

# How Does the ReLU Activation Affect the Implicit Bias of Gradient Descent on High-Dimensional Neural Network Regression?

**Kuo-Wei Lai\***

*School of Electrical and Computer Engineering  
Georgia Institute of Technology, Atlanta, USA*

KLAI36@GATECH.EDU

**Guanghui Wang\***

*College of Computing  
Georgia Institute of Technology, Atlanta, USA*

GWANG369@GATECH.EDU

**Molei Tao**

*School of Mathematics  
Georgia Institute of Technology, Atlanta, USA*

MTAO@GATECH.EDU

**Vidya Muthukumar**

*School of Electrical and Computer Engineering  
H. Milton Stewart School of Industrial and Systems Engineering  
Georgia Institute of Technology, Atlanta, USA*

VMUTHUKUMAR8@GATECH.EDU

**Editors:** Steve Hanneke and Tor Lattimore

## Abstract

Overparameterized ML models, including neural networks, typically induce underdetermined training objectives with multiple global minima. The *implicit bias* refers to the limiting global minimum that is attained by a common optimization algorithm, such as gradient descent (GD). In this paper, we characterize the implicit bias of GD for training a shallow ReLU model with the squared loss on high-dimensional random features. Prior work (Vardi and Shamir, 2021) showed that the implicit bias does not exist in the worst-case, or corresponds exactly to the minimum- $\ell_2$ -norm interpolating solution under *exactly* orthogonal data (Boursier et al., 2022). Our work interpolates between these two extremes and shows that, for sufficiently high-dimensional random data, the implicit bias approximates the minimum- $\ell_2$ -norm solution with high probability with a gap on the order  $\Theta\left(\sqrt{n/\|\lambda\|_1}\right)$ , where  $n$  is the number of training examples and  $\lambda$  denotes the spectrum of the data covariance matrix. Our results are obtained through a novel primal-dual analysis that carefully tracks the evolution of predictions, data-span coefficients, as well as their interactions, and show that the ReLU activation pattern quickly stabilizes with high probability over random data.

**Keywords:** Implicit Bias, Gradient Descent, ReLU, Squared Loss, Regression, Primal-Dual

## 1. Introduction

In many modern machine learning problems, the training objectives are typically *underdetermined*, which implies that they may admit (potentially infinitely) many global minima. Despite this, a large body of empirical results (Neyshabur et al., 2014; Zhang et al., 2021) show that optimization algorithms such as gradient descent frequently converge to solutions that generalize well, even in the absence of any explicit regularization. This phenomenon is commonly referred to as the *implicit*

---

\* Equal contribution; co-first author.

*bias* introduced by gradient descent (Ji and Telgarsky, 2018; Soudry et al., 2018), and understanding the nature of this benign bias has become a central topic of recent research.

The study of the implicit bias of gradient descent originally emerged in the context of linear models. For linear classification with separable data, the seminal work of Soudry et al. (2018); Ji and Telgarsky (2018) shows that, when minimizing exponentially-tailed losses, gradient descent converges in direction to the max-margin solution that minimizes  $\ell_2$ -norm. For linear regression with the squared loss, gradient descent is known to converge to the zero-loss (interpolating) solution with the minimum- $\ell_2$ -norm (Engl et al., 1996). Building on these foundational results, several finer characterizations were derived for linear models, including sharper convergence analyses (Ji and Telgarsky, 2018; Nacson et al., 2019; Ji and Telgarsky, 2021), general classes of first-order methods (Gunasekar et al., 2018; Sun et al., 2022; Wang et al., 2025), and a deeper understanding of the implicit bias on high-dimensional data (Hsu et al., 2021; Lai and Muthukumar, 2025).

Understanding the implicit bias in nonlinear models such as neural networks remains a significant challenge, primarily due to the induced non-convexity of the optimization objective. In this work, we focus on regression with a one-hidden-layer ReLU neural network and the squared loss, which represents one of the most fundamental and natural extensions beyond linearity. Remarkably, Vardi and Shamir (2021) showed that establishing the implicit bias of ReLU models is known to be hard in the worst case, even when the model consists only of a single neuron and assuming global convergence. (On the other hand, the implicit bias of a single neuron with a strictly monotonic activation function (e.g., leaky ReLU) does follow the minimum- $\ell_2$ -norm solution.) They do this through a stylized counterexample of 3 data points with 3-dimensional features, raising the natural question of whether the implicit bias becomes characterizable under typical data ensembles. At the other extreme, Boursier et al. (2022) showed that the implicit bias of gradient flow for *exactly* orthogonal features is exactly the minimum- $\ell_2$ -norm solution. However, an exact orthogonality assumption is restrictive and rarely holds in practice. Notably, high-dimensional random features are *near*-orthogonal, raising the question of whether the implicit bias can be characterized in this more realistic but also more challenging case.

**Our contributions:** In this paper, we establish new insights into the implicit bias induced by gradient descent for ReLU networks trained with the squared loss on sufficiently high-dimensional data. Our main contributions are summarized as follows. First, we completely characterize the expression for the implicit bias of gradient descent dynamics on ReLU models with 1 or 2 neurons for high-dimensional data under sufficient conditions (Theorems 4 and 8). Second, we quantify the relationship between the implicit bias of gradient descent and the global minimum that achieves the minimum- $\ell_2$ -norm. More specifically, we establish both upper and lower bounds on the distance between the gradient descent limit and the minimum- $\ell_2$ -norm solution, showing that it scales as  $\Theta(\sqrt{n}/\|\lambda\|_1)$  where  $n$  is the number of training examples and  $\lambda$  denotes the spectrum of the data covariance matrix (Theorems 6 and 9). Consequently, the solutions are very close, but not identical, for high-dimensional features. Interestingly, a similar phenomenon was also shown to occur with exponentially-tailed losses (Frei et al., 2023a,b) for classification.

**Our techniques in a nutshell:** Our main results are obtained through a novel primal-dual formulation of the gradient descent dynamics under the squared loss with ReLU networks, which is inspired by mirror descent (first studied by Ji and Telgarsky (2021) for linear models). Instead of directly tracking the weight vector in the original parameter space like previous work, we introduce primal variables representing the predictions on training examples, and dual variables capturing

the coefficients in the data span. This representation is particularly well-suited for analyzing ReLU networks because the sign of each primal variable directly determines whether the corresponding example is active, and hence whether its dual variable receives a gradient update. Our analysis reveals that understanding the gradient dynamics hinges on tracking (i) the positivity of the primal variables and (ii) the interactions between training examples. We introduce new tools to carefully control the evolution of positive primal variables and sufficiently negative dual variables (Lemmas 10 and 11, which may be of independent interest). Underlying the proofs of our approximation results to the minimum- $\ell_2$ -norm solutions are novel characterizations of the latter as minimum- $\ell_2$ -norm *linear* interpolations of a (possibly data-dependent) subset of training examples. This data-dependent subset selection is a fundamental difference between the implicit bias of linear models and ReLU models.

### 1.1. Related Work

We now briefly discuss the most closely related work and highlight key differences of our approach. We contextualize our results within the most closely related prior studies on implicit bias of regression models in Table 1. Boursier et al. (2022) study the dynamics of gradient flow on two-layer ReLU networks under an *exact* orthogonality assumption on the data. Exact orthogonality removes interactions between examples and significantly simplifies the activation patterns induced by the ReLU nonlinearity. As a result, their analysis primarily focuses on how the second-layer weights evolve to fit all examples, leading to a multi-phase gradient flow dynamic. Under these assumptions, they show that gradient flow converges to the minimum- $\ell_2$ -norm solution (their Theorem 1). In contrast, our work focuses on understanding how interactions between examples, captured through the Gram matrix, shape the active and inactive patterns in ReLU models under more realistic, controllable high-dimensional settings. Interestingly, we show that in high dimensions, the implicit bias is no longer exactly the minimum- $\ell_2$ -norm solution but remains close to it (Theorems 6 and 9). In comparison, Vardi and Shamir (2021) provide only a multiplicative upper bound on the magnitude of the norm of implicit bias in the worst-case data setting, showing that it is at most *twice* that of the minimum-norm solution. Dana et al. (2025) also analyze the high-dimensional regime and establish global convergence by showing that each example can be fitted by at least one neuron with high probability and all active examples stay active (their Theorem 1). However, their analysis does not address the behavior of inactive examples suppressed by the ReLU nonlinearity and does not shed light on the implicit bias. As a result, their work provides only a partial view of the gradient dynamics. In contrast, we introduce a novel primal-dual framework that allows us to simultaneously track both active and inactive examples (Lemmas 10 and 11). This framework enables a full characterization of the gradient dynamics and, consequently, the implicit bias in high dimensions. We use some of the observations of Dana et al. (2025) as a starting point for our primal-dual characterizations. More generally, most existing analyses (Vardi and Shamir, 2021; Boursier et al., 2022; Dana et al., 2025) rely on gradient flow and continuous-time ODE techniques, which assume infinitesimal step sizes. In contrast, our analysis directly studies gradient descent with finite (though still small) step sizes. This distinction is both theoretically and practically important, as gradient descent is the algorithm used in practice. Our primal-dual approach provides a new framework for analyzing discrete-time optimization dynamics in ReLU networks and opens a complementary direction to existing studies based on gradient flow.

Frei et al. (2023a,b) consider classification in a similar one-hidden-layer setup with the leaky ReLU activation, and also exploit simplified KKT conditions under nearly orthogonal data. How-

	Orthogonal data	High dimensional data	Worst-case data
ReLU models ( $k = 1$ ) $h(\mathbf{x}) := \sum_{k=1}^n s_k \sigma(\mathbf{w}_k^\top \mathbf{x})$	Implicit bias characterization (Boursier et al., 2022) $\mathbf{w}^{(\infty)} = \arg \min_{\mathbf{w} \in \{\mathbf{w}: \mathbf{X}\mathbf{w}=\mathbf{y}\}} \ \mathbf{w} - \mathbf{w}^{(0)}\ _2$	Global convergence only (Dana et al., 2025)	No implicit bias in general (Vardi and Shamir, 2021) $\ \mathbf{w}^{(\infty)} - \mathbf{w}^{(0)}\ _2 \leq 2 \cdot \ \mathbf{w}^* - \mathbf{w}^{(0)}\ _2$
		<b>This work</b> $\ \mathbf{w}^{(\infty)} - \mathbf{w}^*\ _2 \asymp \sqrt{\frac{n}{\ \boldsymbol{\lambda}\ _1}}$	
Linear models $h(\mathbf{x}) := \mathbf{w}^\top \mathbf{x}$	Implicit bias coincides with maximum $\ell_2$ margin SVM (Hsu et al., 2021)		$\mathbf{w}^{(\infty)} = \arg \min_{\mathbf{w} \in \{\mathbf{w}: \mathbf{X}\mathbf{w}=\mathbf{y}\}} \ \mathbf{w} - \mathbf{w}^{(0)}\ _2$ (Engl et al., 1996)

Table 1: Our results contextualized with related literature.

ever, due to the different training objectives underlying classification and regression, the resulting analyses are fundamentally different.

**Notation:** We use lowercase boldface letters (e.g.  $\mathbf{x}$ ) to denote vectors, lowercase letters (e.g.  $y$ ) to denote scalars, and uppercase boldface letters (e.g.  $\mathbf{X}$ ) to denote matrices. We use  $\|\cdot\|_p$  to denote the  $\ell_p$ -norm of a vector for  $p \in [1, \infty)$  and  $\|\cdot\|_2$  to additionally denote the operator norm of a matrix. For a vector  $\mathbf{x} \in \mathbb{R}^d$ , we use  $x_i$  to denote its  $i$ -th component. We use  $[n]$  to denote the set  $\{1, \dots, n\}$ . For a matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , a vector  $\mathbf{y} \in \mathbb{R}^n$ , and any index set  $S \subseteq [n]$ , we use  $\mathbf{X}_S \in \mathbb{R}^{|S| \times d}$  to denote the submatrix of  $\mathbf{X}$  consisting of the rows indexed by  $S$ , and  $\mathbf{y}_S \in \mathbb{R}^{|S|}$  denotes the corresponding subvector. We use  $\mathbf{x} \preceq \mathbf{0}$  (respectively  $\mathbf{x} \succeq \mathbf{0}$ ) to denote that every entry of vector  $\mathbf{x}$  is less than or equal to (respectively greater than or equal to) zero. We use  $C, c$  to denote universal constants that appear in upper and lower bounds, respectively, that may change from line to line. We also use the notation  $C_{(\cdot)}$  to denote universal constants with a specific meaning that *do not* change from line to line. We specifically choose  $C_0 \gtrsim C_\alpha^2$  and  $C_\alpha \gtrsim \max\{C_g^2, C_y C_g\}$  in our analysis.

## 2. Problem Setup

We consider a regression problem with feature vector  $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d$  and label  $y \in \mathcal{Y} \subset \mathbb{R}$ . We consider random feature vectors drawn according to a distribution  $\mathcal{D}$  with zero mean, i.e.,  $\mathbb{E}[\mathbf{x}] = \mathbf{0}$ , and covariance matrix  $\boldsymbol{\Sigma} = \mathbb{E}[\mathbf{x}\mathbf{x}^\top]$ . Let  $\boldsymbol{\Sigma} = \mathbf{V}\boldsymbol{\Lambda}\mathbf{V}^\top \in \mathbb{R}^{d \times d}$  be the eigendecomposition of the covariance matrix, where  $\mathbf{V} \in \mathbb{R}^{d \times d}$  is the matrix of eigenvectors and  $\boldsymbol{\Lambda} \in \mathbb{R}^{d \times d}$  is a diagonal matrix whose entries are the eigenvalues of  $\boldsymbol{\Sigma}$ , arranged in descending order. We make the mild assumption that the feature vector admits the representation  $\mathbf{x} = \mathbf{V}\boldsymbol{\Lambda}^{\frac{1}{2}}\mathbf{z}$  where  $\mathbf{z} \in \mathbb{R}^d$  has independent, mean-zero,  $\sigma_z^2$ -subgaussian components. By the definition of a  $\sigma_z^2$ -subgaussian random variable, each coordinate  $z_j$  satisfies  $\mathbb{E}[\exp(\mathbf{u}^\top \mathbf{z})] \leq \exp(\sigma_z^2 \|\mathbf{u}\|_2^2 / 2)$  for any  $\mathbf{u} \in \mathbb{R}^d$ . For simplicity, we set  $\sigma_z = 1$  throughout the remainder of the analysis. The labels  $y$  are only required to be bounded (see Assumption 1) and may be chosen arbitrarily. In particular,  $y$  does not need to satisfy any particular relationship with respect to  $\mathbf{x}$ .

We observe a dataset  $\{\mathbf{x}_i, y_i\}_{i=1}^n$  where the features  $\{\mathbf{x}_i\}_{i=1}^n$  are drawn i.i.d. from the distribution  $\mathcal{D}$ . We denote the data matrix by  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and the label vector by  $\mathbf{y} \in \mathbb{R}^n$ . Since we operate in a high-dimensional regime ( $d > n$ ), we make the mild assumption that  $\mathbf{X}$  has full row rank, i.e.,  $\text{rank}(\mathbf{X}) = n$ , which is automatically satisfied under the assumptions of all lemmas and theorems in this paper.<sup>1</sup>

<sup>1</sup>The full-rank assumption holds either almost surely or with high probability under random and high-dimensional features; see, e.g. Hsu et al. (2021).

**Assumption 0 (Full-rank Data Matrix)** *The data matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  satisfies  $\text{rank}(\mathbf{X}) = n$ .*

For ease of subsequent notation, we consider without loss of generality the samples with positive labels to appear in the upper block of the data matrix  $\mathbf{X}$ , while those with negative labels appear in the lower block. Let  $n_+$  denote the number of positive labels and  $n_- = n - n_+$  denote the number of negative labels. Accordingly, we write  $\mathbf{X} = [\mathbf{X}_+^\top \ \mathbf{X}_-^\top]^\top$  where  $\mathbf{X}_+ \in \mathbb{R}^{n_+ \times d}$  contains the features corresponding to positive labels and  $\mathbf{X}_- \in \mathbb{R}^{n_- \times d}$  contains the features corresponding to negative labels. We similarly partition the label vector as  $\mathbf{y} = [\mathbf{y}_+^\top \ \mathbf{y}_-^\top]^\top$ .

Next, we introduce our key assumptions on the features and labels. First, we assume that the magnitudes of the labels are bounded away from zero and infinity.

**Assumption 1 (Bounded Labels)** *For all  $i \in [n]$ , the labels satisfy  $y_{\min} \leq |y_i| \leq y_{\max}$  for some  $y_{\min}, y_{\max} \in \mathbb{R}_+$ .*

This assumption ensures that all labels are non-degenerate and have comparable scales, which will be important for controlling the dynamics of gradient-based optimization.

We next impose a high-dimensional assumption on the data features. To characterize the effective dimensionality of the feature distribution, we define two notions of effective dimension based on the spectrum  $\boldsymbol{\lambda} := [\lambda_1, \dots, \lambda_d]^\top$  of the feature covariance matrix  $\boldsymbol{\Sigma}$ , given by  $d_2 := \frac{\|\boldsymbol{\lambda}\|_1^2}{\|\boldsymbol{\lambda}\|_2^2}$ ,  $d_\infty := \frac{\|\boldsymbol{\lambda}\|_1}{\|\boldsymbol{\lambda}\|_\infty}$ . Note that when the covariance is isotropic, i.e.,  $\lambda_1 = \lambda_2 = \dots = \lambda_d = 1$ , we have  $d_2 = d_\infty = d$ , i.e., these reduce to the original data dimension. Our high-dimensional assumption requires that these effective dimensions dominate problem-dependent quantities involving the sample size  $n$  and the range of label magnitudes  $[y_{\min}, y_{\max}]$ . Similar conditions have also appeared in related global convergence analysis under the squared loss (Dana et al., 2025) and implicit bias analyses under the logistic/exponentially-tailed losses (Frei et al., 2023a).

