$

$$r k \mathbf{a}^* k_{\infty} (\epsilon_a E_{(\mathbf{x}, \mathbf{y})} k \mathbf{x} k_2 + \sqrt{2 \log(d) d} \gamma_{\infty}) + j j k \mathbf{a}^* k_{\infty} B + L_{\mathcal{D}}(f^*) \quad (482)$$

$$r B_{a1} (\epsilon_a B_{x1} + \sqrt{2 \log(d) d} \gamma_{\infty}) + j j B_{a1} B + \text{OPT}_{d, r; B_F; S_{p'}^j; B_G; B_{G1}}. \quad (483)$$

We finish the proof by union bound and $\delta = \delta_1 + \delta_2 + \frac{1}{\varrho}$. \square

Lemma E.35 (Empirical Gradient Concentration Bound for Single Coordinate). *For $i \geq 2$ [m], when $n \geq (\log(d))^6$, with probability at least $1 - O\left(\exp\left(-\frac{1}{n^{\frac{1}{3}}}\right)\right)$ over training samples, we have*

$$\left| \frac{\partial \tilde{L}_{\mathcal{Z}}(f)}{\partial \mathbf{w}_{ij}} - \frac{\partial L_{\mathcal{D}}(f)}{\partial \mathbf{w}_{ij}} \right| = O\left(\frac{j \mathbf{a}_j B_{x\infty}}{n^{\frac{1}{3}}}\right), \quad \forall j \geq 2 [d]. \quad (484)$$

Proof of Lemma E.35. First, we define,

$$z_{ij}^{(\ell)} = \ell'(y^{(\ell)} f(\mathbf{x}^{(\ell)})) y^{(\ell)} \left[\sigma'(\langle \mathbf{w}_i, \mathbf{x}^{(\ell)} \rangle - \mathbf{b}_i) \mathbf{x}_j^{(\ell)} \right] \quad (485)$$

$$E_{(\mathbf{x}, \mathbf{y})} [\ell'(y f(\mathbf{x})) y [\sigma'(\langle \mathbf{w}_i, \mathbf{x} \rangle - \mathbf{b}_i)] \mathbf{x}_j]. \quad (486)$$

As $j \ell'(z) j \leq 1, j y j \leq 1, j \sigma'(z) j \leq 1$, we have $z_{ij}^{(\ell)}$ is zero-mean random variable with $|z_{ij}^{(\ell)}| \leq 2B_{x\infty}$

as well as $E \left[\left| z_{ij}^{(\ell)} \right|_2^2 \right] \leq 4B_{x\infty}^2$. Then by Bernstein Inequality, for $0 < z < 2B_{x\infty}$, we have

$$\Pr \left(\left| \frac{\partial \tilde{L}_{\mathcal{Z}}(f)}{\partial \mathbf{w}_{ij}} - \frac{\partial L_{\mathcal{D}}(f)}{\partial \mathbf{w}_{ij}} \right| \geq j \mathbf{a}_j z \right) = \Pr \left(\left| \frac{1}{n} \sum_{i \in [n]} z_{ij}^{(\ell)} \right| \geq z \right) \quad (487)$$

$$\leq \exp \left(-n \frac{z^2}{8B_{x\infty}^2} \right). \quad (488)$$

Thus, for some $i \geq [m]$, when $n \geq (\log(d))^6$, with probability at least $1 - O\left(\exp\left(-\frac{1}{n^{\frac{1}{3}}}\right)\right)$, from a union bound over $j \geq [d]$, we have, for $\delta_j \geq [d]$,

$$\left| \frac{\partial \tilde{L}_{\mathcal{Z}}(f)}{\partial \mathbf{w}_{ij}} - \frac{\partial L_{\mathcal{D}}(f)}{\partial \mathbf{w}_{ij}} \right| = O\left(\frac{j \mathbf{a}_j B_{X\infty}}{n^{\frac{1}{3}}}\right). \quad (489)$$

□

Lemma E.36 (Existence of Good Networks under Empirical Risk. Modified version of Lemma D.13 Under Uniform Parity Setting). *Suppose $n > \left(\left(\frac{B_x}{\sqrt{B_{x2}}} + \log \frac{1}{\rho} + \frac{B_{x1}}{B_{G1}^{\frac{1}{\rho}}(0)} + \frac{B_{x1}}{B_{G1}^{\frac{1}{\rho}}(0)}\right)^3 + (\log(d))^6\right)$. Let $\lambda^{(1)} = \frac{1}{(1)}$. For any $B \geq (0, B_b)$, let $\sigma_a = \left(\frac{b}{-|\cdot|^{\rho}(0)}\right)$ and $\delta = 2re^{-\frac{\rho}{2}} + \frac{1}{d^2}$. Then, with probability at least $1 - \delta$ over the initialization and training samples, there exists \mathbf{a}_j 's such that $f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}(\mathbf{x}) = \sum_{i=1}^{4m} \mathbf{a}_i \sigma\left(\langle \mathbf{w}_i^{(1)}, \mathbf{x} \rangle - \mathbf{b}_i\right)$ satisfies*

$$L_{\mathcal{D}}(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}) \leq r B_{a1} \left(\frac{2B_{x1} B_{G1} B_b}{\rho^{mp} B_G B} + \sqrt{2 \log(d) d} \left(\gamma_{\infty} + O\left(\frac{B_{x\infty}}{B_G j \ell'(0) j n^{\frac{1}{3}}}\right) \right) + B \right) + \text{OPT}_{d; r; B_F; S_{p; \cdot; B_G; B_{G1}}^1}, \quad (490)$$

$$\text{and } \mathbf{a} \mathbf{k}_0 = O\left(r(mp)^{\frac{1}{2}}\right), \mathbf{a} \mathbf{k}_2 = O\left(\frac{B_{a2} B_b}{b(mp)^{\frac{1}{4}}}\right), \mathbf{a} \mathbf{k}_{\infty} = O\left(\frac{B_{a1} B_b}{b(mp)^{\frac{1}{2}}}\right).$$

Proof of Lemma E.36. Denote $\rho = O\left(\exp\left(-\frac{1}{n^{\frac{1}{3}}}\right)\right)$ and $\beta = O\left(\frac{B_{x1}}{n^{\frac{1}{3}}}\right)$. Note that by symmetric initialization, we have $\ell'(y f^{(0)}(\mathbf{x})) = j \ell'(0) j$ for any $\mathbf{x} \geq X$, so that, by Lemma E.35, we have $\left| \tilde{G}(\mathbf{w}_i^{(0)}, \mathbf{b}_i)_j - G(\mathbf{w}_i^{(0)}, \mathbf{b}_i)_j \right| \leq \frac{\rho}{|\cdot|^{\rho}(0)}$ with probability at least $1 - \rho$. Thus, by union bound, we can see that $S_{p; \cdot; B_G; B_{G1}}^{\infty} \leq \frac{\tilde{S}_{p; \cdot; 1 + \frac{B_{x1}}{B_G j \ell'(0) j}; B_G - \frac{1}{j^{\rho}(0) j}; B_{G1} + \frac{1}{j^{\rho}(0) j}}}{\rho}$. Consequently, we have $\text{OPT}_{d; r; B_F; \tilde{S}_{p; \cdot; 1 + \frac{B_{x1}}{B_G j \ell'(0) j}; B_G - \frac{1}{j^{\rho}(0) j}; B_{G1} + \frac{1}{j^{\rho}(0) j}}^1 \leq \frac{\text{OPT}_{d; r; B_F; S_{p; \cdot; 1 + \frac{B_{x1}}{B_G j \ell'(0) j}; B_G - \frac{1}{j^{\rho}(0) j}; B_{G1} + \frac{1}{j^{\rho}(0) j}}^1}{\rho}$. Exactly follow the proof in Lemma D.4 by replacing $S_{p; \cdot; B_G; B_{G1}}^{\infty}$ to $\tilde{S}_{p; \cdot; 1 + \frac{B_{x1}}{B_G j \ell'(0) j}; B_G - \frac{1}{j^{\rho}(0) j}; B_{G1} + \frac{1}{j^{\rho}(0) j}}^{\infty}$. Then, we finish the proof by $\rho \leq \frac{\rho}{2}, \frac{\rho}{|\cdot|^{\rho}(0)} \leq (1 - \frac{\rho}{2}) B_G, \frac{\rho}{|\cdot|^{\rho}(0)} \leq (1 - \frac{\rho}{2}) B_{G1}$. □

Theorem E.37 (Online Convex Optimization under Empirical Risk. Modified version of Theorem D.17 Under Uniform Parity Setting). *Consider training by Algorithm 1, and any $\delta \geq (0, 1)$. Assume $d \geq \log m, \delta \geq O\left(\frac{1}{d^2}\right)$. Set*

$$\sigma_w > 0, \quad \mathbf{b} > 0, \quad \eta^{(t)} = \eta, \quad \lambda^{(t)} = 0 \text{ for all } t \geq 2, 3, \dots, Tg,$$

$$\eta^{(1)} = \left(\frac{\min\{O(\eta), O(\eta \mathbf{b}) g\}}{\ell'(0)(B_{x1} \sigma_w d + \mathbf{b})} \right), \quad \lambda^{(1)} = \frac{1}{\eta^{(1)}}, \quad \sigma_a = \left(\frac{\mathbf{b}(mp)^{\frac{1}{4}}}{\ell'(0) \eta^{(1)} B_{x1} B_G B_b} \right).$$

Let $0 < T \eta B_{x1} \leq o(1)$, $m \geq \left(\frac{1}{\sqrt{\rho}} + \frac{1}{\rho} (\log(\frac{1}{\rho}))^2\right)$ and $n > \left(\left(\frac{B_x}{\sqrt{B_{x2}}} + \log \frac{Tm}{\rho} + \left(1 + \frac{1}{B_G} + \frac{1}{B_{G1}}\right) \frac{B_{x1}}{|\cdot|^{\rho}(0)}\right)^3\right)$. With probability at least $1 - \delta$ over the