**Assumption 2 (High-dimensional Features)** *The data features satisfy  $d_2 \geq C_0^2 \frac{n^2 y_{\max}^2}{y_{\min}^2}$  and  $d_\infty \geq C_0 \frac{n^{1.5} y_{\max}}{y_{\min}}$  for a sufficiently large constant  $C_0 > 1$ .*

This assumption places the problem in a sufficiently high-dimensional regime, ensuring strong concentration properties of the empirical Gram matrix and enabling precise control of the gradient dynamics analyzed in the following sections. We note that our techniques would yield similar guarantees on the implicit bias for deterministic feature vectors that satisfy a near-orthogonality condition adapted from Frei et al. (2023a,b) to real-valued labels.

**Assumption 3 (Deterministic Nearly-orthogonal Features)** *For a sufficiently large constant  $C_0 > 1$ , the data features satisfy  $\min_{i \in [n]} \|\mathbf{x}_i\|_2^2 \geq C_0 n \frac{y_{\max}}{y_{\min}} \max_{i \neq j} |\mathbf{x}_i^\top \mathbf{x}_j|$ .*

This condition ensures that the desired concentration properties of the empirical Gram matrix  $\mathbf{X} \mathbf{X}^\top$  hold for our analysis.

**General ReLU Models and Empirical Risk Minimization.** We denote by  $h_{\boldsymbol{\Theta}} : \mathcal{X} \rightarrow \mathcal{Y}$  the general ReLU model used for the regression task in this work, defined as  $h_{\boldsymbol{\Theta}}(\mathbf{x}) := \sum_{k=1}^m s_k \sigma(\mathbf{w}_k^\top \mathbf{x})$ , where  $\boldsymbol{\Theta}$  denotes the collection of model parameters  $\{\mathbf{w}_k\}_{k=1}^m$  together with fixed signs  $\{s_k\}_{k=1}^m$ . Here,  $s_k \in \{-1, +1\}$  denotes the sign of the  $k$ -th ReLU neuron,  $\sigma(z) := \max\{z, 0\}$  is the ReLU

activation function,  $\mathbf{w}_k \in \mathbb{R}^d$  is its weight vector, and  $m \geq 1$  is the number of neurons. To learn the regression model, we minimize the empirical risk under the squared loss, defined as

$$\mathcal{R}(\Theta) = \frac{1}{2} \sum_{i=1}^n (h_{\Theta}(\mathbf{x}_i) - y_i)^2 = \frac{1}{2} \|h_{\Theta}(\mathbf{X}) - \mathbf{y}\|_2^2, \quad (1)$$

where we define the vector-valued extension  $h_{\Theta}$  as  $h_{\Theta}(\mathbf{X}) := [h_{\Theta}(\mathbf{x}_1), \dots, h_{\Theta}(\mathbf{x}_n)]^{\top} \in \mathbb{R}^n$ . We employ the gradient descent algorithm to minimize (1). To make the dynamics more tractable, we only update the neuron weights  $\{\mathbf{w}_k\}_{k=1}^m$  and fix the signs of the neurons  $\{s_k\}_{k=1}^m$ <sup>2</sup>. When there are  $m > 1$  neurons, we will initialize at least one neuron for a positive sign and one neuron for a negative sign to ensure that the neural network can fit arbitrary labels.

**Gradient Descent and Primal-dual Representation.** For the ReLU model  $h_{\Theta}$ , the gradient of the empirical risk in (1) with respect to  $\mathbf{w}_k$  is given by

$$\nabla_{\mathbf{w}_k} \mathcal{R}(\Theta) = \sum_{i=1}^n (h_{\Theta}(\mathbf{x}_i) - y_i) s_k \cdot \mathbb{1}_{\mathbf{w}_k^{\top} \mathbf{x}_i > 0} \cdot \mathbf{x}_i = s_k \mathbf{X}^{\top} \mathbf{D}(\mathbf{X} \mathbf{w}_k) (h_{\Theta}(\mathbf{X}) - \mathbf{y}),$$

where  $\mathbf{D}(z) : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times n}$  denotes the diagonal matrix with entries  $D_{ii} := \mathbb{1}_{z_i > 0}$ . Accordingly, the gradient descent update for  $\mathbf{w}_k$  takes the form

$$\mathbf{w}_k^{(t+1)} = \mathbf{w}_k^{(t)} - \eta \nabla_{\mathbf{w}_k} \mathcal{R}(\Theta^{(t)}) = \mathbf{w}_k^{(t)} - \eta s_k \mathbf{X}^{\top} \mathbf{D}(\mathbf{X} \mathbf{w}_k^{(t)}) (h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}). \quad (2)$$

To analyze these updates more transparently, we introduce a primal-dual representation used in mirror descent (Ji and Telgarsky, 2021). For all  $k \in [m]$ , we define the primal variable  $\beta_k \in \mathbb{R}^n$  and the dual variable  $\alpha_k \in \mathbb{R}^n$  as

$$\beta_k := \mathbf{X} \mathbf{w}_k, \quad \alpha_k := (\mathbf{X} \mathbf{X}^{\top})^{-1} \mathbf{X} \mathbf{w}_k, \quad \text{and note that } \beta_k = \mathbf{X} \mathbf{X}^{\top} \alpha_k. \quad (3)$$

This representation restricts attention to the components of  $\mathbf{w}_k$  that lie in the span of the data matrix  $\mathbf{X}$ .<sup>3</sup> For ease of notation, we further define  $\beta_{k,+} := \mathbf{X}_+ \mathbf{w}_k$  and decompose the dual variable as  $\alpha_k := [\alpha_{k,+}^{\top} \quad \alpha_{k,-}^{\top}]^{\top}$ , consistent with the partition on labels  $\mathbf{y} = [\mathbf{y}_+^{\top} \quad \mathbf{y}_-^{\top}]^{\top}$ . Under this parameterization, the gradient descent update (2) can be expressed in primal-dual form as

$$\text{(Primal)} \quad \beta_k^{(t+1)} = \beta_k^{(t)} - \eta s_k \mathbf{X} \mathbf{X}^{\top} \mathbf{D}(\beta_k^{(t)}) (h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}), \quad (4a)$$

$$\text{(Dual)} \quad \alpha_k^{(t+1)} = \alpha_k^{(t)} - \eta s_k \mathbf{D}(\beta_k^{(t)}) (h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}). \quad (4b)$$

This primal-dual formulation plays a central role in our analysis. In particular, the sign of each primal coordinate  $\beta_{k,i}^{(t)}$  determines whether the corresponding dual variable  $\alpha_{k,i}^{(t+1)}$  is updated through the diagonal matrix  $\mathbf{D}(\beta_k^{(t)})$ . Consequently, understanding the positivity pattern of  $\beta_k^{(t)}$  and the resulting dynamics of  $\alpha_k^{(t)}$  is key to characterizing the behavior and implicit bias of gradient descent.

<sup>2</sup>This is a reasonable approximation for the dynamics when both layers are trained via the well-known *balancedness* condition, but balancedness is typically formally shown only under gradient flow (see, e.g. Du et al. 2018, Theorem 2.1).

<sup>3</sup>In general,  $\mathbf{w}_k$  may contain components orthogonal to  $\text{span}(\{\mathbf{x}_i\}_{i=1}^n)$ , i.e.,  $\mathbf{w}_k = \mathbf{X}^{\top} \alpha_k + \sum_{j=n+1}^d \tilde{\alpha}_{k,j} \tilde{\mathbf{x}}_j$  where  $\tilde{\alpha}_{k,j} \in \mathbb{R}$  and we define the vector  $\tilde{\mathbf{x}}_j \perp \mathbf{x}_i$  for all  $i \in [n]$  such that  $\{\mathbf{x}_i\}_{i=1}^n \cup \{\tilde{\mathbf{x}}_j\}_{j=n+1}^d$  forms a complete basis of  $\mathbb{R}^d$ . However, since the gradient updates act only within  $\text{span}(\{\mathbf{x}_i\}_{i=1}^n)$ , the orthogonal components remain unchanged throughout training. Our results can be easily extended by adding back in this orthogonal component.

**Minimum- $\ell_2$ -norm Solution.** It is well known that, for linear regression with zero initialization, i.e.,  $h(\mathbf{x}) := \mathbf{w}^\top \mathbf{x}$  with  $\mathbf{w}^{(0)} = \mathbf{0}$ , gradient descent converges to the minimum- $\ell_2$ -norm interpolation, which is given by  $\mathbf{w}_{\text{linear-MNI}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$ , s.t.  $\mathbf{w}^\top \mathbf{x}_i = y_i$ , for all  $i \in [n]$ . This solution admits the closed-form expression  $\mathbf{w}_{\text{linear-MNI}} = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{y}$ . Motivated by this classical result, we consider the minimum- $\ell_2$ -norm solution for the general ReLU regression problem that we study, defined as:

$$\begin{aligned} \{\mathbf{w}_k^*\}_{k=1}^m &= \arg \min_{\{\mathbf{w}_k\}_{k=1}^m} \frac{1}{2} \sum_{k=1}^m \|\mathbf{w}_k\|_2^2 \\ \text{s.t. } \sum_{k=1}^m s_k \sigma(\mathbf{w}_k^\top \mathbf{x}_i) &= y_i, \text{ for all } i \in [n]. \end{aligned} \quad (5)$$

Let  $\Theta_g$  denote the set of global minimizers of the empirical risk (1). Note that for  $m = 1$ , the empirical risk can often not be driven to zero; labels that are opposite in sign to the sign of the neuron  $s_1$  cannot be fit. For networks with  $m > 1$  neurons, we will consider at least one neuron to be positively signed and negatively signed respectively, ensuring that the global minimizers will achieve zero empirical risk and interpolate the training data (i.e.  $h_{\Theta}(\mathbf{x}_i) = y_i$ ,  $\forall i \in [n]$ ).

### 3. Implicit Bias of Single ReLU Models ( $m = 1$ ) Under Gradient Descent

We begin by analyzing the case of the single positive ReLU neuron model ( $m = 1$ ). Specifically, we consider  $h_{\Theta}(\mathbf{x}) := s_1 \sigma(\mathbf{w}^\top \mathbf{x})$  where  $\mathbf{w} \in \mathbb{R}^d$  is the only model parameter. We will also assume that  $s_1 = +1$  as will become clear through this section, the single neuron is only capable of fitting positive labels in this case. A symmetric version of our results in this section will hold in the opposite case where  $s_1 = -1$ , with all instances of positive labels replaced by negative labels. We omit this case for brevity.

#### 3.1. Gradient Descent Updates and Convergence

For the single ReLU model ( $m = 1$ ), the gradient descent update in (2) simplifies to

$$\begin{aligned} \mathbf{w}^{(t+1)} &= \mathbf{w}^{(t)} - \eta \nabla_{\mathbf{w}} \mathcal{R}(\mathbf{w}^{(t)}) = \mathbf{w}^{(t)} - \eta \mathbf{X}^\top \mathbf{D}(\mathbf{X} \mathbf{w}^{(t)}) (\sigma(\mathbf{X} \mathbf{w}^{(t)}) - \mathbf{y}) \\ &= \mathbf{w}^{(t)} - \eta \mathbf{X}^\top \mathbf{D}(\mathbf{X} \mathbf{w}^{(t)}) (\mathbf{X} \mathbf{w}^{(t)} - \mathbf{y}), \end{aligned} \quad (6)$$

where we write the vector-valued extension of the ReLU as  $\sigma(\mathbf{z}) := [\sigma(z_1), \dots, \sigma(z_n)]^\top \in \mathbb{R}^n$ , and the second equality follows from the fact that the diagonal matrix  $\mathbf{D}(\mathbf{X} \mathbf{w}^{(t)})$  enforces the ReLU activation pattern. Specifically, since  $\mathbf{D}(\mathbf{X} \mathbf{w}^{(t)})$  contains indicators of positive pre-activations, the explicit nonlinearity  $\sigma(\cdot)$  can be removed once it is applied.

Compared to linear regression, the key difference in the gradient update for a single ReLU model is the presence of the diagonal matrix  $\mathbf{D}(\mathbf{X} \mathbf{w}^{(t)})$ . This matrix effectively selects a subset of examples — those with positive pre-activations — to contribute to each gradient update. As a result, the optimization trajectory becomes both data-dependent and time-varying, with the active set of samples evolving during training.

### 3.1.1. SUFFICIENT CONDITIONS FOR GRADIENT DESCENT CONVERGENCE

According to Equation (6), convergence of gradient descent occurs when  $\nabla_{\mathbf{w}} \mathcal{R}(\mathbf{w}^{(t)}) = \mathbf{0}$ . This condition implies that, for every  $i \in [n]$ , either  $\mathbf{x}_i^\top \mathbf{w}^{(t)} \leq 0$  or  $\mathbf{x}_i^\top \mathbf{w}^{(t)} = y_i$ . In other words, at convergence, each training example is either inactive due to the ReLU nonlinearity or is fit exactly. These criteria can define either a global or local minimum depending on the activation pattern.

In general, the optimization trajectory and loss landscape induced by gradient descent, even on a single ReLU model, are difficult to characterize, primarily due to this data-dependent activation pattern. However, suppose there exists an iteration  $t_0 \geq 0$  such that, for all  $t \geq t_0$ , the set of active examples  $S := \{i \in [n] : \mathbf{x}_i^\top \mathbf{w}^{(t_0)} > 0\}$ , remains unchanged. In this “final phase”, we expect the gradient descent dynamics of the single ReLU model to reduce to those of linear regression restricted to the active subset of samples. We formalize this observation in the following lemma, which is proved in Appendix B.1.

**Lemma 1** *Suppose there exists  $t_0 \geq 0$  such that  $\mathbf{D}(\mathbf{X}\mathbf{w}^{(t_0)}) = \mathbf{D}(\mathbf{X}\mathbf{w}^{(t)})$  for all  $t \geq t_0$ . Define the subset of examples  $S := \{i \in [n] : \mathbf{x}_i^\top \mathbf{w}^{(t_0)} > 0\}$ . Then, for all  $t \geq t_0$ , the gradient descent dynamics of the single ReLU model are equivalent to gradient descent applied to a linear model initialized at  $\mathbf{w}^{(t_0)}$  and trained only on the subset  $S$ .*

As a direct consequence, convergence of the single ReLU model in the final phase follows from classical convergence guarantees for linear regression. In particular, it is straightforward to show that if the step size  $\eta \leq \frac{1}{\mu_1(\mathbf{X}\mathbf{X}^\top)}$ , then gradient descent converges in the final phase under our conditions on the training data. We state and prove this result for completeness in Appendix B.1.

**Lemma 2** *Suppose there exists  $t_0 \geq 0$  such that  $\mathbf{D}(\mathbf{X}\mathbf{w}^{(t_0)}) = \mathbf{D}(\mathbf{X}\mathbf{w}^{(t)})$  for all  $t \geq t_0$ , if the step size satisfies  $\eta \leq \frac{1}{\mu_1(\mathbf{X}\mathbf{X}^\top)}$ , gradient descent applied to the single ReLU model converges to  $\mathbf{w}^{(\infty)} = \arg \min_{\mathbf{w} \in \{\mathbf{w} : \mathbf{X}_S \mathbf{w} = \mathbf{y}_S\}} \|\mathbf{w} - \mathbf{w}^{(t_0)}\|_2$ , where  $S := \{i \in [n] : \mathbf{x}_i^\top \mathbf{w}^{(t_0)} > 0\}$ .*

## 3.2. Minimum- $\ell_2$ -norm Solution of Single ReLU Models

In Section 2, we discussed the minimum- $\ell_2$ -norm solution for linear regression (called  $\mathbf{w}_{\text{linear-MNI}}$ ). In contrast, due to the presence of the ReLU activation, single ReLU models can only produce nonnegative outputs. As a result, such models can minimize the empirical risk only by exactly fitting all samples with positive labels and outputting zero on samples with nonpositive labels. It is natural to consider the minimum- $\ell_2$ -norm solution for the single ReLU model subject to these constraints. Interestingly, this can be written as a convex optimization problem (despite the empirical risk itself being nonconvex) in which the constraints associated with nonpositive labels are written as linear inequalities, as below:

$$\begin{aligned} \mathbf{w}^\star &= \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2 & (7) \\ \text{s.t. } \mathbf{w}^\top \mathbf{x}_i &= y_i, \text{ for all } y_i > 0, \\ \mathbf{w}^\top \mathbf{x}_j &\leq 0, \text{ for all } y_j \leq 0. \end{aligned}$$

We show that the solution of (7) coincides with the minimum- $\ell_2$ -norm solution associated with *linearly* fitting a suitable subset of training examples, after setting all negative labels to zero. We define

the linear MNI solution associated with the subset  $S \subseteq [n]$  as  $\mathbf{w}_{\text{linear-MNI},S} = \mathbf{X}_S^\top (\mathbf{X}_S \mathbf{X}_S^\top)^{-1} \tilde{\mathbf{y}}_S$ , where  $\tilde{\mathbf{y}}_S \in \mathbb{R}^{|S|}$  denotes the corresponding modified label subvector with all negative entries replaced by zero.

**Lemma 3** *Consider a single ReLU model  $h_\Theta(\mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x})$ . The minimum- $\ell_2$ -norm solution  $\mathbf{w}^*$  of  $h_\Theta(\mathbf{x})$  defined in Equation (7) satisfies  $\mathbf{w}^* = \mathbf{w}_{\text{linear-MNI},S}$  for some index subset  $S \subseteq [n]$  that necessarily contains all indices  $i$  such that  $y_i > 0$ , where the corresponding labels are given by  $\tilde{y}_{S,i} = \max\{y_i, 0\}$ .*

Lemma 3 is proved in Appendix B.1 through the Karush-Kahn-Tucker (KKT) conditions. It is important to note that  $\mathbf{w}^*$  is a fundamentally different inductive bias from  $\mathbf{w}_{\text{linear-MNI}}$  as the subset  $S$  does not have an explicit formula, and is training data-dependent.