initialization and training samples, there exists $t \geq 2 [T]$ such that

$$L_{\mathcal{D}}(f^{(t)}) \tag{491}$$

$$\text{OPT}_{d;r;B_F;S_{p_1};B_G} + rB_{a1} \left(\frac{2^{\rho} \overline{B_{x1} B_{G1}}}{(mp)^{\frac{1}{4}}} \sqrt{\frac{B_b}{B_G}} + \sqrt{2 \log(d)d} \left(\gamma_{\infty} + O\left(\frac{B_{x\infty}}{B_G j_{\ell'}(0) j n^{\frac{1}{3}}}\right) \right) \right)$$

$$+ \eta \left(\overline{rB_{a2} B_b T \eta B_{x1}^2} + m \overline{b} \right) O\left(\frac{\overline{b^{\rho} \log m B_{x1}} (mp)^{\frac{1}{4}}}{\overline{B_b B_G}} + 1 \right) + O\left(\frac{B_{a2}^2 B_b^2}{\eta T \overline{b}^2 (mp)^{\frac{1}{2}}} \right) \tag{492}$$

$$+ \frac{1}{n^{\frac{1}{3}}} O\left(\left(\frac{rB_{a1} B_b}{\overline{b}} + m \left(\frac{\overline{b^{\rho} \log m} (mp)^{\frac{1}{4}}}{\overline{B_b B_G}} + \frac{\overline{b}}{B_{x1}} \right) \right) \right) \tag{493}$$

$$\left(\left(\frac{\overline{b^{\rho} \log m} (mp)^{\frac{1}{4}}}{\overline{B_b B_G}} + T \eta^2 B_{x1} \overline{b} \right) B_x + \overline{b} \right) + 2 \right) \tag{494}$$

$$+ \frac{1}{n^{\frac{1}{3}}} O\left(m \eta \left(\frac{\overline{b^{\rho} \log m} (mp)^{\frac{1}{4}}}{\overline{B_b B_G}} + T \eta^2 B_{x1} \overline{b} \right) \sqrt{B_{x2}} \right). \tag{495}$$

Furthermore, for any $\epsilon \geq 2(0, 1)$, set

$$\overline{b} = \left(\frac{B_G^{\frac{1}{4}} B_{a2} B_b^{\frac{3}{4}}}{r B_{a1}} \right), \quad m = \left(\frac{1}{p \epsilon^4} \left(r B_{a1} \sqrt{B_{x1} B_{G1}} \sqrt{\frac{B_b}{B_G}} \right)^4 + \frac{1}{\delta} + \frac{1}{p} \left(\log\left(\frac{r}{\delta}\right) \right)^2 \right),$$

$$\eta = \left(\frac{\epsilon}{\left(\frac{\sqrt{r} B_{a2} B_b B_{x1}}{(mp)^{\frac{1}{4}}} + m \overline{b} \right) \left(\frac{\sqrt{\log m B_{x1}} (mp)^{\frac{1}{4}}}{\sqrt{B_b B_G}} + 1 \right)} \right), \quad T = \left(\frac{1}{\eta B_{x1} (mp)^{\frac{1}{4}}} \right),$$

$$n = \left(\left(\frac{m B_x B_{a2}^2 \overline{B_b} (mp)^{\frac{1}{2}} \log m}{\epsilon r B_{a1} \overline{B_G}} \right)^3 + \left(\frac{B_x}{B_{x2}} + \log \frac{Tm}{p \delta} + \left(1 + \frac{1}{B_G} + \frac{1}{B_{G1}} \right) \frac{B_{x\infty}}{j_{\ell'}(0) j} \right)^3 \right),$$

we have there exists $t \geq 2 [T]$ with

$$\Pr[\text{sign}(f^{(t)})(\mathbf{x}) \neq y] \leq L_{\mathcal{D}}(f^{(t)}) \tag{496}$$

$$\text{OPT}_{d;r;B_F;S_{p_1};B_G;B_{G1}} + rB_{a1} \sqrt{2 \log(d)d} \left(\gamma_{\infty} + O\left(\frac{B_{x\infty}}{B_G j_{\ell'}(0) j n^{\frac{1}{3}}}\right) \right) + \epsilon. \tag{497}$$

Proof of Theorem E.37. Proof of the theorem and parameter choices remain the same as Theorem D.17 except for setting $B = \frac{\sqrt{B_{x1} B_{G1}}}{(mp)^{\frac{1}{4}}} \sqrt{\frac{B_b}{B_G}}$ and apply Lemma E.36. \square

E.6.2 Feature Learning of Uniform Parity Functions

We denote

$$g_{i;j} = \mathbb{E}_{(\mathbf{x};y)} \left[y \sigma' \left[\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right] \mathbf{x}_j \right] \tag{498}$$

$$\xi_k = \binom{k-1}{n-k} \frac{\binom{\frac{n-1}{2}}{\frac{k-1}{2}}}{\binom{n-1}{k}} 2^{-(n-1)} \binom{n-1}{\frac{n-1}{2}}. \tag{499}$$

Lemma E.38 (Uniform Parity Functions: Gradient Feature Learning. Corollary of Lemma 3 in [15]). *Assume that $n \geq 2(k+1)^2$. Then, the following holds:*

If $j \in A$, then

$$g_{i;j} = \xi_{k-1} \prod_{l \in A \setminus \{j\}} (\mathbf{w}_{i;l}^{(0)}). \tag{500}$$

If $i \notin A$, then

$$g_{i:j} = \xi_{k-1} \prod_{l \in A \cup \{j\}} (\mathbf{w}_{i:l}^{(0)}). \quad (501)$$

Lemma E.39 (Uniform Parity Functions: Existence of Good Networks (Alternative)). *Assume the same condition as in Lemma E.38. Define*

$$D = \frac{\sum_{l \in A} \mathbf{M}_l}{k \sum_{l \in A} \mathbf{M}_l k_2} \quad (502)$$

and

$$f^*(\mathbf{x}) = \sum_{i=0}^k \binom{k}{i} \rho_{\frac{k}{k}}^i \left[\sigma \left(\langle \mathbf{h}D, \mathbf{x} \rangle - \frac{2i}{\rho_{\frac{k}{k}}} \right) - 2\sigma \left(\langle \mathbf{h}D, \mathbf{x} \rangle - \frac{2i}{\rho_{\frac{k}{k}}} \right) + \sigma \left(\langle \mathbf{h}D, \mathbf{x} \rangle - \frac{2i}{\rho_{\frac{k}{k}}} \right) \right]. \quad (503)$$

For $D_{\text{parity-uniform}}$ setting, we have $f^* \in \mathcal{F}_{d;3(k+1);B_F;S_{p_1}^1;B_G;B_{G1}}$ where $B_F = (B_{a1}, B_{a2}, B_b) = (2^{\frac{\rho_{\frac{k}{k}}}{k}}, 2\sqrt{(k(k+1))}, \frac{k+1}{\sqrt{k}})$, $p = \frac{1}{2^{k-1}}$, $\gamma_\infty = O\left(\frac{\sqrt{k}}{d-k}\right)$, $B_G = (B_{G1}) = (d^{-k})$ and $B_{x1} = \frac{\rho_{\frac{k}{k}}}{d}$, $B_{x2} = d$. We also have $\text{OPT}_{d;3(k+1);B_F;S_{p_1}^1;B_G;B_{G1}} = 0$.

Proof of Lemma E.39. Fix index i , with probability $p_1 = \frac{1}{2^{k-1}}$, we will have $\mathbf{w}_{i:j}^{(0)} = \text{sign}(\mathbf{a}_i^{(0)}) \text{sign}(\xi_{k-1})$, for δ_j . For $\mathbf{w}_i^{(0)}$ that satisfy these conditions, we will have:

$$\text{sign}(\mathbf{a}_i^{(0)})g_{i:j} = j\xi_{k-1}j, \quad \delta_j \in A \quad (504)$$

$$\text{sign}(\mathbf{a}_i^{(0)})g_{i:j} = j\xi_{k+1}j, \quad \delta_j \notin A. \quad (505)$$

Then by Lemma 4 in [15], we have

$$\left\| \frac{\text{sign}(\mathbf{a}_i^{(0)})G(\mathbf{w}_i^{(0)}, \mathbf{b})}{kG(\mathbf{w}_i^{(0)}, \mathbf{b})k} - D \right\|_\infty \max \left\{ \left| \frac{1}{k\sqrt{\frac{1}{k} + \frac{1}{d-k}}} - \frac{1}{\rho_{\frac{k}{k}}} \right|, \left| \frac{1}{(d-k)\sqrt{\frac{1}{k} + \frac{1}{d-k}}} \right| \right\} \quad (506)$$

$$\frac{\rho_{\frac{k}{k}}}{d-k} \quad (507)$$

and

$$k\text{sign}(\mathbf{a}_i^{(0)})G(\mathbf{w}_i^{(0)}, \mathbf{b})k_2 = \sqrt{kj\xi_{k-1}j^2 + (d-k)j\xi_{k+1}j^2} = (d-k)^{(k)}. \quad (508)$$

From here, we can see that if we set $\gamma_\infty = \frac{\sqrt{k}}{d-k}$, $B_G = B_{G1} = \sqrt{kj\xi_{k-1}j^2 + (d-k)j\xi_{k+1}j^2}$, $p = p_1$, we will have $(D, +1), (D, -1) \in \mathcal{S}_{p_1;B_G;B_{G1}}^\infty$ by our symmetric initialization. As a result, we have $f^* \in \mathcal{F}_{d;3(k+1);B_F;S_{p_1}^1;B_G;B_{G1}}$. Finally, it is easy to verify that $f^*(\mathbf{x}) = \text{XOR}(\mathbf{x}_A)$, thus $\text{OPT}_{d;3(k+1);B_F;S_{p_1}^1;B_G;B_{G1}} = 0$. \square

Theorem E.40 (Uniform Parity Functions: Main Result (Alternative)). *For $D_{\text{parity-uniform}}$ setting, for any $\delta \in (0, 1)$ satisfying $\delta = O\left(\frac{1}{d^2}\right)$ and for any $\epsilon \in (0, 1)$ when*

$$m = \text{poly} \left(\log \left(\frac{1}{\delta} \right), \frac{1}{\epsilon}, 2^{(k)}, d \right), T = (d-k)^{(k)}, n = (d-k)^{(k)} \quad (509)$$

trained by Algorithm 1 with hinge loss, with probability at least $1 - \delta$ over the initialization, with proper hyper-parameters, there exists $t \geq [T]$ such that

$$\Pr[\text{sign}(f_{(t)}(\mathbf{x})) \neq y] \leq \frac{k^2 \sqrt{d \log(d)}}{d-k} + \epsilon. \quad (510)$$

Proof of Theorem E.40. Plug the values of parameters into Theorem E.37 and directly get the result. \square

E.7 Multiple Index Model with Low Degree Polynomial

E.7.1 Problem Setup

The multiple-index data problem has been used for studying network learning [18, 32]. We consider proving guarantees for the setting in [32], following our framework. We use the properties of the problem to prove the key lemma (i.e., the existence of good networks) in our framework and then derive the final guarantee from our theorem of the simple setting (Theorem 3.4).