### 3.3. High-dimensional Implicit Bias of Single ReLU Models

Our first main result, stated below, characterizes the gradient descent dynamics of single ReLU models on high-dimensional data.

**Theorem 4** *Consider Assumptions 1 and 2, suppose the initialization is  $\mathbf{w}^{(0)} = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \boldsymbol{\epsilon}$ , where  $0 < \epsilon_i \leq \frac{1}{C_\alpha} y_{\min}$  for all  $i \in [n]$ , and the step size to satisfy  $\frac{1}{CC_g \|\boldsymbol{\lambda}\|_1} \leq \eta \leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1}$ . Then, the gradient descent limit  $\mathbf{w}^{(\infty)}$  for the single ReLU model coincides with the solution obtained by linear regression trained only on the positively labeled examples with initialization  $\mathbf{w}^{(1)} = \eta \mathbf{X}^\top \left( \mathbf{y} - \boldsymbol{\epsilon} + \frac{1}{\eta} (\mathbf{X} \mathbf{X}^\top)^{-1} \boldsymbol{\epsilon} \right)$  with probability at least  $1 - 2 \exp(-cn)$ . Formally, we have  $\mathbf{w}^{(\infty)} = \arg \min_{\mathbf{w} \in \{\mathbf{w}: \mathbf{X}_+ \mathbf{w} = \mathbf{y}_+\}} \|\mathbf{w} - \mathbf{w}^{(1)}\|_2$  and  $\mathbf{X}_- \mathbf{w}^{(\infty)} \preceq \mathbf{0}$ .*

Theorem 4 is proved in Appendix B.2 and characterizes a regime of gradient descent in which the convergence behavior is tractable. Due to the presence of the ReLU activation, the main challenge lies in monitoring which examples are active and which are inactive during gradient descent. Under our assumption of sufficiently high-dimensional data, we show, through careful tracking of the primal and dual variables, that examples with positive labels remain active throughout the optimization process (see Lemma 10), while examples with negative labels eventually become and remain inactive (see Lemma 11). Therefore, the limiting solution fits all positive labels exactly and produces predictions equal to zero for samples with negative labels. Consequently, this solution achieves the minimum empirical risk, i.e. is a specific global minimizer of (1).

**Remark 5** *In addition to Assumptions 1 and 2, Theorem 4 assumes a sufficiently small initialization where all the training examples are active (to see this, note that  $\mathbf{X} \mathbf{w}^{(0)} = \boldsymbol{\epsilon} \succ \mathbf{0}$ )<sup>4</sup>. Essentially, the primal variables are initialized in the positive orthant. The sufficiently small initialization is also assumed in previous work (Boursier et al., 2022; Dana et al., 2025). The positivity assumption is made to ensure high-probability convergence to a global minimum. On the other hand, a random initialization would map to both positive and negative primal variables. Our simulations in Appendix F (in particular, Figure 5) demonstrate that in this case, a positively labeled but initially*

<sup>4</sup>The initialization expression  $\mathbf{w}^{(0)} = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \boldsymbol{\epsilon}$  with  $\boldsymbol{\epsilon} \in \mathbb{R}^n$  is isomorphic to an arbitrary initialization in the span of the data  $\{\mathbf{x}_i\}_{i=1}^n$ , owing to the full-row-rank assumption on  $\mathbf{X}$ . Note in particular that such an initialization can be arbitrarily far from  $\mathbf{w}^*$ , and does not constitute a “local” initialization.

inactive example remains inactive, meaning that gradient descent can only converge to a local minimum<sup>5</sup>. Additionally, we provide a simple explicit counterexample showing that such initializations result in convergence to a local minimum that is not globally optimal in Appendix B.4.

### 3.4. Approximation to Minimum- $\ell_2$ -norm Solution in High Dimensions

We now show that the limiting solution obtained from Theorem 4 is different from, but close to the minimum- $\ell_2$ -norm solution in high dimensions. Specifically, the following theorem upper and lower bounds the Euclidean distance between  $\mathbf{w}^{(\infty)}$  and  $\mathbf{w}^*$  as a function of the number of negative examples  $n_-$ , effective dimension and label magnitude.

**Theorem 6** Consider Assumptions 1 and 2, suppose the initialization is  $\mathbf{w}^{(0)} = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \boldsymbol{\epsilon}$ , where  $0 < \epsilon_i \leq \frac{1}{C_\alpha} y_{\min}$  for all  $i \in [n]$ , and the step size to satisfy  $\frac{1}{CC_g \|\boldsymbol{\lambda}\|_1} \leq \eta \leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1}$ . Then, we have  $\sqrt{\frac{n_- y_{\min}^2}{CC_g \|\boldsymbol{\lambda}\|_1}} \leq \|\mathbf{w}^{(\infty)} - \mathbf{w}^*\|_2 \leq \sqrt{\frac{16n_- y_{\max}^2}{C_g \|\boldsymbol{\lambda}\|_1}}$  with probability at least  $1 - 2 \exp(-cn)$ .

Theorem 6 is proved in Appendix B.3 and heavily leverages our characterization of the minimum- $\ell_2$ -norm solution in Lemma 3. Our simulation in Figure 2 shows excellent agreement with Theorem 6. Note that Theorem 6 implies that  $\mathbf{w}^{(\infty)} = \mathbf{w}^* = \mathbf{w}_{\text{linear-MNI}}$  in the case where all labels are positive!

## 4. Implicit Bias of Two ReLU Models ( $m = 2$ ) Under Gradient Descent

We now extend our analysis to a 2-ReLU model ( $m = 2$ ), which combines one positive ReLU neuron and one negative ReLU neuron. More specifically, we define  $h_{\Theta}(\mathbf{x}) = \sigma(\mathbf{w}_{\oplus}^\top \mathbf{x}) - \sigma(\mathbf{w}_{\ominus}^\top \mathbf{x})$ , where  $\Theta$  is a set of model parameters such that  $\Theta := \{\mathbf{w}_{\oplus}, \mathbf{w}_{\ominus}\}$  and  $\mathbf{w}_{\oplus}, \mathbf{w}_{\ominus} \in \mathbb{R}^d$ . As mentioned in Section 2, this model is more expressive than the single ReLU model as it can perfectly fit arbitrary labels (both positive and negative). For the 2-ReLU model, the gradient descent update in (2) simplifies to

$$\begin{aligned} \mathbf{w}_{\oplus}^{(t+1)} &= \mathbf{w}_{\oplus}^{(t)} - \eta \nabla_{\mathbf{w}_{\oplus}} \mathcal{R}(\Theta^{(t)}) = \mathbf{w}_{\oplus}^{(t)} - \eta \mathbf{X}^\top \mathbf{D}(\mathbf{X} \mathbf{w}_{\oplus}^{(t)}) (h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}). \\ \mathbf{w}_{\ominus}^{(t+1)} &= \mathbf{w}_{\ominus}^{(t)} - \eta \nabla_{\mathbf{w}_{\ominus}} \mathcal{R}(\Theta^{(t)}) = \mathbf{w}_{\ominus}^{(t)} + \eta \mathbf{X}^\top \mathbf{D}(\mathbf{X} \mathbf{w}_{\ominus}^{(t)}) (h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}). \end{aligned}$$

### 4.1. Minimum- $\ell_2$ -norm Solution of 2-ReLU Models

First, we characterize the minimum- $\ell_2$ -norm solution for the 2-ReLU model, defined below:

$$\begin{aligned} \mathbf{w}_{\oplus}^*, \mathbf{w}_{\ominus}^* &= \arg \min_{\mathbf{w}_{\oplus}, \mathbf{w}_{\ominus}} \frac{1}{2} \|\mathbf{w}_{\oplus}\|_2^2 + \frac{1}{2} \|\mathbf{w}_{\ominus}\|_2^2 \\ \text{s.t. } &\sigma(\mathbf{w}_{\oplus}^\top \mathbf{x}_i) - \sigma(\mathbf{w}_{\ominus}^\top \mathbf{x}_i) = y_i, \text{ for all } i \in [n]. \end{aligned} \quad (8)$$

Unlike in the case of the single ReLU model, (8) cannot be written as a convex program. To analyze (8), we show that the optimal solution is also the optimal solution to a *restricted convex program* obtained by fixing the activation pattern of the two ReLU units across the training examples. To state this result, we define some additional notation. Let  $S_+ = \{i : y_i > 0, \text{ for all } i \in [n]\}$ ,  $S_- = \{j : y_j < 0, \text{ for all } j \in [n]\}$ , so that  $S_+ \cup S_- = [n]$  and  $S_+ \cap S_- = \emptyset$ .

<sup>5</sup>In fact, this is why Dana et al. (2025) need to assume a sufficiently large number of neurons  $m$  to ensure global convergence under random initialization.

**Lemma 7** *The feasible set of (8) is nonempty, and there exist partitions  $S_1 \cup S_2 = S_+$ ,  $S_1 \cap S_2 = \emptyset$ , and  $S_3 \cup S_4 = S_-$ ,  $S_3 \cap S_4 = \emptyset$  such that the optimal solution  $\{\mathbf{w}_\oplus^*, \mathbf{w}_\ominus^*\}$  of (8) is also an optimal solution of the following convex program:*

$$\begin{aligned} \mathbf{w}_\oplus^*, \mathbf{w}_\ominus^* &= \arg \min_{\mathbf{w}_\oplus, \mathbf{w}_\ominus} \frac{1}{2} \|\mathbf{w}_\oplus\|_2^2 + \frac{1}{2} \|\mathbf{w}_\ominus\|_2^2 & (9) \\ \text{s.t. } \mathbf{w}_\oplus^\top \mathbf{x}_i &= y_i, & \mathbf{w}_\ominus^\top \mathbf{x}_i \leq 0, & \text{ for all } i \in S_1, \\ \mathbf{w}_\oplus^\top \mathbf{x}_i - \mathbf{w}_\ominus^\top \mathbf{x}_i &= y_i, & -\mathbf{w}_\ominus^\top \mathbf{x}_i \leq 0, & \text{ for all } i \in S_2, \\ & -\mathbf{w}_\ominus^\top \mathbf{x}_i = y_i, & \mathbf{w}_\oplus^\top \mathbf{x}_i \leq 0, & \text{ for all } i \in S_3, \\ \mathbf{w}_\oplus^\top \mathbf{x}_i - \mathbf{w}_\ominus^\top \mathbf{x}_i &= y_i, & -\mathbf{w}_\oplus^\top \mathbf{x}_i \leq 0, & \text{ for all } i \in S_4. \end{aligned}$$

Lemma 7 is proved in Appendix C.1. Note that, in general, it is not possible to explicitly identify or characterize the exact activation patterns and corresponding partitions. However, the mere existence of such a partition is sufficient for our purposes and allows us to derive the desired approximation results for the implicit bias in Section 4.3.

## 4.2. High-dimensional Implicit Bias of 2-ReLU Models

Next, we characterize the gradient descent dynamics of two ReLU models in the high-dimensional regime in a manner analogous to the single-ReLU model (Theorem 4).

**Theorem 8** *Consider Assumptions 1 and 2, suppose the initialization  $\mathbf{w}_\oplus^{(0)} = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \boldsymbol{\epsilon}_\oplus$  and  $\mathbf{w}_\ominus^{(0)} = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \boldsymbol{\epsilon}_\ominus$ , where  $0 < \epsilon_{\oplus, i}, \epsilon_{\ominus, i} \leq \frac{1}{2C_\alpha} y_{\min}$  for all  $i \in [n]$ , and the step size to satisfy  $\frac{1}{CC_g \|\boldsymbol{\lambda}\|_1} \leq \eta \leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1}$ . Then, with probability at least  $1 - 2 \exp(-cn)$ , we have: The gradient descent limit  $\mathbf{w}_\oplus^{(\infty)}$  coincides with the solution obtained by linear regression trained only on the positively labeled examples, with the initialization  $\mathbf{w}_\oplus^{(1)} = \eta \mathbf{X}^\top (\mathbf{y} - \boldsymbol{\epsilon}_\oplus + \boldsymbol{\epsilon}_\ominus + \frac{1}{\eta} (\mathbf{X} \mathbf{X}^\top)^{-1} \boldsymbol{\epsilon}_\oplus)$ , and  $\mathbf{w}_\oplus^{(\infty)} = \arg \min_{\mathbf{w} \in \{\mathbf{w}: \mathbf{X}_+ \mathbf{w} = \mathbf{y}_+\}} \|\mathbf{w} - \mathbf{w}_\oplus^{(1)}\|_2$ ; the gradient descent limit  $\mathbf{w}_\ominus^{(\infty)}$  coincides with the solution obtained by linear regression trained only on the negatively labeled examples, with the initialization  $\mathbf{w}_\ominus^{(1)} = \eta \mathbf{X}^\top (-\mathbf{y} + \boldsymbol{\epsilon}_\oplus - \boldsymbol{\epsilon}_\ominus + \frac{1}{\eta} (\mathbf{X} \mathbf{X}^\top)^{-1} \boldsymbol{\epsilon}_\ominus)$  and  $\mathbf{w}_\ominus^{(\infty)} = \arg \min_{\mathbf{w} \in \{\mathbf{w}: \mathbf{X}_- \mathbf{w} = -\mathbf{y}_-\}} \|\mathbf{w} - \mathbf{w}_\ominus^{(1)}\|_2$ , with  $\mathbf{X}_- \mathbf{w}_\oplus^{(\infty)} \preceq \mathbf{0}$  and  $\mathbf{X}_+ \mathbf{w}_\ominus^{(\infty)} \preceq \mathbf{0}$ .*

Theorem 8 is proved in Appendix C.2 in a manner similar to the proof of Theorem 4. Our main additional insight is that, in high dimensions, the optimization dynamics naturally decouple:  $\mathbf{w}_\oplus$  learns to fit all positively labeled examples, while  $\mathbf{w}_\ominus$  learns to fit all negatively labeled examples.

## 4.3. Approximation to Minimum- $\ell_2$ -norm Solution in High Dimensions

Finally, we show, in a result analogous to Theorem 6, that the limiting solution of Theorem 8 is close to the minimum- $\ell_2$ -norm solution  $\{\mathbf{w}_\oplus^*, \mathbf{w}_\ominus^*\}$ .

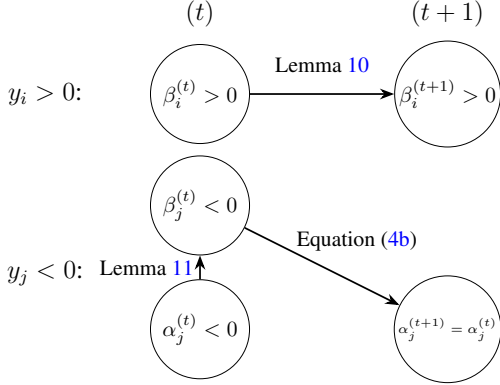


Figure 1: Gradient descent transition diagram for the  $k$ -th neuron.

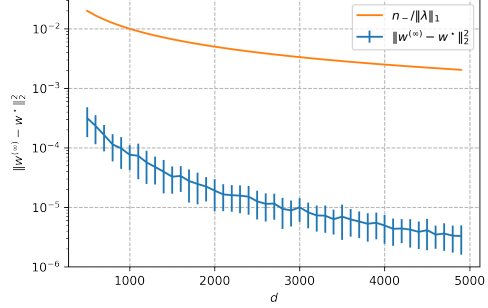


Figure 2: Approximation error between the implicit bias of the single ReLU model  $\mathbf{w}^{(\infty)}$  and the minimum- $\ell_2$ -norm solution  $\mathbf{w}^*$ .

**Theorem 9** Consider Assumptions 1 and 2, suppose the initialization is  $\mathbf{w}_{\oplus}^{(0)} = \mathbf{X}^{\top} (\mathbf{X} \mathbf{X}^{\top})^{-1} \boldsymbol{\epsilon}_{\oplus}$ ,  $\mathbf{w}_{\ominus}^{(0)} = \mathbf{X}^{\top} (\mathbf{X} \mathbf{X}^{\top})^{-1} \boldsymbol{\epsilon}_{\ominus}$ , where  $0 < \epsilon_{\oplus, i}, \epsilon_{\ominus, i} \leq \frac{1}{2C_{\alpha}} y_{\min}$  for all  $i \in [n]$ , and the step size satisfies  $\frac{1}{CC_g \|\boldsymbol{\lambda}\|_1} \leq \eta \leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1}$ . Then, we have  $\sqrt{\frac{n - y_{\min}^2}{CC_g \|\boldsymbol{\lambda}\|_1}} \leq \|\mathbf{w}_{\oplus}^{(\infty)} - \mathbf{w}_{\oplus}^*\|_2 \leq \sqrt{\frac{16n - y_{\max}^2}{C_g \|\boldsymbol{\lambda}\|_1}}$  and  $\sqrt{\frac{n + y_{\min}^2}{CC_g \|\boldsymbol{\lambda}\|_1}} \leq \|\mathbf{w}_{\ominus}^{(\infty)} - \mathbf{w}_{\ominus}^*\|_2 \leq \sqrt{\frac{16n + y_{\max}^2}{C_g \|\boldsymbol{\lambda}\|_1}}$  with probability at least  $1 - 2 \exp(-cn)$ .

Theorem 9 is proved in Appendix C.3 and leverages the restricted convex program that we derived in Lemma 7. Due to the relative complexity of (9), the proof becomes more involved than that of Theorem 6, but the underlying basic ideas are similar. Note that, because one of  $n_+, n_- > 0$ , the implicit bias of 2-ReLU cannot exactly coincide with the minimum- $\ell_2$ -norm solution.