Data Distributions. We draw input from the distribution $D_{\mathcal{X}} = N(0, I_{d \times d})$, and we assume the target function is $g^*(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}$, where g^* is a degree τ polynomial normalized so that $\mathbb{E}_{\mathbf{x} \sim D_{\mathcal{X}}} [g^*(\mathbf{x})^2] = 1$.

Assumption E.41. *There exists linearly independent vectors u_1, \dots, u_r such that $g^*(\mathbf{x}) = g(\langle \mathbf{x}, u_1 \rangle, \dots, \langle \mathbf{x}, u_r \rangle)$. $H := \mathbb{E}_{\mathbf{x} \sim D_{\mathcal{X}}} [r^2 g^*(\mathbf{x})]$ has rank r , where H is a Hessian matrix.*

Definition E.42. *Denote the normalized condition number of H by*

$$\kappa := \frac{kH^\dagger k}{r}. \quad (511)$$

Initialization and Loss. For $\delta_i \geq [m]$, we use the following initialization:

$$\mathbf{a}_i^{(0)} = \frac{1}{\sqrt{m}}, \mathbf{w}_i^{(0)} \sim N\left(0, \frac{1}{d} I_{d \times d}\right) \text{ and } \mathbf{b}_i = 0. \quad (512)$$

For this regression problem, we use mean square loss:

$$L_{D_{\mathcal{X}}}(f) = \mathbb{E}_{\mathbf{x} \sim D_{\mathcal{X}}} [(f(\mathbf{x}) - g^*(\mathbf{x}))^2]. \quad (513)$$

Training Process. We use the following one-step training algorithm for this specific data distribution.

Algorithm 5 Network Training via Gradient Descent [32]. Special case of Algorithm 2

Initialize $(\mathbf{a}^{(0)}, \mathbf{W}^{(0)}, \mathbf{b})$ as in Equation (8) and Equation (512); Sample $Z \sim D_{\mathcal{X}}^n$
 $\rho = \frac{1}{n} \sum_{\mathbf{x} \in Z} g^*(\mathbf{x}), \beta = \frac{1}{n} \sum_{\mathbf{x} \in Z} g^*(\mathbf{x})\mathbf{x}$
 $y = g^*(\mathbf{x}) - \rho - \beta \mathbf{x}$
 $\mathbf{W}^{(1)} = \mathbf{W}^{(0)} - \eta^{(1)} (r \mathbf{W}^{(0)} \tilde{L}_Z(f^{(0)}) + \lambda^{(1)} \mathbf{W}^{(0)})$
 Re-initialize $\mathbf{b}_i \sim N(0, 1)$
for $t = 2$ **to** T **do**
 $\mathbf{a}^{(t)} = \mathbf{a}^{(t-1)} - \eta^{(t)} r \mathbf{a} \tilde{L}_Z(f^{(t-1)})$
end for

Lemma E.43 (Multiple Index Model with Low Degree Polynomial: Existence of Good Networks. Rephrase of Lemma 25 in [32]). *Assume $n \geq d^2 r \kappa^2 (C_l \log(nmd))^{-1}$, $d \leq C_d \kappa r^{3=2}$, and $m \geq r \kappa^2 (C_l \log(nmd))^{6+1}$ for sufficiently large constants C_d, C_l , and let $\eta^{(1)} = \sqrt{\frac{d}{(C_l \log(nmd))^3}}$ and $\lambda^{(1)} = \frac{1}{(1)}$. Then with probability $1 - \frac{1}{\text{poly}(m;d)}$, there exists $\mathbf{a} \geq \mathbb{R}^m$ such that $f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}$ satisfies*

$$L_{D_{\mathcal{X}}}(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}) = O\left(\frac{1}{n} + \frac{r \kappa^2 (C_l \log(nmd))^{6+1}}{m}\right) \quad (514)$$

and

$$\|\mathbf{a}\|_2^2 = O\left(\frac{r \kappa^2 (C_l \log(nmd))^{6+1}}{m}\right). \quad (515)$$

E.7.2 Multiple Index Model: Final Guarantee

Considering training by Algorithm 5, we have the following results.