## 5. Main Proof Ideas

Our analysis hinges on **precisely tracking the activation patterns of ReLU neurons across all training examples**. By controlling which examples remain active or inactive throughout training, we are able to understand the resulting gradient dynamics and, consequently, the implicit bias of the converged solution. To establish these results, we introduce two key lemmas. Lemma 10, inspired by ideas in Dana et al. (2025), shows that once the primal variable  $\beta_{k,i}$  corresponding to the  $k$ -th neuron and the  $i$ -th example is active—and the sign of the neuron  $s_k$  agrees with the label  $y_i$ —it remains active in the next iteration. This ensures that such an example is not suppressed by the ReLU nonlinearity and continues to contribute to the gradient updates.

**Lemma 10** Under Assumptions 1 and 2, suppose the gradient descent step size satisfies  $\eta \leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1}$ . Consider the  $k$ -th ReLU neuron in  $h_{\ominus}$ . For any  $t \geq 0$  and any index  $i \in [n]$  such that  $s_k \cdot y_i > 0$ , if  $\beta_{k,i}^{(t)} > 0$ ,  $\beta_{k,i}^{(t)} \geq s_k \cdot h_{\ominus}^{(t)}(\mathbf{x}_i)$ , and  $\|h_{\ominus}^{(t)}(\mathbf{X}) - \mathbf{y}\|_2 \leq C_y \|\mathbf{y}\|_2$ , then  $\beta_{k,i}^{(t+1)} > 0$  with probability at least  $1 - 2 \exp(-cn)$ .

This lemma is proved in Appendix A.2. The main idea behind Lemma 10 is that as long as the primal variable  $\beta_{k,i}^{(t)}$  is positive and the empirical risk remains uniformly bounded, the gradient update of  $\beta_{k,i}^{(t)}$  is dominated by its self-interaction term for high-dimensional data. The reason, at a high level, is that cross-sample interactions can be bounded in high dimensions (due to the concentration of the random Gram matrix  $\mathbf{X}\mathbf{X}^\top$  around  $\|\boldsymbol{\lambda}\|_1 \mathbf{I}$ ). As a result, the magnitude of the update is strictly smaller than  $\beta_{k,i}^{(t)}$ , ensuring that  $\beta_{k,i}^{(t+1)}$  remains positive.

Lemma 11 concerns the behavior of inactive examples. It shows that once a dual variable  $\alpha_{k,j}$  associated with the  $k$ -th neuron and the  $j$ -th example becomes sufficiently negative, the corresponding primal variable  $\beta_{k,j}$  remains inactive. Consequently, the dual variable is no longer updated and stays frozen throughout training. This mechanism effectively removes certain examples from the optimization dynamics.

**Lemma 11** *Under Assumptions 1 and 2, consider the  $k$ -th ReLU neuron in  $h_{\Theta}$ . For any  $t \geq 0$  and any index  $j \in [n]$ , if  $\alpha_{k,j}^{(t)} \leq -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1}$  and  $\|\boldsymbol{\alpha}_k^{(t)}\|_2 \leq \frac{C_\alpha \sqrt{ny_{\max}}}{\|\boldsymbol{\lambda}\|_1}$ , then  $\beta_{k,j}^{(t)} \leq 0$  and  $\alpha_{k,j}^{(t+1)} = \alpha_{k,j}^{(t)}$  with probability at least  $1 - 2 \exp(-cn)$ .*

The proof of Lemma 11 (see Appendix A.3) relies on the primal-dual relationship  $\beta_k = \mathbf{X}\mathbf{X}^\top \boldsymbol{\alpha}_k$  from Equation (3), together with concentration results for the Gram matrix. Specifically, if a dual variable is sufficiently negative, then the corresponding primal variable  $\beta_{k,j}^{(t)}$  is strictly negative. According to the dual update rule in Equation (4b), a negative  $\beta_{k,j}^{(t)}$  implies that the ReLU is inactive and the dual coordinate receives no further updates. As a result,  $\alpha_{k,j}^{(t+1)} = \alpha_{k,j}^{(t)}$ , and sufficiently negative dual variables remain frozen throughout training. Figure 1 depicts the transition of primal-dual updates in Lemma 10 and Lemma 11. We also provide deterministic-feature counterparts of Lemmas 10 and 11, stated as Lemmas 15 and 16, respectively. Their proofs are given in Appendix A.4.

In the following paragraphs, we outline the proof sketch for single ReLU models. The proof ideas for the 2-ReLU case follow analogously.

**Proof Sketch of Theorem 4:** The proof combines the insights from Lemma 10 and Lemma 11 to obtain a complete picture of how activation patterns evolve during training. Together, these lemmas allow us to track which examples remain active or inactive throughout gradient descent. Our goal is to reach—and maintain—a configuration in which positive-labeled examples remain active while negative-labeled examples remain inactive, as formalized by the sufficient conditions in Lemma 17 in Appendix B.2. To achieve this, we leverage two key properties of the initialization. First, the positive initialization guarantees that every example initially has at least one active neuron capable of fitting it. Second, the small initialization ensures that, after the first gradient step, positive-labeled examples remain in the active regime while negative-labeled examples acquire sufficiently negative dual variables and become inactive. Together, these properties place positive and negative examples into their respective regimes after a single update. We then apply Lemma 17 to show that this configuration is stable under subsequent iterations. As a result, the activation pattern becomes fixed after the first step, and the dynamics enter the final phase described in Lemma 1.

**Proof Sketch of Theorem 6:** To compare the gradient descent limit  $w^{(\infty)}$  with the minimum- $\ell_2$ -norm solution  $w^*$ , we relate their distance in parameter space to their distance in prediction space. Since both solutions interpolate all positive-labeled examples exactly, any discrepancy between them must arise from their predictions on negative-labeled examples. We bound this discrepancy

using the KKT conditions characterizing  $w^*$ , as established in Lemma 3. These conditions precisely describe how  $w^*$  treats negative-labeled examples and allow us to control the prediction distance in terms of the distance between the primal and dual variables. In particular, the KKT conditions imply that this gap is nonzero, showing that  $w^{(\infty)} \neq w^*$ . Translating our bounds back to parameter space yields matching upper and lower bounds on  $\|w^{(\infty)} - w^*\|_2$ .

## 6. Discussion

We showed that the implicit bias of single and 2-ReLU models, under appropriate initialization, is remarkably close to the minimum-norm solution if the features are sufficiently high-dimensional (and under appropriate initialization to ensure global convergence). Natural open questions include: 1) characterizing the dynamics for  $m > 2$  neurons, and 2) studying the effect of moderate dimension where  $d > n$  but not  $d \gg n$ . We provide partial extensions of our results to  $m > 2$  neurons in Appendices D and E that require a specific “disjoint” initialization, i.e., neurons are partitioned into sets such that they are active on disjoint examples. Handling more realistic initializations is an important direction for future work. We also simulate the effect of moderate-dimensional data on the dynamics in Appendix F and observe that the primal and dual variables intricately influence each other. We hope to characterize these more complex dynamics in future work, for which we will likely require different mathematical tools.

## Acknowledgments

KL gratefully acknowledges the support of the ARC-ACO Fellowship provided by Georgia Tech. GW gratefully acknowledges the Apple Scholars in AI/ML PhD fellowship by Apple and ARC-ACO fellowship provided by Georgia Tech. MT gratefully acknowledges the partial support of NSF Grant DMS-2513699, DOE Grants NA0004261, SC0026274, and Richard Duke Fellowship. VM gratefully acknowledges the support of the NSF (through award CCF-2239151 and award IIS-2212182), an Adobe Data Science Research Award, and an Amazon Research Award.

## References

- Peter L. Bartlett, Philip M. Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences (PNAS)*, 117(48):30063–30070, 2020.
- Rajendra Bhatia and Fuad Kittaneh. On the singular values of a product of operators. *SIAM Journal on Matrix Analysis and Applications*, 11(2):272–277, 1990.
- Etienne Boursier, Loucas Pillaud-Vivien, and Nicolas Flammarion. Gradient flow dynamics of shallow relu networks for square loss and orthogonal inputs. *Advances in Neural Information Processing Systems*, 35:20105–20118, 2022.
- Stephen Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- Léo Dana, Francis Bach, and Loucas Pillaud-Vivien. Convergence of shallow relu networks on weakly interacting data. *arXiv preprint arXiv:2502.16977*, 2025.

- Simon S Du, Wei Hu, and Jason D Lee. Algorithmic regularization in learning deep homogeneous models: Layers are automatically balanced. *Advances in neural information processing systems*, 31, 2018.
- Heinz Werner Engl, Martin Hanke, and Andreas Neubauer. *Regularization of inverse problems*, volume 375. Springer Science & Business Media, 1996.
- Spencer Frei, Gal Vardi, Peter Bartlett, and Nathan Srebro. Benign overfitting in linear classifiers and leaky relu networks from kkt conditions for margin maximization. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 3173–3228. PMLR, 2023a.
- Spencer Frei, Gal Vardi, Peter L Bartlett, and Wei Hu. Implicit bias in leaky relu networks trained on high-dimensional data. *ICLR*, 2023b.
- 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.
- Daniel Hsu, Vidya Muthukumar, and Ji Xu. On the proliferation of support vectors in high dimensions. In *International Conference on Artificial Intelligence and Statistics*, pages 91–99. PMLR, 2021.
- Ziwei Ji and Matus Telgarsky. Risk and parameter convergence of logistic regression. *arXiv preprint arXiv:1803.07300*, 2018.
- Ziwei Ji and Matus Telgarsky. Characterizing the implicit bias via a primal-dual analysis. In *Algorithmic Learning Theory*, pages 772–804. PMLR, 2021.
- Kuo-Wei Lai and Vidya Muthukumar. General loss functions lead to (approximate) interpolation in high dimensions. *Journal of Machine Learning Research*, 26(244):1–72, 2025.
- 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*, pages 3420–3428, 2019.
- Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias: On the role of implicit regularization in deep learning. *arXiv preprint arXiv:1412.6614*, 2014.
- 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*, 19(1):2822–2878, 2018.
- Haoyuan Sun, Kwangjun Ahn, Christos Thrampoulidis, and Navid Azizan. Mirror descent maximizes generalized margin and can be implemented efficiently. *Advances in Neural Information Processing Systems*, 35:31089–31101, 2022.
- Gal Vardi and Ohad Shamir. Implicit regularization in relu networks with the square loss. In *Conference on Learning Theory*, pages 4224–4258. PMLR, 2021.

Guanghai Wang, Zihao Hu, Claudio Gentile, Vidya Muthukumar, and Jacob Abernethy. Faster margin maximization rates for generic and adversarially robust optimization methods. *Mathematical Programming*, pages 1–41, 2025.

Ke Wang and Christos Thrampoulidis. Binary classification of Gaussian mixtures: Abundance of support vectors, benign overfitting, and regularization. *SIAM Journal on Mathematics of Data Science (SIMODS)*, 4:260–284, 2022.

Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.

# Appendix

## Table of Contents

---

<b>A Proofs of Key Lemmas Tracking Primal-Dual Gradient Dynamics</b>	<b>18</b>
A.1 Concentration of Random Gram Matrices in High Dimensions . . . . .	18
A.2 Proof of Lemma 10 (Primal Variable Gradient Dynamics in High Dimensions) . . .	19
A.3 Proof of Lemma 11 (Dual Variable Gradient Dynamics in High Dimensions) . . . .	20
A.4 Lemma 15 and Lemma 16 (Gradient Dynamics under Deterministic Features) . . . .	21
<b>B Proofs for the Single ReLU model (<math>m = 1</math>) Trained with Gradient Descent</b>	<b>24</b>
B.1 Proofs of Lemmas 1, 2 and 3 (Gradient Descent Convergence and $w^*$ ) . . . . .	24
B.2 Proof of Theorem 4 (High-dimensional Implicit Bias) . . . . .	27
B.3 Proof of Theorem 6 (Implicit Bias Approximation to $w^*$ ) . . . . .	34
B.4 Local Minimum Convergence Under a Not-All-Positive Initialization . . . . .	38
<b>C Proofs for the Two ReLU Model (<math>m = 2</math>) Trained with Gradient Descent</b>	<b>41</b>
C.1 Proof of Lemma 7 (Characterization of $w^*$ ) . . . . .	41
C.2 Proof of Theorem 8 (High-dimensional Implicit Bias) . . . . .	42
C.3 Proof of Theorem 9 (Implicit Bias Approximation to $w^*$ ) . . . . .	49
<b>D Implicit Bias of Multiple ReLU Models (<math>m &gt; 2</math>) Under Gradient Descent</b>	<b>59</b>
D.1 Gradient Descent Updates and Convergence . . . . .	59
D.2 Minimum- $\ell_2$ -norm Solution of Multiple ReLU Models . . . . .	60
D.3 High-dimensional Implicit Bias of Multiple ReLU Models . . . . .	60
D.4 Approximation to Minimum- $\ell_2$ -norm Solution in High Dimensions . . . . .	61
<b>E Proofs for Multiple ReLU Models (<math>m &gt; 2</math>) Trained with Gradient Descent</b>	<b>62</b>
E.1 Proof of Lemma 20 (Gradient Descent Convergence) . . . . .	62
E.2 Proof of Theorem 21 (High-dimensional Implicit Bias) . . . . .	64
E.3 Proof of Theorem 22 (Implicit Bias Approximation to $w^*$ ) . . . . .	72
<b>F Simulations</b>	<b>73</b>
F.1 Moderate-Dimensional Data and Single ReLU Model . . . . .	73
F.2 Gradient Descent Dynamics of Two ReLU Models . . . . .	74
F.3 Gradient Descent Dynamics of Multiple ReLU Models . . . . .	77

---

## Appendix A. Proofs of Key Lemmas Tracking Primal-Dual Gradient Dynamics

In this section, we present the proofs of the key lemmas used to track the gradient dynamics of the primal and dual variables. The central factor governing these dynamics is the sign pattern of the primal variables, which determines whether individual examples are active or inactive under the ReLU nonlinearity and, consequently, whether the corresponding dual variables are updated.

Before presenting the proofs, we first recall two key technical lemmas: 1) a concentration result on the eigenvalues of random Gram matrices in high dimensions from [Bartlett et al. \(2020\)](#); 2) a concentration bound on the operator norm of random Gram matrices from [Hsu et al. \(2021\)](#). Both these lemmas play a crucial role throughout the analysis.

### A.1. Concentration of Random Gram Matrices in High Dimensions

Our analysis relies heavily on properties of the Gram matrix on high-dimensional data. These concentration results allow us to control cross-sample interactions and isolate the dominant self-interaction terms that drive the gradient updates. As a result, we can rigorously characterize how positivity and negativity patterns in the primal and dual variables evolve over time.

In [Lemma 12](#), we characterize the typical behavior of the eigenvalues of a weighted sum of outer products of independent subgaussian vectors. Recall from [Section 2](#) that the feature vector  $\mathbf{x} \in \mathbb{R}^d$  admits the representation  $\mathbf{x} = \mathbf{V}\mathbf{\Lambda}^{\frac{1}{2}}\mathbf{z}$ , where  $\mathbf{z} \in \mathbb{R}^d$  has independent, mean-zero,  $\sigma_z^2$ -subgaussian components, and we take  $\sigma_z = 1$ . Under this model, the empirical Gram matrix can be written as  $\mathbf{X}\mathbf{X}^\top = \sum_{j=1}^d \lambda_j \mathbf{v}_j \mathbf{v}_j^\top$  where each  $\mathbf{v}_j \in \mathbb{R}^n$  is an independent random vector with independent, mean-zero, subgaussian entries. Concretely, [Lemma 12](#) provides high-probability bounds on the extreme eigenvalues of  $\mathbf{X}\mathbf{X}^\top$ .

**Lemma 12 ([Bartlett et al., 2020](#), [Lemma 9](#), [Wang and Thrampoulidis, 2022](#), [Lemma 12](#))**

*There exists a constant  $c$  such that with probability at least  $1 - 2e^{-n/c}$ , we have*

$$\frac{1}{c} \sum_{j=1}^d \lambda_j - c\lambda_1 n \leq \mu_n(\mathbf{X}\mathbf{X}^\top) \leq \mu_1(\mathbf{X}\mathbf{X}^\top) \leq c \left( \sum_{j=1}^d \lambda_j + \lambda_1 n \right).$$

*Moreover, if the effective dimension satisfies  $d_\infty = \frac{\sum_{j=1}^d \lambda_j}{\lambda_1} \geq bn$  for some constant  $b \geq 1$ , then there exists a constant  $C_g \geq 1$  such that*

$$\frac{1}{C_g} \sum_{j=1}^d \lambda_j \leq \mu_n(\mathbf{X}\mathbf{X}^\top) \leq \mu_1(\mathbf{X}\mathbf{X}^\top) \leq C_g \sum_{j=1}^d \lambda_j.$$

*with probability at least  $1 - 2e^{-n/C_g}$ .*

Next, [Lemma 13](#) provides a high-probability bound on the operator norm deviation between the Gram matrix  $\mathbf{X}\mathbf{X}^\top$  and  $\|\boldsymbol{\lambda}\|_1 \mathbf{I}$ , which is fruitful for high-dimensional data, and [Corollary 14](#) shows that the typical value of this deviation can be expressed in terms of  $n$  and effective dimensions  $d_2, d_\infty$ .