Theorem E.44 (Multiple Index Model with Low Degree Polynomial: Main Result). Assume $n \geq (d^2 r \kappa^2 (C_l \log(nmd))^{+1} + m)$, $d \leq C_d \kappa r^{3=2}$, and $m \leq (1-r \kappa^2 (C_l \log(nmd))^6 + 1)$ for sufficiently large constants C_d, C_l . Let $\eta^{(1)} = \sqrt{\frac{d}{(C_l \log(nmd))^3}}$ and $\lambda^{(1)} = \frac{1}{(1)}$, and $\eta = \eta^{(t)} = (m^{-1})$, for all $t \geq 2, 3, \dots, T$. For any $\epsilon \geq (0, 1)$, if $T \geq \left(\frac{m^2}{\epsilon}\right)$, then with properly set parameters and Algorithm 5, with high probability that there exists $t \geq [T]$ such that

$$L_{\mathcal{D}_X} f_{(\mathbf{a}^{(t)}, \mathbf{W}^{(1)}, \mathbf{b})} \leq \epsilon. \quad (516)$$

Proof of Theorem E.44. By Lemma E.43, we have for properly chosen hyper-parameters,

$$\text{OPT}_{\mathbf{W}^{(1)}, \mathbf{b}; B_{a_2}} L_{\mathcal{D}_X} (f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}) \leq O\left(\frac{1}{n} + \frac{r \kappa^2 (C_l \log(nmd))^6 + 1}{m}\right) \quad (517)$$

$$\leq \frac{\epsilon}{3}. \quad (518)$$

We compute the L -smooth constant of $\tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})})$ to \mathbf{a} .

$$\left\| r_{\mathbf{a}} \tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}_1, \mathbf{W}^{(1)}, \mathbf{b})}) - r_{\mathbf{a}} \tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}_2, \mathbf{W}^{(1)}, \mathbf{b})}) \right\|_2 \quad (519)$$

$$= \left\| \frac{1}{n} \sum_{\mathbf{x} \in \mathcal{Z}} \left[2 \left(f_{(\mathbf{a}_1, \mathbf{W}^{(1)}, \mathbf{b})}(\mathbf{x}) - g^* - f_{(\mathbf{a}_2, \mathbf{W}^{(1)}, \mathbf{b})}(\mathbf{x}) + g^* \right) \sigma(\mathbf{W}^{(1)\top} \mathbf{x} - \mathbf{b}) \right] \right\|_2 \quad (520)$$

$$\left\| \frac{1}{n} \sum_{\mathbf{x} \in \mathcal{Z}} \left[2 \left(\mathbf{a}_1^\top \sigma(\mathbf{W}^{(1)\top} \mathbf{x} - \mathbf{b}) - \mathbf{a}_2^\top \sigma(\mathbf{W}^{(1)\top} \mathbf{x} - \mathbf{b}) \right) \sigma(\mathbf{W}^{(1)\top} \mathbf{x} - \mathbf{b}) \right] \right\|_2 \quad (521)$$

$$\frac{1}{n} \sum_{\mathbf{x} \in \mathcal{Z}} \left[2 k \mathbf{a}_1 - \mathbf{a}_2 k_2 \left\| \sigma(\mathbf{W}^{(1)\top} \mathbf{x} - \mathbf{b}) \right\|_2^2 \right]. \quad (522)$$

By the proof of Lemma 25 in [32], we have for $\delta_i \geq [4m]$, with probability at least $1 - \frac{1}{\text{poly}(m, d)}$, $\| \mathbf{W}^{(1)\top} \mathbf{x} - \mathbf{b} \|_2 \leq 1$, with some large polynomial $\text{poly}(m, d)$. As a result, we have

$$\frac{1}{n} \sum_{\mathbf{x} \in \mathcal{Z}} \left\| \mathbf{W}^{(1)\top} \mathbf{x} - \mathbf{b} \right\|_2^2 \leq m + \frac{1}{\text{poly}(m, d)} = O(m). \quad (523)$$

Thus, we have,

$$L = O\left(\frac{1}{n} \sum_{\mathbf{x} \in \mathcal{Z}} \left\| \sigma(\mathbf{W}^{(1)\top} \mathbf{x} - \mathbf{b}) \right\|_2^2\right) \quad (524)$$

$$O\left(\frac{1}{n} \sum_{\mathbf{x} \in \mathcal{Z}} \left\| \mathbf{W}^{(1)\top} \mathbf{x} - \mathbf{b} \right\|_2^2\right) \quad (525)$$

$$O\left(\frac{1}{n} \sum_{\mathbf{x} \in \mathcal{Z}} \left\| \mathbf{W}^{(1)\top} \mathbf{x} \right\|_2^2 + k \mathbf{b} k_2^2\right) \quad (526)$$

$$O(m). \quad (527)$$

This means that we can let $\eta = (m^{-1})$ and we will get our convergence result. We can bound $k \mathbf{a}^{(1)} k_2$ and $k \mathbf{a} k_2$ by $k \mathbf{a}^{(1)} k_2 = O\left(\frac{r \kappa^2 (C_l \log(nmd))^6}{m}\right)$ and $k \mathbf{a} k_2 = O(\epsilon)$. So, if we choose $T \geq \left(\frac{m}{\epsilon}\right)$, there exists $t \geq [T]$ such that $\tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}^{(t)}, \mathbf{W}^{(1)}, \mathbf{b})}) - \tilde{L}_{\mathcal{Z}}(f_{(\mathbf{a}, \mathbf{W}^{(1)}, \mathbf{b})}) \leq O\left(\frac{L \|\mathbf{a}^{(1)} - \mathbf{a}\|_2^2}{T}\right) \leq \epsilon/3$.