**Lemma 13 ([Hsu et al., 2021](#), [Lemma 8](#))** *There exists a universal constant  $c > 0$ , for any  $\tau > 0$ ,*

$$\Pr \left( \left\| \mathbf{X}\mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right\|_2 \geq \tau \right) \leq 2 \cdot 9^n \cdot \exp \left( -c \cdot \min \left\{ \frac{\tau^2}{\|\boldsymbol{\lambda}\|_2^2}, \frac{\tau}{\|\boldsymbol{\lambda}\|_\infty} \right\} \right),$$

where  $\|\boldsymbol{\lambda}\|_1 := \sum_{j=1}^d \lambda_j$ ,  $\|\boldsymbol{\lambda}\|_2^2 := \sum_{j=1}^d \lambda_j^2$ , and  $\|\boldsymbol{\lambda}\|_\infty := \max_{j \in [d]} \lambda_j$ .

**Corollary 14** *With the choice of  $\tau = C \cdot \max(\|\boldsymbol{\lambda}\|_2 \sqrt{n}, \|\boldsymbol{\lambda}\|_\infty n)$  and the constant  $C \cdot c > \ln 9$ , we have*

$$\left\| \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{X} \mathbf{X}^\top - \mathbf{I} \right\|_2 \leq C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right),$$

with probability at least  $1 - 2 \exp(-n(Cc - \ln 9))$ , where we have defined  $d_2 := \frac{\|\boldsymbol{\lambda}\|_1^2}{\|\boldsymbol{\lambda}\|_2^2}$ ,  $d_\infty := \frac{\|\boldsymbol{\lambda}\|_1}{\|\boldsymbol{\lambda}\|_\infty}$ . Similarly, we have

$$\left\| \|\boldsymbol{\lambda}\|_1 \left( \mathbf{X} \mathbf{X}^\top \right)^{-1} - \mathbf{I} \right\|_2 \leq C_g C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right),$$

with probability at least  $1 - 2 \exp(-n(Cc - \ln 9))$ .

## A.2. Proof of Lemma 10 (Primal Variable Gradient Dynamics in High Dimensions)

In this proof, we show that under the assumptions of the lemma, if the sign of any ReLU neuron agrees with the label of an example, then the corresponding primal variable remains positive after one gradient descent step.

**Proof** (Lemma 10) According to the primal gradient descent update in Equation (4a) for the  $k$ -th neuron, we have

$$\beta_k^{(t+1)} = \beta_k^{(t)} - \eta s_k \mathbf{X} \mathbf{X}^\top \mathbf{D}(\beta_k^{(t)})(h_{\boldsymbol{\Theta}^{(t)}}(\mathbf{X}) - \mathbf{y}).$$

We aim to separate the gradient contribution arising from the diagonal and off-diagonal components of the Gram matrix and to show that the updated primal coordinate remains positive, i.e.,  $\beta_{k,i}^{(t+1)} > 0$ . Fix any  $t \geq 0$  and any index  $i$  such that  $s_k \cdot y_i > 0$  and  $\beta_{k,i}^{(t)} > 0$ . Then, the update of the  $i$ -th coordinate can be written as

$$\begin{aligned} \beta_{k,i}^{(t+1)} &= \beta_{k,i}^{(t)} - \eta s_k \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top \mathbf{D}(\beta_k^{(t)})(h_{\boldsymbol{\Theta}^{(t)}}(\mathbf{X}) - \mathbf{y}) \\ &= \beta_{k,i}^{(t)} - \eta s_k \mathbf{e}_i^\top \left[ \|\boldsymbol{\lambda}\|_1 \mathbf{I} + \left( \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right) \right] \mathbf{D}(\beta_k^{(t)})(h_{\boldsymbol{\Theta}^{(t)}}(\mathbf{X}) - \mathbf{y}) \\ &= \left[ \beta_{k,i}^{(t)} - \eta \|\boldsymbol{\lambda}\|_1 (s_k h_{\boldsymbol{\Theta}^{(t)}}(\mathbf{x}_i) - s_k y_i) \right] \\ &\quad - \eta s_k \mathbf{e}_i^\top \left( \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right) \mathbf{D}(\beta_k^{(t)})(h_{\boldsymbol{\Theta}^{(t)}}(\mathbf{X}) - \mathbf{y}), \end{aligned} \quad (10)$$

where the last equality uses the assumption  $\beta_{k,i}^{(t)} > 0$ , which implies  $D_{ii} = \mathbb{1}_{\beta_{k,i}^{(t)} > 0} = 1$ . We now lower bound  $\beta_{k,i}^{(t+1)}$ . By the step size condition  $\eta \leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1}$  and the assumption  $\beta_{k,i}^{(t)} \geq s_k \cdot h_{\boldsymbol{\Theta}^{(t)}}(\mathbf{x}_i)$ , the first term in Equation (10) satisfies

$$\beta_{k,i}^{(t)} - \eta \|\boldsymbol{\lambda}\|_1 (s_k h_{\boldsymbol{\Theta}^{(t)}}(\mathbf{x}_i) - s_k y_i) \geq \eta \|\boldsymbol{\lambda}\|_1 |y_i|.$$

Substituting this into Equation (10) yields

$$\begin{aligned}
(10) &\geq \eta \|\boldsymbol{\lambda}\|_1 |y_i| - \eta s_k e_i^\top \left( \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right) \mathbf{D}(\boldsymbol{\beta}_k^{(t)})(h_{\boldsymbol{\Theta}^{(t)}}(\mathbf{X}) - \mathbf{y}) \\
&\geq \eta \|\boldsymbol{\lambda}\|_1 |y_i| - \eta \left\| \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right\|_2 \|h_{\boldsymbol{\Theta}^{(t)}}(\mathbf{X}) - \mathbf{y}\|_2, \tag{11}
\end{aligned}$$

where the last inequality follows from the Cauchy–Schwarz inequality and the sub-multiplicativity of the operator norm. Next, we upper bound the second term in Equation (11) using Corollary 14. With probability at least  $1 - 2 \exp(-n(Cc - \ln 9))$ , we obtain

$$\begin{aligned}
(11) &\geq \eta \|\boldsymbol{\lambda}\|_1 \left[ |y_i| - C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \|h_{\boldsymbol{\Theta}^{(t)}}(\mathbf{X}) - \mathbf{y}\|_2 \right] \\
&\stackrel{(i)}{\geq} \eta \|\boldsymbol{\lambda}\|_1 \left[ y_{\min} - C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \cdot C_y \sqrt{n} y_{\max} \right] \\
&\stackrel{(ii)}{\geq} \eta \|\boldsymbol{\lambda}\|_1 \left[ y_{\min} - C \cdot C_y \cdot \frac{y_{\min}}{C_0 y_{\max}} \cdot y_{\max} \right] \\
&> 0.
\end{aligned}$$

Inequality (i) applies the lemma assumption that  $\|h_{\boldsymbol{\Theta}^{(t)}}(\mathbf{X}) - \mathbf{y}\|_2 \leq C_y \|\mathbf{y}\|_2 \leq C_y \sqrt{n} y_{\max}$ . Inequality (ii) follows from Assumption 2, which guarantees that  $d_2 \geq C_0^2 \frac{n^2 y_{\max}^2}{y_{\min}^2}$  and  $d_\infty \geq C_0 \frac{n^{1.5} y_{\max}}{y_{\min}}$  with large enough  $C_0 > C \cdot C_y$ . This completes the proof of the lemma.  $\blacksquare$

### A.3. Proof of Lemma 11 (Dual Variable Gradient Dynamics in High Dimensions)

In this proof, we show that under the assumptions of the lemma, if the dual variable  $\alpha_{k,j}^{(t)}$  for the  $k$ -th neuron and  $j$ -th example is sufficiently negative, then it remains unchanged in the next iteration, i.e.,  $\alpha_{k,j}^{(t+1)} = \alpha_{k,j}^{(t)}$ .

**Proof** (Lemma 11) By the definition of primal and dual variables in Equation (3), we have

$$\boldsymbol{\beta}_k^{(t)} = \mathbf{X} \mathbf{X}^\top \boldsymbol{\alpha}_k^{(t)}.$$

According to the dual gradient update in Equation (4b), we have

$$\boldsymbol{\alpha}_k^{(t+1)} = \boldsymbol{\alpha}_k^{(t)} - \eta \mathbf{D}(\boldsymbol{\beta}_k^{(t)})(h_{\boldsymbol{\Theta}^{(t)}}(\mathbf{X}) - \mathbf{y}).$$

This update reveals that each coordinate  $\alpha_{k,j}^{(t)}$  evolves independently and is governed by the sign of the corresponding primal variable  $\beta_{k,j}^{(t)}$ . In particular, if  $\beta_{k,j}^{(t)} \leq 0$ , then the  $j$ -th diagonal entry of  $\mathbf{D}(\boldsymbol{\beta}_k^{(t)})$  vanishes, and consequently  $\alpha_{k,j}^{(t+1)} = \alpha_{k,j}^{(t)}$ .

We therefore establish a sufficient condition under which  $\beta_{k,j}^{(t)} \leq 0$  in terms of the dual variable  $\alpha_{k,j}^{(t)}$ . Specifically, we separate the diagonal and off-diagonal components of the Gram matrix as

$$\begin{aligned}
 \beta_{k,j}^{(t)} &= \mathbf{e}_j^\top \mathbf{X} \mathbf{X}^\top \boldsymbol{\alpha}_k^{(t)} \\
 &= \mathbf{e}_j^\top \left[ \|\boldsymbol{\lambda}\|_1 \mathbf{I} + (\mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}) \right] \boldsymbol{\alpha}_k^{(t)} \\
 &= \|\boldsymbol{\lambda}\|_1 \alpha_{k,j}^{(t)} + \mathbf{e}_j^\top (\mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}) \boldsymbol{\alpha}_k^{(t)} \\
 &\leq \|\boldsymbol{\lambda}\|_1 \alpha_{k,j}^{(t)} + \left\| \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right\|_2 \left\| \boldsymbol{\alpha}_k^{(t)} \right\|_2, \tag{12}
 \end{aligned}$$

where the last inequality follows from the sub-multiplicativity of the operator norm. Next, we upper bound the two terms  $\left\| \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right\|_2$  and  $\left\| \boldsymbol{\alpha}_k^{(t)} \right\|_2$  appearing in Equation (12). Following the same argument as in the proof of Lemma 10, we apply Corollary 14. Consequently, with probability at least  $1 - 2 \exp(-n(Cc - \ln 9))$ , we obtain

$$(12) \leq \|\boldsymbol{\lambda}\|_1 \left[ \alpha_{k,j}^{(t)} + C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \left\| \boldsymbol{\alpha}_k^{(t)} \right\|_2 \right]. \tag{13}$$

Finally, substituting the upper bound of  $\alpha_{k,j}^{(t)}$  and  $\left\| \boldsymbol{\alpha}_k^{(t)} \right\|_2$  in lemma assumptions into Equation (13), we obtain

$$\begin{aligned}
 (13) &\leq \|\boldsymbol{\lambda}\|_1 \left[ -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1} + C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \frac{C_\alpha \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1} \right] \\
 &= \|\boldsymbol{\lambda}\|_1 \left[ -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1} + C \cdot \frac{y_{\min}}{C_0 y_{\max}} \frac{C_\alpha y_{\max}}{\|\boldsymbol{\lambda}\|_1} \right] \\
 &\leq 0,
 \end{aligned}$$

where the last inequality follows from Assumption 2, which ensures  $d_2 \geq C_0^2 \frac{n^2 y_{\max}^2}{y_{\min}^2}$  and  $d_\infty \geq C_0 \frac{n^{1.5} y_{\max}}{y_{\min}}$  with large enough  $C_0 > C \cdot C_\alpha^2$ . We have thus shown that if  $\alpha_{k,j}^{(t)}$  is sufficiently negative, then  $\beta_{k,j}^{(t)} \leq 0$ , and consequently  $\alpha_{k,j}^{(t+1)} = \alpha_{k,j}^{(t)}$ . This completes the proof of the lemma.  $\blacksquare$

#### A.4. Lemma 15 and Lemma 16 (Gradient Dynamics under Deterministic Features)

In this section, we present the gradient dynamics of the primal and dual variables under deterministic feature assumptions. Lemma 15 serves as the deterministic-feature counterpart of Lemma 10.

**Lemma 15** *Under Assumptions 1 and 3, suppose the gradient descent step size satisfies  $\eta \leq \frac{1}{\mu_1(\mathbf{X} \mathbf{X}^\top)}$ . Consider the  $k$ -th ReLU neuron in  $h_\Theta$ . For any  $t \geq 0$  and any index  $i \in [n]$  such that  $s_k \cdot y_i > 0$ , if  $\beta_{k,i}^{(t)} > 0$ ,  $\beta_{k,i}^{(t)} \geq s_k \cdot h_{\Theta^{(t)}}(\mathbf{x}_i)$ , and  $\|h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}\|_2 \leq C_y \|\mathbf{y}\|_2$ , then  $\beta_{k,i}^{(t+1)} > 0$ .*

**Proof** (Lemma 15) According to the primal gradient descent update in Equation (4a) for the  $k$ -th neuron, we have

$$\boldsymbol{\beta}_k^{(t+1)} = \boldsymbol{\beta}_k^{(t)} - \eta s_k \mathbf{X} \mathbf{X}^\top \mathbf{D}(\boldsymbol{\beta}_k^{(t)})(h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}).$$

We aim to separate the gradient contribution arising from the diagonal and off-diagonal components of the Gram matrix and to show that the updated primal coordinate remains positive, i.e.,  $\beta_{k,i}^{(t+1)} > 0$ . Fix any  $t \geq 0$  and any index  $i$  such that  $s_k \cdot y_i > 0$  and  $\beta_{k,i}^{(t)} > 0$ . Then, the update of the  $i$ -th coordinate can be written as

$$\begin{aligned} \beta_{k,i}^{(t+1)} &= \beta_{k,i}^{(t)} - \eta s_k \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top \mathbf{D}(\boldsymbol{\beta}_k^{(t)})(h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}) \\ &= \left[ \beta_{k,i}^{(t)} - \eta s_k \|\mathbf{x}_i\|_2^2 \mathbb{1}_{\beta_{k,i}^{(t)} > 0} (h_{\Theta^{(t)}}(\mathbf{x}_i) - y_i) \right] - \eta s_k \sum_{j \neq i} \mathbf{x}_i^\top \mathbf{x}_j \mathbb{1}_{\beta_{k,j}^{(t)} > 0} (h_{\Theta^{(t)}}(\mathbf{x}_j) - y_j), \end{aligned} \quad (14)$$

where we have  $\mathbb{1}_{\beta_{k,i}^{(t)} > 0} = 1$  according to the assumption  $\beta_{k,i}^{(t)} > 0$ . We now lower bound  $\beta_{k,i}^{(t+1)}$ . By the step size condition  $\eta \leq \frac{1}{\mu_1(\mathbf{X} \mathbf{X}^\top)} \leq \frac{1}{\|\mathbf{x}_i\|_2^2}$  for all  $i \in [n]$  and the assumption  $\beta_{k,i}^{(t)} \geq s_k \cdot h_{\Theta^{(t)}}(\mathbf{x}_i)$ , the first term in Equation (14) satisfies

$$\beta_{k,i}^{(t)} - \eta \|\mathbf{x}_i\|_2^2 (s_k h_{\Theta^{(t)}}(\mathbf{x}_i) - s_k y_i) \geq \eta \|\mathbf{x}_i\|_2^2 |y_i|.$$

Substituting this into Equation (14) yields

$$\begin{aligned} (14) &\geq \eta \|\mathbf{x}_i\|_2^2 |y_i| - \eta s_k \sum_{j \neq i} \mathbf{x}_i^\top \mathbf{x}_j \mathbb{1}_{\beta_{k,j}^{(t)} > 0} (h_{\Theta^{(t)}}(\mathbf{x}_j) - y_j) \\ &\geq \eta \|\mathbf{x}_i\|_2^2 |y_i| - \eta \max_{i \neq j} |\mathbf{x}_i^\top \mathbf{x}_j| \sum_{j \neq i} |h_{\Theta^{(t)}}(\mathbf{x}_j) - y_j| \\ &\geq \eta \|\mathbf{x}_i\|_2^2 |y_i| - \eta \max_{i \neq j} |\mathbf{x}_i^\top \mathbf{x}_j| \sqrt{n} \|h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}\|_2, \end{aligned} \quad (15)$$

where we take absolute values in the second inequality, and the last inequality follows from the inequality between  $\ell_1$  and  $\ell_2$  norms. Next, we upper bound the second term in Equation (15) by the lemma assumption  $\|h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}\|_2 \leq C_y \|\mathbf{y}\|_2 \leq C_y \sqrt{n} y_{\max}$ . We obtain

$$\begin{aligned} (15) &\geq \eta \|\mathbf{x}_i\|_2^2 |y_i| - \eta \max_{i \neq j} |\mathbf{x}_i^\top \mathbf{x}_j| C_y n y_{\max} \\ &\geq \eta \left[ \min_{i \in [n]} \|\mathbf{x}_i\|_2^2 y_{\min} - C_y n \max_{i \neq j} |\mathbf{x}_i^\top \mathbf{x}_j| y_{\max} \right] \\ &> 0. \end{aligned}$$

The second inequality follows from taking the minimum for  $i \in [n]$ , and the last inequality follows from Assumption 3 with  $C_0 > C_y$ . This completes the proof of the lemma.  $\blacksquare$

Next, we present the gradient dynamics of the dual variables under deterministic feature assumptions. Lemma 16 serves as the deterministic-feature counterpart of Lemma 11.