We also have $\sqrt{\frac{\|\mathbf{a}\|_2^2 (\|\mathbf{W}^{(1)}\|_F^2 B_X^2 + \|\mathbf{b}\|_2^2)}{n}} \leq \frac{\epsilon}{3}$. Then our theorem gets proved by Theorem 3.4. \square

Discussion. We would like to unify [32], which are very closely related to our framework: their analysis for multiple index data follows the same principle and analysis approach as our general framework, although it does not completely fit into our Theorem 3.12 due to some technical differences. We can cover it with our Theorem 3.4.

Our work and [32] share the same principle and analysis approach. [32] shows that the first layer learns good features by one gradient step update, which can approximate the true labels by a low-degree polynomial function. Then, a classifier (the second layer) is trained on top of the learned first layer which leads to the final guarantees. This is consistent with our framework: we first show that the first layer learns good features by one gradient step update, which can approximate the true labels, and then show a good classifier can be learned on the first layer.

Our work and [32] have technical differences. First, in the second stage, [32] fix the first layer and only update the top layer which is a convex optimization. Our framework allows updates in the first layer and uses online convex learning techniques for the analysis. Second, they consider the square loss (this is used to calculate Hermite coefficients explicitly for gradients, which are useful in the low-degree polynomial function approximation). While in our online convex learning analysis, we need boundedness of the derivative of the loss to show that the first layer weights' changes are bounded in the second stage. Given the above two technicalities, we analyze their training algorithm (Algorithm 2) which fixes the first layer weights and fits into our Theorem 3.4.

F Auxiliary Lemmas

In this section, we present some Lemmas used frequently.

Lemma F.1 (Lemmas on Gradients).

$$r_{\mathbf{w}} L_{(\mathbf{x};y)}(f) = \left[\frac{\partial L_{(\mathbf{x};y)}(f)}{\partial \mathbf{w}_1}, \dots, \frac{\partial L_{(\mathbf{x};y)}(f)}{\partial \mathbf{w}_i}, \dots, \frac{\partial L_{(\mathbf{x};y)}(f)}{\partial \mathbf{w}_{4m}} \right], \quad (528)$$

$$\frac{\partial L_{(\mathbf{x};y)}(f)}{\partial \mathbf{w}_i} = \mathbf{a}_i \ell'(yf(\mathbf{x})) y [\sigma'(h\mathbf{w}_i, \mathbf{x} - \mathbf{b}_i)] \mathbf{x}, \quad (529)$$

$$r_{\mathbf{w}} L_{\mathcal{D}}(f) = \left[\frac{\partial L_{\mathcal{D}}(f)}{\partial \mathbf{w}_1}, \dots, \frac{\partial L_{\mathcal{D}}(f)}{\partial \mathbf{w}_i}, \dots, \frac{\partial L_{\mathcal{D}}(f)}{\partial \mathbf{w}_{4m}} \right], \quad (530)$$

$$\frac{\partial L_{\mathcal{D}}(f)}{\partial \mathbf{w}_i} = \mathbf{a}_i E_{(\mathbf{x};y)} [\ell'(yf(\mathbf{x})) y [\sigma'(h\mathbf{w}_i, \mathbf{x} - \mathbf{b}_i)] \mathbf{x}], \quad (531)$$

$$\frac{\partial L_{\mathcal{D}}(f)}{\partial \mathbf{a}_i} = E_{(\mathbf{x};y)} [\ell'(yf(\mathbf{x})) y [\sigma'(h\mathbf{w}_i, \mathbf{x} - \mathbf{b}_i)]]. \quad (532)$$

Proof. These can be verified by direct calculation. \square

Lemma F.2 (Property of Symmetric Initialization). *For any $\mathbf{x} \in \mathbb{R}^d$, we have $f_{(0)}(\mathbf{x}) = 0$. For all $i \in [2m]$, we have $\mathbf{w}_i^{(1)} = \mathbf{w}_{i+2m}^{(1)}$. When input data is symmetric, i.e., $E_{(\mathbf{x};y)}[y\mathbf{x}] = \mathbf{0}$, for all $i \in [m]$, we have $\mathbf{w}_i^{(1)} = \mathbf{w}_{i+m}^{(1)}$.*

Proof of Lemma F.2. By symmetric initialization, we have $f_{(0)}(\mathbf{x}) = 0$. For all $i \in [2m]$, we have

$$\mathbf{w}_i^{(1)} = \eta^{(1)} \ell'(0) \mathbf{a}_i^{(0)} E_{(\mathbf{x};y)} \left[y \sigma' \left[\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_i \right] \mathbf{x} \right] \quad (533)$$

$$= \eta^{(1)} \ell'(0) \mathbf{a}_{i+2m}^{(0)} E_{(\mathbf{x};y)} \left[y \sigma' \left[\left\langle \mathbf{w}_{i+2m}^{(0)}, \mathbf{x} \right\rangle - \mathbf{b}_{i+2m} \right] \mathbf{x} \right] \quad (534)$$

$$= \mathbf{w}_{i+2m}^{(1)}. \quad (535)$$

When $\mathbb{E}_{(\mathbf{x},y)}[y\mathbf{x}] = \mathbf{0}$, for all $i \geq [m]$, we have