**Lemma 16** *Under Assumptions 1 and 3, consider the  $k$ -th ReLU neuron in  $h_{\Theta}$ . For any  $t \geq 0$  and any index  $j \in [n]$ , if  $\alpha_{k,j}^{(t)} \leq -\frac{y_{\min}}{C_{\alpha}\|\boldsymbol{\lambda}\|_1}$  and  $\|\boldsymbol{\alpha}_k^{(t)}\|_2 \leq \frac{C_{\alpha}\sqrt{n}y_{\max}}{\|\boldsymbol{\lambda}\|_1}$ , then  $\beta_{k,j}^{(t)} \leq 0$  and  $\alpha_{k,j}^{(t+1)} = \alpha_{k,j}^{(t)}$ .*

**Proof** (Lemma 16) By the definition of primal and dual variables in Equation (3), we have

$$\boldsymbol{\beta}_k^{(t)} = \mathbf{X}\mathbf{X}^{\top}\boldsymbol{\alpha}_k^{(t)}.$$

According to the dual gradient update in Equation (4b), we have

$$\boldsymbol{\alpha}_k^{(t+1)} = \boldsymbol{\alpha}_k^{(t)} - \eta \mathbf{D}(\boldsymbol{\beta}_k^{(t)})(h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}).$$

This update reveals that each coordinate  $\alpha_{k,j}^{(t)}$  evolves independently and is governed by the sign of the corresponding primal variable  $\beta_{k,j}^{(t)}$ . In particular, if  $\beta_{k,j}^{(t)} \leq 0$ , then the  $j$ -th diagonal entry of  $\mathbf{D}(\boldsymbol{\beta}_k^{(t)})$  vanishes, and consequently  $\alpha_{k,j}^{(t+1)} = \alpha_{k,j}^{(t)}$ .

We therefore establish a sufficient condition under which  $\beta_{k,j}^{(t)} \leq 0$  in terms of the dual variable  $\alpha_{k,j}^{(t)}$ . Specifically, we separate the diagonal and off-diagonal components of the Gram matrix as

$$\begin{aligned} \beta_{k,j}^{(t)} &= \mathbf{e}_j^{\top} \mathbf{X}\mathbf{X}^{\top}\boldsymbol{\alpha}_k^{(t)} \\ &= \|\mathbf{x}_j\|_2^2 \alpha_{k,j}^{(t)} + \sum_{j \neq i} \mathbf{x}_j^{\top} \mathbf{x}_i \alpha_{k,i}^{(t)} \\ &\leq \|\mathbf{x}_j\|_2^2 \alpha_{k,j}^{(t)} + \max_{j \neq i} |\mathbf{x}_j^{\top} \mathbf{x}_i| \sum_{j \neq i} |\alpha_{k,i}^{(t)}|, \end{aligned} \quad (16)$$

where the last inequality takes the maximum and absolute values for the second term. Furthermore, by applying the inequality between  $\ell_2$  and  $\ell_1$  norms and taking the minimum for  $j \in [n]$  in the first term, we obtain

$$(16) \leq \left( \min_{j \in [n]} \|\mathbf{x}_j\|_2^2 \right) \alpha_{k,j}^{(t)} + \max_{j \neq i} |\mathbf{x}_j^{\top} \mathbf{x}_i| \sqrt{n} \|\boldsymbol{\alpha}_k^{(t)}\|_2. \quad (17)$$

Finally, substituting the upper bound of  $\alpha_{k,j}^{(t)}$  and  $\|\boldsymbol{\alpha}_k^{(t)}\|_2$  in lemma assumptions into Equation (17), we obtain

$$\begin{aligned} (17) &\leq - \left( \min_{j \in [n]} \|\mathbf{x}_j\|_2^2 \right) \frac{y_{\min}}{C_{\alpha}\|\boldsymbol{\lambda}\|_1} + \max_{j \neq i} |\mathbf{x}_j^{\top} \mathbf{x}_i| \sqrt{n} \frac{C_{\alpha}\sqrt{n}y_{\max}}{\|\boldsymbol{\lambda}\|_1} \\ &\leq 0, \end{aligned}$$

where the last inequality follows from Assumption 3 with large enough  $C_0 > C \cdot C_{\alpha}^2$ . We have thus shown that if  $\alpha_{k,j}^{(t)}$  is sufficiently negative, then  $\beta_{k,j}^{(t)} \leq 0$ , and consequently  $\alpha_{k,j}^{(t+1)} = \alpha_{k,j}^{(t)}$ . This completes the proof of the lemma.  $\blacksquare$

## Appendix B. Proofs for the Single ReLU model ( $m = 1$ ) Trained with Gradient Descent

In this section, we present the proofs concerning the behavior of the single ReLU model trained with gradient descent.

### B.1. Proofs of Lemmas 1, 2 and 3 (Gradient Descent Convergence and $w^*$ )

We present complete proofs of the gradient descent convergence for single ReLU models in Lemmas 1 and 2, as well as a characterization of the minimum- $\ell_2$ -norm solution in Lemma 3.

**Proof** (Lemma 1) We prove this lemma by showing that after iteration  $t_0 \geq 0$ , since the activation pattern is fixed, the gradient of the single ReLU model is equivalent to the gradient of a linear model using only a subset of examples. Consider a linear model

$$h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x},$$

where  $\mathbf{w} \in \mathbb{R}^d$  is the linear model parameter (also called weight). Let  $S \subseteq [n]$  denote the active set for the single ReLU model at iteration  $t_0$ , defined by  $S := \{i \in [n] : \mathbf{x}_i^\top \mathbf{w}^{(t_0)} > 0\}$ . We write the empirical risk with the linear model using only the examples in  $S$  as

$$\mathcal{R}_S(\mathbf{w}) = \frac{1}{2} \sum_{i \in S} (\mathbf{w}^\top \mathbf{x}_i - y_i)^2.$$

The gradient descent update for this linear model is

$$\begin{aligned} \mathbf{w}^{(t+1)} &= \mathbf{w}^{(t)} - \eta \nabla \mathcal{R}_S(\mathbf{w}^{(t)}) \\ &= \mathbf{w}^{(t)} - \eta \sum_{i \in S} (\mathbf{w}^{(t)\top} \mathbf{x}_i - y_i) \mathbf{x}_i. \end{aligned} \quad (18)$$

On the other hand, the original gradient descent dynamic for the single ReLU model (Equation 6) tells us that

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \mathbf{X}^\top \mathbf{D}(\mathbf{X} \mathbf{w}^{(t)}) (\mathbf{X} \mathbf{w}^{(t)} - \mathbf{y}).$$

Under the lemma assumption,  $\mathbf{D}(\mathbf{X} \mathbf{w}^{(t_0)}) = \mathbf{D}(\mathbf{X} \mathbf{w}^{(t)})$  for all  $t \geq t_0$ . Thus, we know that  $D_{ii} = \mathbb{1}_{i \in S}$  for all  $t \geq t_0$ . Therefore, for  $t \geq t_0$ , we can write the gradient update of the original single ReLU model as

$$\begin{aligned} \mathbf{w}^{(t+1)} &= \mathbf{w}^{(t)} - \eta \mathbf{X}^\top \mathbf{D}(\mathbf{X} \mathbf{w}^{(t)}) (\mathbf{X} \mathbf{w}^{(t)} - \mathbf{y}) \\ &= \mathbf{w}^{(t)} - \eta \sum_{i \in S} (\mathbf{w}^{(t)\top} \mathbf{x}_i - y_i) \mathbf{x}_i. \end{aligned}$$

This gradient update is equivalent to the gradient update of the linear model in Equation (18) for all  $t \geq t_0$ . As a result, for  $t \geq t_0$ , the gradient update of the single ReLU model is equivalent to a linear model using only data in  $S$ . This completes the proof of the lemma.  $\blacksquare$

**Proof** (Lemma 2) By Lemma 1, the activation pattern is fixed for all  $t \geq t_0$ , so the gradient descent update reduces to linear regression restricted to the active subset  $S$ , given by

$$\begin{aligned} \mathbf{w}^{(t+1)} &= \mathbf{w}^{(t)} - \eta \mathbf{X}^\top \mathbf{D}(\mathbf{X}\mathbf{w}^{(t)})(\mathbf{X}\mathbf{w}^{(t)} - \mathbf{y}) \\ &= \mathbf{w}^{(t)} - \eta \sum_{i \in S} (\mathbf{w}^{(t)\top} \mathbf{x}_i - y_i) \mathbf{x}_i \\ &= \mathbf{w}^{(t)} - \eta \mathbf{X}_S^\top (\mathbf{X}_S \mathbf{w}^{(t)} - \mathbf{y}_S). \end{aligned}$$

The final phase empirical risk is given by

$$\mathcal{R}(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}_S \mathbf{w} - \mathbf{y}_S\|_2^2 + \frac{1}{2} \|\mathbf{y}_{S^c}\|_2^2,$$

where the second term comes from the examples in  $S^c$  with negative pre-activations, and it does not depend on  $\mathbf{w}$  because the activation pattern does not change after  $t_0$ . Note that  $\mathcal{R}(\mathbf{w})$  is a convex quadratic with

$$\nabla \mathcal{R}(\mathbf{w}) = \mathbf{X}_S^\top (\mathbf{X}_S \mathbf{w} - \mathbf{y}_S), \quad \nabla^2 \mathcal{R}(\mathbf{w}) = \mathbf{X}_S^\top \mathbf{X}_S.$$

Therefore,  $\mathcal{R}$  is  $L$ -smooth with

$$L = \|\nabla^2 \mathcal{R}(\mathbf{w})\|_2 = \|\mathbf{X}_S^\top \mathbf{X}_S\|_2 = \mu_1(\mathbf{X}_S \mathbf{X}_S^\top).$$

A standard smoothness/descent result (e.g., [Boyd and Vandenberghe 2004](#), Equation 9.17) implies that for any  $\eta \leq \frac{1}{L}$ ,

$$\mathcal{R}(\mathbf{w}^{(t+1)}) \leq \mathcal{R}(\mathbf{w}^{(t)}) - \frac{\eta}{2} \|\nabla \mathcal{R}(\mathbf{w}^{(t)})\|_2^2,$$

and in particular,  $\mathcal{R}(\mathbf{w}^{(t)})$  is non-increasing for all  $t \geq t_0$ .

It remains to upper bound  $L$ . Since  $S$  is a subset of the training indices,  $|S| \leq n$ . Since  $L = \mu_1(\mathbf{X}_S \mathbf{X}_S^\top) \leq \mu_1(\mathbf{X} \mathbf{X}^\top)$ , choosing  $\eta \leq \frac{1}{\mu_1(\mathbf{X} \mathbf{X}^\top)}$  guarantees that  $\mathcal{R}(\mathbf{w}^{(t)})$  is non-increasing for all  $t \geq t_0$ . This establishes the desired step size condition in the final phase (and thus convergence in function value for the single ReLU dynamics after  $t_0$ ).

Finally, according to [Gunasekar et al. \(2018, Section 2.1\)](#), the set of minimizers of  $\mathcal{R}(\mathbf{w})$  is the affine subspace,

$$\mathcal{W}_S = \{\mathbf{w} : \mathbf{X}_S \mathbf{w} = \mathbf{y}_S\},$$

and gradient descent with constant step size converges to the Euclidean projection of the initialization  $\mathbf{w}^{(t_0)}$  onto this subspace  $\mathbf{w}^{(\infty)} = \arg \min_{\mathbf{w} \in \mathcal{W}_S} \|\mathbf{w} - \mathbf{w}^{(t_0)}\|_2$ . This completes the proof of the lemma. ■

**Proof** (Lemma 3) We prove the lemma by showing that the optimal solution  $\mathbf{w}^*$  of the original convex program for single ReLU models also solves a reduced convex program whose solution is

the minimum- $\ell_2$ -norm interpolation (MNI) over an index subset  $S \subseteq [n]$  with modified labels. First, we restate the convex program (7) and its KKT conditions below:

$$\begin{aligned} \mathbf{w}^* &\in \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t. } \mathbf{w}^\top \mathbf{x}_i &= y_i, \text{ for all } i \in S_1, \\ \mathbf{w}^\top \mathbf{x}_j &\leq 0, \text{ for all } j \in S_2, \end{aligned}$$

where we denote  $S_1 = \{i : y_i > 0, \text{ for all } i \in [n]\}$ ,  $S_2 = \{j : y_j \leq 0, \text{ for all } j \in [n]\}$  and  $S_1 \cup S_2 = [n]$ . Since  $n \leq d$  and we have assumed  $\text{rank}(\mathbf{X}) = n$ , we can always find a feasible solution satisfying all  $n$  equality constraints. This implies that the solution set is nonempty, and  $\mathbf{w}^*$  always exists. Hence, the following KKT conditions are necessary (and also sufficient) to  $\mathbf{w}^*$  for some  $\boldsymbol{\lambda}^* \in \mathbb{R}^{|S_1|}$  and  $\boldsymbol{\mu}^* \in \mathbb{R}^{|S_2|}$ :

**Stationarity:**

$$\mathbf{w}^* + \sum_{i \in S_1} \lambda_i^* \mathbf{x}_i + \sum_{j \in S_2} \mu_j^* \mathbf{x}_j = 0 \Leftrightarrow \mathbf{w}^* = - \sum_{i \in S_1} \lambda_i^* \mathbf{x}_i - \sum_{j \in S_2} \mu_j^* \mathbf{x}_j.$$

**Primal feasibility:**

$$\begin{aligned} \mathbf{w}^{*\top} \mathbf{x}_i &= y_i, \text{ for all } i \in S_1, \\ \mathbf{w}^{*\top} \mathbf{x}_j &\leq 0, \text{ for all } j \in S_2. \end{aligned}$$

**Dual feasibility:**

$$\begin{aligned} \lambda_i^* &\in \mathbb{R}, \text{ for all } i \in S_1, \\ \mu_j^* &\geq 0, \text{ for all } j \in S_2. \end{aligned}$$

**Complementary slackness:**

$$\sum_{j \in S_2} \mu_j^* (\mathbf{w}^{*\top} \mathbf{x}_j) = 0.$$

Next, we further denote a subset  $\tilde{S}_2 \subseteq S_2$  such that  $\tilde{S}_2 = \{j : \mu_j^* > 0 \text{ for all } j \in S_2\}$  (note that  $\tilde{S}_2$  can be empty). By the KKT conditions, it is necessary for  $\mathbf{w}^*$  to satisfy the following:

$$\mathbf{w}^* = - \sum_{i \in S_1} \lambda_i^* \mathbf{x}_i - \sum_{j \in \tilde{S}_2} \mu_j^* \mathbf{x}_j, \text{ with } \lambda_i^* \in \mathbb{R} \text{ and } \mu_j^* > 0, \quad (19a)$$

$$\mathbf{w}^{*\top} \mathbf{x}_i = y_i, \text{ for all } i \in S_1, \quad (19b)$$

$$\mathbf{w}^{*\top} \mathbf{x}_j = 0, \text{ for all } j \in \tilde{S}_2. \quad (19c)$$

Now, we consider a reduced convex program:

$$\begin{aligned} \tilde{\mathbf{w}} &\in \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t. } \mathbf{w}^\top \mathbf{x}_i &= y_i, \text{ for all } i \in S_1, \\ \mathbf{w}^\top \mathbf{x}_j &= 0, \text{ for all } j \in \tilde{S}_2. \end{aligned} \quad (20)$$

Its KKT conditions are given below.

**Stationarity:**

$$\tilde{\mathbf{w}} + \sum_{i \in S_1} \tilde{\lambda}_i \mathbf{x}_i + \sum_{j \in \tilde{S}_2} \tilde{\lambda}_j \mathbf{x}_j = 0 \Leftrightarrow \tilde{\mathbf{w}} = - \sum_{i \in S_1} \tilde{\lambda}_i \mathbf{x}_i - \sum_{j \in \tilde{S}_2} \tilde{\lambda}_j \mathbf{x}_j.$$

**Primal feasibility:**

$$\begin{aligned} \tilde{\mathbf{w}}^\top \mathbf{x}_i &= y_i, \text{ for all } i \in S_1, \\ \tilde{\mathbf{w}}^\top \mathbf{x}_j &= 0, \text{ for all } j \in \tilde{S}_2. \end{aligned}$$

**Dual feasibility:**

$$\begin{aligned} \tilde{\lambda}_i &\in \mathbb{R}, \text{ for all } i \in S_1, \\ \tilde{\lambda}_j &\in \mathbb{R}, \text{ for all } j \in \tilde{S}_2. \end{aligned}$$

Since  $\mathbf{w}^*$  satisfies all the conditions in Equation (19), it also satisfies the KKT conditions for the reduced convex program (20). Thus,  $\mathbf{w}^*$  is also the optimal solution of the reduced convex program. Finally, we have a closed-form solution for the reduced convex program such that  $\mathbf{w}^* = \tilde{\mathbf{w}} = \mathbf{w}_{\text{linear-MNI},S} = \mathbf{X}_S^\top (\mathbf{X}_S \mathbf{X}_S^\top)^{-1} \tilde{\mathbf{y}}_S$  where  $S = S_1 \cup \tilde{S}_2$  and  $\tilde{\mathbf{y}}_S$  denotes the corresponding label subvector with all negative entries replaced by zero. This completes the proof of the lemma.  $\blacksquare$

## B.2. Proof of Theorem 4 (High-dimensional Implicit Bias)

In this section, we present the proof of Theorem 4. For the single ReLU model ( $m = 1$ ), the primal-dual gradient update in (4) simplifies to

$$\text{(Primal)} \quad \boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} - \eta \mathbf{X} \mathbf{X}^\top \mathbf{D}(\boldsymbol{\beta}^{(t)}) (\boldsymbol{\beta}^{(t)} - \mathbf{y}), \quad (21a)$$

$$\text{(Dual)} \quad \boldsymbol{\alpha}^{(t+1)} = \boldsymbol{\alpha}^{(t)} - \eta \mathbf{D}(\boldsymbol{\beta}^{(t)}) (\boldsymbol{\beta}^{(t)} - \mathbf{y}). \quad (21b)$$

Before proceeding to the proof, we introduce a set of sufficient conditions under which the signs of the primal variables agree with the signs of the labels at iteration  $t$ . Moreover, these conditions are preserved at iteration  $t + 1$ .

**Lemma 17** *Under Assumptions 1 and 2, suppose the gradient descent step size satisfies  $\eta \leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1}$ . For any single ReLU model, if the following six conditions hold at some iteration  $t \geq 0$ , then they also hold at iteration  $t + 1$ .*

- a.  $\beta_i^{(t)} > 0$ , for all  $i \in [n]$  with  $y_i > 0$ .
- b.  $-\frac{3y_{\max}}{C_g \|\boldsymbol{\lambda}\|_1} \leq \alpha_j^{(t)} \leq -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1}$ , for all  $j \in [n]$  with  $y_j < 0$ .
- c.  $\left\| \boldsymbol{\beta}_+^{(t)} - \mathbf{y}_+ \right\|_2 \leq C_y \|\mathbf{y}_+\|_2$ .
- d.  $\left\| \boldsymbol{\alpha}^{(t)} \right\|_2 \leq \frac{C_\alpha \sqrt{ny_{\max}}}{\|\boldsymbol{\lambda}\|_1}$ .

e.  $\beta_j^{(t)} \leq 0$ , for all  $j \in [n]$  with  $y_j < 0$ .

f.  $\sigma(\boldsymbol{\beta}^{(t)}) = \begin{bmatrix} \boldsymbol{\beta}_+^{(t)} \\ \mathbf{0} \end{bmatrix}$ .

Consequently, the set of active examples consists exactly of the positively labeled examples, and the activation pattern remains unchanged, i.e.,  $\mathbf{D}(\boldsymbol{\beta}^{(t)}) = \mathbf{D}(\boldsymbol{\beta}^{(t+1)})$ .

**Proof** (Lemma 17) In the following, we show that if the six sufficient conditions hold at some iteration  $t \geq 0$ , then they also hold at iteration  $t + 1$ .

Part (a): By conditions (c) and (f) at iteration  $t$ , we have  $\|h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}\|_2^2 = \|\sigma(\boldsymbol{\beta}^{(t)}) - \mathbf{y}\|_2^2 = \|\boldsymbol{\beta}_+^{(t)} - \mathbf{y}_+\|_2^2 + \|\mathbf{y}_-\|_2^2 \leq C_y^2 \|\mathbf{y}\|_2^2$ . Together with  $h_{\Theta^{(t)}}(\mathbf{x}_i) = \beta_i^{(t)}$  and condition (a), all the assumptions of Lemma 10 are satisfied for all  $i$  with  $y_i > 0$ . Consequently, we obtain  $\beta_i^{(t+1)} > 0$  for all  $i \in [n]$  with  $y_i > 0$ , and thus condition (a) holds at iteration  $t + 1$ .

Part (b): According to the dual gradient update in Equation (21b), and using condition (e) at iteration  $t$ , we conclude that the dual variables corresponding to negatively labeled examples remain unchanged, i.e.,  $\alpha_j^{(t+1)} = \alpha_j^{(t)}$  for all  $j \in [n]$  with  $y_j < 0$ . Therefore, condition (b) continues to hold at iteration  $t + 1$ .

Part (c): By conditions (a) and (e), the gradient update at iteration  $t$  depends only on the positively labeled examples. Consequently, the update is equivalent to a linear regression gradient descent step using only the positive-labeled subset. As similarly argued in the proof of Lemma 2, since the step size satisfies  $\eta \leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1}$ , the squared loss is monotonically non-increasing, and we obtain  $\|\boldsymbol{\beta}_+^{(t+1)} - \mathbf{y}_+\|_2 \leq \|\boldsymbol{\beta}_+^{(t)} - \mathbf{y}_+\|_2 \leq C_y \|\mathbf{y}_+\|_2$  by condition (c) at iteration  $t$ . Therefore, condition (c) holds at iteration  $t + 1$ .

Part (d): For this part, we use conditions (b) and (c) at iteration  $t + 1$ . By the triangle inequality, we have

$$\|\boldsymbol{\alpha}^{(t+1)}\|_2 \leq \|\boldsymbol{\alpha}_+^{(t+1)}\|_2 + \|\boldsymbol{\alpha}_-^{(t+1)}\|_2.$$

By condition (b) at iteration  $t + 1$ , it follows that  $\|\boldsymbol{\alpha}_-^{(t+1)}\|_2 \leq \frac{3\sqrt{n}y_{\max}}{C_g \|\boldsymbol{\lambda}\|_1}$ . It therefore remains to upper bound  $\|\boldsymbol{\alpha}_+^{(t+1)}\|_2$ . By condition (c) at iteration  $t + 1$ , we have  $\|\boldsymbol{\beta}_+^{(t+1)}\|_2 \leq C_y \|\mathbf{y}_+\|_2 + \|\mathbf{y}_+\|_2 \leq (C_y + 1) \|\mathbf{y}\|_2$ . Moreover, we have

$$\begin{aligned} \|\boldsymbol{\beta}_+^{(t+1)}\|_2 &= \|\mathbf{X}_+ \mathbf{X}_+^\top \boldsymbol{\alpha}^{(t+1)}\|_2 \\ &= \left\| \mathbf{X}_+ \begin{bmatrix} \mathbf{X}_+^\top \mathbf{X}_+^\top \\ \mathbf{X}_-^\top \mathbf{X}_-^\top \end{bmatrix} \begin{bmatrix} \boldsymbol{\alpha}_+^{(t+1)} \\ \boldsymbol{\alpha}_-^{(t+1)} \end{bmatrix} \right\|_2 \\ &= \left\| \mathbf{X}_+ \mathbf{X}_+^\top \boldsymbol{\alpha}_+^{(t+1)} + \mathbf{X}_+ \mathbf{X}_-^\top \boldsymbol{\alpha}_-^{(t+1)} \right\|_2. \end{aligned}$$

Applying the triangle inequality yields

$$\begin{aligned} \left\| \mathbf{X}_+ \mathbf{X}_+^\top \boldsymbol{\alpha}_+^{(t+1)} \right\|_2 &\leq \left\| \boldsymbol{\beta}_+^{(t+1)} \right\|_2 + \left\| \mathbf{X}_+ \mathbf{X}_-^\top \boldsymbol{\alpha}_-^{(t+1)} \right\|_2 \\ &\leq (C_y + 1) \|\mathbf{y}\|_2 + \left\| \mathbf{X}_+ \mathbf{X}_-^\top \boldsymbol{\alpha}_-^{(t+1)} \right\|_2. \end{aligned}$$

Since  $\mathbf{X}_+ \mathbf{X}_+^\top \in \mathbb{R}^{n_+ \times n_+}$  is full rank, we obtain

$$\left\| \boldsymbol{\alpha}_+^{(t+1)} \right\|_2 \leq \frac{(C_y + 1) \|\mathbf{y}\|_2 + \left\| \mathbf{X}_+ \mathbf{X}_-^\top \boldsymbol{\alpha}_-^{(t+1)} \right\|_2}{\mu_{n_+}(\mathbf{X}_+ \mathbf{X}_+^\top)}.$$

For the denominator, the variational formulation for eigenvalues of a submatrix and Lemma 12 imply that

$$\mu_{n_+}(\mathbf{X}_+ \mathbf{X}_+^\top) \geq \mu_n(\mathbf{X} \mathbf{X}^\top) \geq \frac{1}{C_g} \sum_{j=1}^d \lambda_j = \frac{\|\boldsymbol{\lambda}\|_1}{C_g},$$

with probability at least  $1 - 2e^{-n/C_g}$ . For the numerator, we have  $(C_y + 1) \|\mathbf{y}\|_2 \leq (C_y + 1) \sqrt{n} y_{\max}$ . Moreover, by [Bhatia and Kittaneh \(1990, Theorem 1\)](#), we have

$$\begin{aligned} \left\| \mathbf{X}_+ \mathbf{X}_-^\top \right\|_2 &\leq \frac{1}{2} \left\| \mathbf{X}_+ \mathbf{X}_+^\top + \mathbf{X}_- \mathbf{X}_-^\top \right\|_2 \\ &\leq \frac{1}{2} \left( \left\| \mathbf{X}_+ \mathbf{X}_+^\top \right\|_2 + \left\| \mathbf{X}_- \mathbf{X}_-^\top \right\|_2 \right) \\ &\leq C_g \sum_{j=1}^d \lambda_j \\ &= C_g \|\boldsymbol{\lambda}\|_1, \end{aligned}$$

where the last inequality follows from Lemma 12. Combining these bounds yields

$$\left\| \boldsymbol{\alpha}_+^{(t+1)} \right\|_2 \leq \frac{(C_y + 1) \sqrt{n} y_{\max} + C_g \|\boldsymbol{\lambda}\|_1 \cdot \frac{3\sqrt{n} y_{\max}}{C_g \|\boldsymbol{\lambda}\|_1}}{\|\boldsymbol{\lambda}\|_1 / C_g} = ((C_y + 1) C_g + 3C_g) \frac{\sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1}.$$

Consequently, we have

$$\left\| \boldsymbol{\alpha}^{(t+1)} \right\|_2 \leq ((C_y + 1) C_g + 3C_g) \frac{\sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1} + \frac{3\sqrt{n} y_{\max}}{C_g \|\boldsymbol{\lambda}\|_1} \leq \frac{C_\alpha \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1},$$

with  $C_\alpha \gtrsim \max\{C_g^2, C_y C_g\}$ , and thus condition (d) holds at iteration  $t + 1$ .

Part (e): By Lemma 11, and since conditions (b) and (d) hold at iteration  $t + 1$ , we conclude that  $\beta_j^{(t+1)} \leq 0$  for all  $j \in [n]$  with  $y_j < 0$ . Thus, condition (e) holds at iteration  $t + 1$ .

Part (f): By conditions (a) and (e) at iteration  $t + 1$ , the signs of the primal variables continue to agree with the signs of the labels. Consequently,  $\sigma(\boldsymbol{\beta}^{(t+1)}) = \begin{bmatrix} \boldsymbol{\beta}_+^{(t+1)} \\ \mathbf{0} \end{bmatrix}$ , and thus condition (f) holds at iteration  $t + 1$ .

We have shown that the six sufficient conditions hold at iteration  $t + 1$ . Consequently, the signs of the primal variables continue to agree with the signs of the labels, and hence  $\mathbf{D}(\boldsymbol{\beta}^{(t)}) = \mathbf{D}(\boldsymbol{\beta}^{(t+1)})$ . This completes the proof.  $\blacksquare$

Equipped with Lemma 17, we are now ready to prove Theorem 4.

**Proof** (Theorem 4) In the proof, we first show that after the first gradient step, the iterate at  $t = 1$  satisfies the conditions in Lemma 17. Next, since the conditions hold at  $t = 1$  and are preserved from  $t = \tilde{t}$  to  $t = \tilde{t} + 1$  by Lemma 17, we fully characterize the gradient descent dynamics by induction.

We begin by verifying that the iterate at  $t = 1$  satisfies the sufficient conditions in Lemma 17. With the initialization  $\mathbf{w}^{(0)} = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \boldsymbol{\epsilon}$ , we have  $\boldsymbol{\beta}^{(0)} = \mathbf{X} \mathbf{w}^{(0)} = \boldsymbol{\epsilon}$ . Therefore, using the primal gradient update in Equation (21a), we obtain

$$\begin{aligned} \boldsymbol{\beta}^{(1)} &= \boldsymbol{\beta}^{(0)} - \eta \mathbf{X} \mathbf{X}^\top \mathbf{D}(\boldsymbol{\beta}^{(0)}) (\boldsymbol{\beta}^{(0)} - \mathbf{y}) \\ &= \boldsymbol{\epsilon} - \eta \mathbf{X} \mathbf{X}^\top (\boldsymbol{\epsilon} - \mathbf{y}) \\ &= \mathbf{X} \mathbf{X}^\top \left[ \underbrace{\eta \left( \mathbf{y} - \boldsymbol{\epsilon} + \frac{1}{\eta} (\mathbf{X} \mathbf{X}^\top)^{-1} \boldsymbol{\epsilon} \right)}_{=: \boldsymbol{\alpha}^{(1)}} \right]. \end{aligned} \quad (22)$$

We denote  $\boldsymbol{\alpha}^{(1)} := \eta \left( \mathbf{y} - \boldsymbol{\epsilon} + \frac{1}{\eta} (\mathbf{X} \mathbf{X}^\top)^{-1} \boldsymbol{\epsilon} \right)$  according to the primal-dual formulation  $\boldsymbol{\beta}^{(1)} = \mathbf{X} \mathbf{X}^\top \boldsymbol{\alpha}^{(1)}$  in Equation (3). In the below, we show that at iteration  $t = 1$ , the variables  $\boldsymbol{\beta}^{(1)}$  and  $\boldsymbol{\alpha}^{(1)}$  satisfy all the conditions in Lemma 17.

Part (a): For all  $i \in [n]$  with  $y_i > 0$ , we apply Lemma 10. Since  $\beta_i^{(0)} = \epsilon_i > 0$ ,  $\beta_i^{(0)} = h_{\Theta^{(0)}}(\mathbf{x}_i)$  and  $\left\| \sigma(\boldsymbol{\beta}^{(0)}) - \mathbf{y} \right\|_2 \leq \|\boldsymbol{\epsilon}\|_2 + \|\mathbf{y}\|_2 \leq \frac{\sqrt{n}}{C_\alpha} y_{\min} + \|\mathbf{y}\|_2 \leq C_y \|\mathbf{y}\|_2$  with  $C_y > 1 + \frac{1}{C_\alpha}$ , it follows that  $\beta_i^{(1)} > 0$  for all  $i \in [n]$  with  $y_i > 0$ .

Part (b): For all  $j \in [n]$  with  $y_j < 0$ , we verify that  $\alpha_j^{(1)}$  satisfies the required upper and lower bounds. For the upper bound, recall that

$$\begin{aligned} \alpha_j^{(1)} &= \eta \left( y_j - \epsilon_j + \frac{1}{\eta} \mathbf{e}_j^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \boldsymbol{\epsilon} \right) \\ &= \eta \left( y_j - \epsilon_j + \frac{1}{\eta} \mathbf{e}_j^\top \left[ \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} + \left( (\mathbf{X} \mathbf{X}^\top)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right) \right] \boldsymbol{\epsilon} \right) \\ &= \eta \left( y_j - \epsilon_j + \frac{\epsilon_j}{\eta \|\boldsymbol{\lambda}\|_1} + \frac{1}{\eta} \mathbf{e}_j^\top \left( (\mathbf{X} \mathbf{X}^\top)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right) \boldsymbol{\epsilon} \right) \\ &\stackrel{(i)}{\leq} \eta \left( y_j + \frac{\epsilon_j}{\eta \|\boldsymbol{\lambda}\|_1} + \frac{1}{\eta} \mathbf{e}_j^\top \left( (\mathbf{X} \mathbf{X}^\top)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right) \boldsymbol{\epsilon} \right) \\ &\stackrel{(ii)}{\leq} \eta \left( y_j + \frac{\epsilon_j}{\eta \|\boldsymbol{\lambda}\|_1} + \frac{1}{\eta} \left\| \left( (\mathbf{X} \mathbf{X}^\top)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right) \boldsymbol{\epsilon} \right\|_2 \right), \end{aligned}$$

where inequality (i) drops the negative term  $-\epsilon_j$ , and inequality (ii) follows from the submultiplicativity of the operator norm. By Corollary 14, we have

$$\left\| \left( \mathbf{X} \mathbf{X}^\top \right)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right\|_2 \leq \frac{C_g C}{\|\boldsymbol{\lambda}\|_1} \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right),$$

with probability at least  $1 - 2 \exp(-n(Cc - \ln 9))$ . Moreover, by the theorem assumptions,  $\|\boldsymbol{\epsilon}\|_2 \leq \frac{\sqrt{n}}{C_\alpha} y_{\min}$  and  $\frac{1}{\eta} \leq CC_g \|\boldsymbol{\lambda}\|_1$ . Combining these bounds yields

$$\begin{aligned} \alpha_j^{(1)} &\leq \frac{1}{CC_g \|\boldsymbol{\lambda}\|_1} \left( -y_{\min} + \frac{CC_g}{C_\alpha} y_{\min} + C^2 C_g^2 \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \cdot \frac{\sqrt{n}}{C_\alpha} y_{\min} \right) \\ &\leq \frac{1}{CC_g \|\boldsymbol{\lambda}\|_1} \left( -y_{\min} + \frac{CC_g}{C_\alpha} y_{\min} + C^2 C_g^2 \cdot \frac{y_{\min}}{C_0 y_{\max}} \cdot \frac{1}{C_\alpha} y_{\min} \right) \\ &= -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1} \left( \frac{C_\alpha}{CC_g} - 1 - \frac{CC_g y_{\min}}{C_0 y_{\max}} \right) \\ &\leq -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1}. \end{aligned}$$

The second inequality follows from  $d_2 \geq C_0^2 \frac{n^2 y_{\max}^2}{y_{\min}^2}$  and  $d_\infty \geq C_0 \frac{n^{1.5} y_{\max}}{y_{\min}}$  in Assumption 2, and the last inequality uses the following relationships between constants:  $C_0 > C \cdot C_\alpha^2$  and  $C_\alpha > C \cdot \max\{C_g^2, C_y C_g\}$ . For the lower bound, we have

$$\begin{aligned} \alpha_j^{(1)} &= \eta \left( y_j - \epsilon_j + \frac{\epsilon_j}{\eta \|\boldsymbol{\lambda}\|_1} + \frac{1}{\eta} \mathbf{e}_j^\top \left( \left( \mathbf{X} \mathbf{X}^\top \right)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right) \boldsymbol{\epsilon} \right) \\ &\geq \eta \left( -y_{\max} - \epsilon_j - \frac{1}{\eta} \left\| \left( \mathbf{X} \mathbf{X}^\top \right)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right\|_2 \|\boldsymbol{\epsilon}\|_2 \right) \\ &\geq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1} \left( -y_{\max} - \frac{1}{C_\alpha} y_{\min} - C^2 C_g^2 \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \cdot \frac{\sqrt{n}}{C_\alpha} y_{\min} \right) \\ &\geq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1} \left( -y_{\max} - \frac{1}{C_\alpha} y_{\min} - C^2 C_g^2 \cdot \frac{y_{\min}}{C_0 y_{\max}} \cdot \frac{1}{C_\alpha} y_{\min} \right) \\ &\geq \frac{-3y_{\max}}{C_g \|\boldsymbol{\lambda}\|_1}, \end{aligned}$$

by the same arguments. Thus,  $\alpha_j^{(1)}$  satisfies both the required upper and lower bounds for all  $j$  with  $y_j < 0$ .