$$\mathbf{w}_i^{(1)} = \eta^{(1)} \ell'(0) \mathbf{a}_i^{(0)} \mathbb{E}_{(\mathbf{x},y)} \left[y \sigma' \left[\left\langle \mathbf{w}_i^{(0)}, \mathbf{x} \right\rangle + \mathbf{b}_i \right] \mathbf{x} \right] \quad (536)$$

$$= \eta^{(1)} \ell'(0) \mathbf{a}_{i+m}^{(0)} \mathbb{E}_{(\mathbf{x},y)} \left[y \sigma' \left[\left\langle \mathbf{w}_{i+m}^{(0)}, \mathbf{x} \right\rangle + \mathbf{b}_{i+m} \right] \mathbf{x} \right] \quad (537)$$

$$= \eta^{(1)} \ell'(0) \mathbf{a}_{i+m}^{(0)} \mathbb{E}_{(\mathbf{x},y)} \left[y \sigma' \left[\left\langle \mathbf{w}_{i+m}^{(0)}, \mathbf{x} \right\rangle + \mathbf{b}_{i+m} \right] \mathbf{x} - y \mathbf{x} \right] \quad (538)$$

$$= \eta^{(1)} \ell'(0) \mathbf{a}_{i+m}^{(0)} \mathbb{E}_{(\mathbf{x},y)} \left[y \sigma' \left[\left\langle \mathbf{w}_{i+m}^{(0)}, \mathbf{x} \right\rangle + \mathbf{b}_{i+m} \right] \mathbf{x} \right] \quad (539)$$

$$= \mathbf{w}_{i+m}^{(1)}. \quad (540)$$

□

Lemma F.3 (Property of Direction Neighborhood). *If $\mathbf{w} \in \mathcal{C}_{D;}$, we have $\rho \mathbf{w} \in \mathcal{C}_{D;}$ for any $\rho \neq 0$. We also have $\mathbf{0} \notin \mathcal{C}_{D;}$. Also, if $(D, s) \in \mathcal{S}_{p; ; B_G}$, we have $(\rho D, s) \in \mathcal{S}_{p; ; B_G}$.*

Proof. These can be verified by direct calculation. □

Lemma F.4 (Maximum Gaussian Tail Bound). *M_n is the maximum of n i.i.d. standard normal Gaussian. Then*

$$\Pr \left(M_n \geq \sqrt{2 \log n} + \frac{z}{\sqrt{2 \log n}} \right) \leq e^{-z}. \quad (541)$$

Proof. These can be verified by direct calculation. □

Lemma F.5 (Chi-squared Tail Bound). *If X is a $\chi^2(k)$ random variable. Then, $z \in \mathbb{R}$, we have*

$$\Pr(X \geq k + 2\sqrt{kz} + 2z) \leq e^{-z}. \quad (542)$$

Proof. These can be verified by direct calculation. □

Lemma F.6 (Gaussian Tail Bound). *If g is standard Gaussian and $z > 0$, we have*

$$\frac{1}{\sqrt{2\pi}} \frac{z}{z^2 + 1} e^{-z^2/2} < \Pr_{g \sim \mathcal{N}(0,1)}[g > z] < \frac{1}{\sqrt{2\pi}} \frac{1}{z} e^{-z^2/2}. \quad (543)$$

Proof. These can be verified by direct calculation. □

Lemma F.7 (Gaussian Tail Expectation Bound). *If g is standard Gaussian and $z \in \mathbb{R}$, we have*

$$z \mathbb{E}_{g \sim \mathcal{N}(0,1)}[|g| \mathbb{1}[g > z]] < 2 \Pr_{g \sim \mathcal{N}(0,1)}[g > z]^{0.9}. \quad (544)$$

Proof of Lemma F.7. For any $p \in (0, 1)$, we have

$$\left| \int_{-\infty}^{\sqrt{2} \text{erf}^{-1}(2p-1)} \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx \right| < 2p^{0.9}, \quad (545)$$

where $\sqrt{2} \text{erf}^{-1}(2p-1)$ is the quantile function of the standard Gaussian. We finish the proof by replacing p to be $\Pr_{g \sim \mathcal{N}(0,1)}[g > z]$. □

Lemma F.8. *If a function g satisfy $h(n+2) = 2h(n+1) - (1-\rho^2)h(n) + \beta$ for $n \in \mathbb{N}_+$ where $\rho, \beta > 0$, then $h(n) = \frac{\beta}{1-\rho^2} + c_1(1-\rho)^n + c_2(1+\rho)^n$, where c_1, c_2 only depends on $h(1)$ and $h(2)$.*

Proof. These can be verified by direct calculation. □

Lemma F.9 (Rademacher Complexity Bounds. Rephrase of Lemma 48 in [32]). *For fixed \mathbf{W}, \mathbf{b} , let $F = f_{(\mathbf{a}, \mathbf{W}; \mathbf{b})} : \mathcal{K} \times \mathcal{K} \rightarrow \mathbb{R}$. Then,*

$$\mathfrak{R}(F) \leq \sqrt{\frac{B_{a2}^2 (k \mathbf{W} k_F^2 B_x^2 + k \mathbf{b} k_2^2)}{n}}. \quad (546)$$