Part (c): We now verify that the primal variables corresponding to positively labeled examples minus  $\mathbf{y}_+$  satisfy the norm bound in Lemma 17. Specifically, we show that

$\left\| \boldsymbol{\beta}_+^{(1)} - \mathbf{y}_+ \right\|_2^2 \leq C_y^2 \|\mathbf{y}_+\|_2^2$ . According to Equation (22), we have

$$\begin{aligned} \left\| \boldsymbol{\beta}_+^{(1)} - \mathbf{y}_+ \right\|_2^2 &= \sum_{i: y_i > 0} \left( \beta_i^{(1)} - y_i \right)^2 \\ &= \sum_{i: y_i > 0} \left( \underbrace{\epsilon_i - \eta \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top (\boldsymbol{\epsilon} - \mathbf{y}) - y_i}_{=: T_i} \right)^2. \end{aligned} \quad (23)$$

Next, we bound the term  $T_i := \epsilon_i - \eta \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top (\boldsymbol{\epsilon} - \mathbf{y}) - y_i$  for all  $i \in [n]$  with  $y_i > 0$ . We have

$$\begin{aligned} T_i &= \epsilon_i - \eta \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top (\boldsymbol{\epsilon} - \mathbf{y}) - y_i \\ &= (\epsilon_i - y_i) - \eta \mathbf{e}_i^\top \left[ \|\boldsymbol{\lambda}\|_1 \mathbf{I} + \left( \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right) \right] (\boldsymbol{\epsilon} - \mathbf{y}) \\ &= (1 - \eta \|\boldsymbol{\lambda}\|_1) (\epsilon_i - y_i) - \eta \mathbf{e}_i^\top \left( \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right) (\boldsymbol{\epsilon} - \mathbf{y}). \end{aligned}$$

Since the step size assumption guarantees that  $\frac{1}{CC_g \|\boldsymbol{\lambda}\|_1} \leq \eta \leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1}$ , and  $\epsilon_i \leq \frac{1}{C_\alpha} y_{\min}$ , the term  $(1 - \eta \|\boldsymbol{\lambda}\|_1) (\epsilon_i - y_i)$  is strictly negative. Hence, in order to upper bound  $T_i^2$ , it suffices to find the lower bound for  $T_i$ . We have

$$\begin{aligned} T_i &= (1 - \eta \|\boldsymbol{\lambda}\|_1) (\epsilon_i - y_i) - \eta \mathbf{e}_i^\top \left( \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right) (\boldsymbol{\epsilon} - \mathbf{y}) \\ &\geq -y_i - \eta \left\| \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right\|_2 \|\boldsymbol{\epsilon} - \mathbf{y}\|_2, \end{aligned}$$

where the inequality drops the positive terms  $(1 - \eta \|\boldsymbol{\lambda}\|_1) \epsilon_i$  and  $\eta \|\boldsymbol{\lambda}\|_1 y_i$ . We again upper bound  $\left\| \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right\|_2$  by Corollary 14. With probability at least  $1 - 2 \exp(-n(Cc - \ln 9))$ , we have

$$\begin{aligned} T_i &\geq -y_i - \eta \cdot C \|\boldsymbol{\lambda}\|_1 \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \|\boldsymbol{\epsilon} - \mathbf{y}\|_2 \\ &\geq -y_i - \frac{C}{C_g} \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \|\boldsymbol{\epsilon} - \mathbf{y}\|_2, \end{aligned}$$

by applying  $\eta \leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1}$ . Finally, we apply the upper bounds for  $\|\boldsymbol{\epsilon}\|_2$  and  $\|\mathbf{y}\|_2$ , and Assumption 2 ensures that  $d_2 \geq C_0^2 \frac{n^2 y_{\max}^2}{y_{\min}^2}$  and  $d_\infty \geq C_0 \frac{n^{1.5} y_{\max}}{y_{\min}}$ . We have

$$\begin{aligned} T_i &\geq -y_i - \frac{C}{C_g} \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \left( \frac{\sqrt{n}}{C_\alpha} y_{\min} + \sqrt{n} y_{\max} \right) \\ &\geq -y_i - \frac{C y_{\min}}{C_g C_0 y_{\max}} \left( \frac{1}{C_\alpha} y_{\min} + y_{\max} \right) \\ &\geq -y_i \left( 1 + \frac{2C}{C_g C_0} \right) \\ &\geq -C_y y_i, \end{aligned}$$

with the choice of  $C_y \geq 2$ . Substituting  $T_i^2 \leq C_y^2 y_i^2$  into Equation (23), we have

$$\left\| \boldsymbol{\beta}_+^{(1)} - \mathbf{y}_+ \right\|_2^2 \leq \sum_{i: y_i > 0} C_y^2 y_i^2 = C_y^2 \|\mathbf{y}_+\|_2^2.$$

As a result, we conclude that  $\left\| \boldsymbol{\beta}_+^{(1)} - \mathbf{y}_+ \right\|_2 \leq C_y \|\mathbf{y}_+\|_2$  as required.

Part (d): We next verify that  $\alpha^{(1)}$  satisfies the required norm bound. Recall that

$$\alpha^{(1)} = \eta \left( \mathbf{y} - \epsilon + \frac{1}{\eta} (\mathbf{X} \mathbf{X}^\top)^{-1} \epsilon \right).$$

Taking the  $\ell_2$  norm and applying the triangle inequality yields

$$\begin{aligned} \|\alpha^{(1)}\|_2 &= \left\| \eta \left( \mathbf{y} - \epsilon + \frac{1}{\eta} (\mathbf{X} \mathbf{X}^\top)^{-1} \epsilon \right) \right\|_2 \\ &\leq \eta \left[ \|\mathbf{y}\|_2 + \|\epsilon\|_2 + \frac{1}{\eta} \left\| (\mathbf{X} \mathbf{X}^\top)^{-1} \right\|_2 \|\epsilon\|_2 \right]. \end{aligned}$$

We now bound each term on the right-hand side. We apply the label bound,  $\|\mathbf{y}\|_2 \leq \sqrt{n} y_{\max}$  and the construction of the initialization,  $\|\epsilon\|_2 \leq \frac{\sqrt{n}}{C_\alpha} y_{\min}$ . Moreover, Lemma 12 implies  $\left\| (\mathbf{X} \mathbf{X}^\top)^{-1} \right\|_2 \leq \frac{C_g}{\|\boldsymbol{\lambda}\|_1}$  with probability at least  $1 - 2e^{-n/C_g}$ , and the step size condition ensures  $\frac{1}{C_g \|\boldsymbol{\lambda}\|_1} \leq \eta \leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1}$ . Substituting these bounds, we obtain

$$\begin{aligned} \|\alpha^{(1)}\|_2 &\leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1} \left[ \sqrt{n} y_{\max} + \frac{\sqrt{n}}{C_\alpha} y_{\min} + C C_g \|\boldsymbol{\lambda}\|_1 \cdot \frac{C_g}{\|\boldsymbol{\lambda}\|_1} \cdot \frac{\sqrt{n}}{C_\alpha} y_{\min} \right] \\ &\leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1} (3\sqrt{n} y_{\max}) \\ &\leq \frac{C_\alpha \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1}, \end{aligned}$$

with  $C_\alpha \gtrsim \max\{C_g^2, C_y C_g\}$ . Therefore,  $\alpha^{(1)}$  satisfies the required norm bound.

Part (e): Since we have shown that  $\alpha_j^{(1)} \leq -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1}$  and  $\|\alpha^{(1)}\|_2 \leq \frac{C_\alpha \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1}$  for all  $j \in [n]$  with  $y_j < 0$ , it follows from Lemma 11 that  $\beta_j^{(1)} \leq 0$  for all  $j \in [n]$  with  $y_j < 0$ .

Part (f): Since we have shown that  $\beta_i^{(1)} > 0$  for all  $i \in [n]$  with  $y_i > 0$  and  $\beta_j^{(1)} \leq 0$  for all  $j \in [n]$  with  $y_j < 0$ , the signs of the primal variables coincide with the signs of the labels. Consequently,  $\sigma(\beta^{(1)}) = \begin{bmatrix} \beta_+^{(1)} \\ \mathbf{0} \end{bmatrix}$ .

We have shown that at iteration  $t = 1$ , all conditions in Lemma 17 are satisfied. Consequently, all positively labeled examples are active, while all negatively labeled examples are inactive. We now complete the proof by induction and characterize the gradient descent dynamics for all subsequent iterations. By Lemma 17, since the conditions hold at  $t = 1$ , they also hold at  $t = 2$ . More generally, the same lemma implies that if the conditions hold at  $t = \tilde{t}$  then they continue to hold at  $t = \tilde{t} + 1$ . This completes the induction argument.

As a result, for all  $t \geq 1$ , the activation pattern remains fixed, i.e.,  $D(\beta^{(t)}) = D(\beta^{(1)})$ , and all negative labeled examples are inactive, i.e.,  $\mathbf{X}_- \mathbf{w}^{(t)} \leq 0$ . By Lemma 1, the gradient descent dynamics from this point onward are equivalent to those of linear regression trained on the positively labeled examples, with initialization  $\mathbf{w}^{(1)}$ . Finally, by Lemma 2, the  $\mathbf{w}^{(\infty)}$  satisfies

$$\mathbf{w}^{(\infty)} = \arg \min_{\mathbf{w} \in \{\mathbf{w} : \mathbf{X}_+ \mathbf{w} = \mathbf{y}_+\}} \left\| \mathbf{w} - \mathbf{w}^{(1)} \right\|_2,$$

where we have  $\mathbf{w}^{(1)} = \eta \mathbf{X}^\top \left( \mathbf{y} - \boldsymbol{\epsilon} + \frac{1}{\eta} (\mathbf{X} \mathbf{X}^\top)^{-1} \boldsymbol{\epsilon} \right)$ . This completes the proof of Theorem 4. ■

### B.3. Proof of Theorem 6 (Implicit Bias Approximation to $\mathbf{w}^*$ )

In this section, we present the proof of implicit bias approximation to  $\mathbf{w}^*$  for single ReLU models.

**Proof** (Theorem 6) We restate the definition of  $\mathbf{w}^*$  in Equation (7) below.

$$\begin{aligned} \mathbf{w}^* &= \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t. } \mathbf{w}^\top \mathbf{x}_i &= y_i, \text{ for all } y_i > 0 \\ \mathbf{w}^\top \mathbf{x}_j &\leq 0, \text{ for all } y_j \leq 0. \end{aligned}$$

Recall that the gradient descent limit  $\mathbf{w}^{(\infty)}$  satisfies the same set of constraints: it interpolates all positively labeled examples and produces negative predictions for negatively labeled examples. Consequently, both  $\mathbf{w}^{(\infty)}$  and  $\mathbf{w}^*$  are feasible solutions achieving the minimum empirical risk.

We start with showing the upper bound on  $\|\mathbf{w}^{(\infty)} - \mathbf{w}^*\|_2$ . We first relate the distance between the predictors  $\mathbf{w}^{(\infty)}$  and  $\mathbf{w}^*$  to the distance in their predictions, i.e.,  $\|\mathbf{X} \mathbf{w}^{(\infty)} - \mathbf{X} \mathbf{w}^*\|_2$ . Since both vectors lie in the span of the data  $\{\mathbf{x}_i\}_{i=1}^n$ , their difference has no component in the null space corresponding to the smallest  $d - n$  eigenvalues of  $\mathbf{X}^\top \mathbf{X}$ . Therefore, we have

$$\begin{aligned} \|\mathbf{X} \mathbf{w}^{(\infty)} - \mathbf{X} \mathbf{w}^*\|_2^2 &= \|\mathbf{X} (\mathbf{w}^{(\infty)} - \mathbf{w}^*)\|_2^2 \geq \mu_n(\mathbf{X}^\top \mathbf{X}) \|\mathbf{w}^{(\infty)} - \mathbf{w}^*\|_2^2 \\ &= \mu_n(\mathbf{X} \mathbf{X}^\top) \|\mathbf{w}^{(\infty)} - \mathbf{w}^*\|_2^2. \end{aligned} \quad (24)$$

As a result, to derive an upper bound for  $\|\mathbf{w}^{(\infty)} - \mathbf{w}^*\|_2$ , it suffices to upper bound the distance between their prediction  $\|\mathbf{X} \mathbf{w}^{(\infty)} - \mathbf{X} \mathbf{w}^*\|_2$ . We begin with analyzing  $\mathbf{w}^{(\infty)}$ . By Theorem 4,  $\mathbf{w}^{(\infty)}$  satisfies the following:

$$\begin{aligned} \mathbf{w}^{(\infty)\top} \mathbf{x}_i &= y_i && \text{for all } y_i > 0, \\ \alpha_j^{(\infty)} = \alpha_j^{(1)} &= \eta \left( y_j - \epsilon_j + \frac{1}{\eta} \mathbf{e}_j^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \boldsymbol{\epsilon} \right) && \text{for all } y_j < 0, \end{aligned}$$

and also all the conditions in Lemma 17. On the other hand, according to the necessary conditions in Equation (19) in Lemma 3,  $\mathbf{w}^*$  satisfies

$$\begin{aligned} \mathbf{w}^* &= - \sum_{i \in S_1} \lambda_i^* \mathbf{x}_i - \sum_{j \in \tilde{S}_2} \mu_j^* \mathbf{x}_j, \text{ with } \lambda_i^* \in \mathbb{R} \text{ and } \mu_j^* > 0, \\ \mathbf{w}^{*\top} \mathbf{x}_i &= y_i, \text{ for all } i \in S_1, \\ \mathbf{w}^{*\top} \mathbf{x}_j &= 0, \text{ for all } j \in \tilde{S}_2, \end{aligned}$$

where we have denoted  $S_1 = \{i : y_i > 0, \text{ for all } i \in [n]\}$ ,  $S_2 = \{j : y_j \leq 0, \text{ for all } j \in [n]\}$ ,  $\tilde{S}_2 \subseteq S_2$  (note that  $\tilde{S}_2$  can be empty) and  $S = S_1 \cup \tilde{S}_2$ . Based on these necessary conditions, we

can define  $\mathbf{w}^* = \mathbf{X}^\top \boldsymbol{\alpha}^*$  where

$$\boldsymbol{\alpha}_i^* = \begin{cases} -\lambda_i^* & \text{for all } i \in S_1 \\ -\mu_i^* & \text{for all } i \in \tilde{S}_2 \\ 0 & \text{for all } i \in S_2 \cup \tilde{S}_2^c =: S_3 \end{cases}.$$

Let  $\mathbf{X}_S \in \mathbb{R}^{|S| \times d}$  denote the submatrix of  $\mathbf{X}$  consisting of the rows indexed by  $S$  (taken in increasing order), and let  $\mathbf{y}_S \in \mathbb{R}^{|S|}$  denote the corresponding label subvector with all negative entries replaced by zero. We have

$$\mathbf{y}_S = \mathbf{X}_S \mathbf{X}_S^\top \boldsymbol{\alpha}_S^*,$$

and similarly, by taking the norm and using the matrix norm lower bound of the smallest eigenvalue of  $\mathbf{X}_S \mathbf{X}_S^\top$ , we have

$$\begin{aligned} \|\mathbf{y}_S\|_2 &= \left\| \mathbf{X}_S \mathbf{X}_S^\top \boldsymbol{\alpha}_S^* \right\|_2 \\ &\geq \mu_{|S|}(\mathbf{X}_S \mathbf{X}_S^\top) \|\boldsymbol{\alpha}_S^*\|_2. \end{aligned}$$

Consequently, we have

$$\|\boldsymbol{\alpha}^*\|_2 = \|\boldsymbol{\alpha}_S^*\|_2 \leq \frac{\|\mathbf{y}_S\|_2}{\mu_{|S|}(\mathbf{X}_S \mathbf{X}_S^\top)} \leq \frac{\sqrt{n} y_{\max}}{\mu_n(\mathbf{X} \mathbf{X}^\top)} \leq \frac{C_g \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1}, \quad (25)$$

where the second inequality follows from the variational formulation of submatrix, and the last inequality follows from Lemma 12 with probability at least  $1 - 2e^{-n/C_g}$ .

We know that for all  $i \in S_1$ ,  $\mathbf{w}^{(\infty)\top} \mathbf{x}_i = \mathbf{w}^{*\top} \mathbf{x}_i = y_i$ , and  $\mathbf{w}^{*\top} \mathbf{x}_j = 0$  for all  $j \in \tilde{S}_2$ . Therefore, we can write

$$\begin{aligned} \left\| \mathbf{X} \mathbf{w}^{(\infty)} - \mathbf{X} \mathbf{w}^* \right\|_2^2 &= \sum_{i=1}^n \left( \mathbf{w}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}^{*\top} \mathbf{x}_i \right)^2 \\ &= \sum_{i \in \tilde{S}_2} \left( \mathbf{w}^{(\infty)\top} \mathbf{x}_i \right)^2 + \sum_{i \in S_3} \left( \mathbf{w}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}^{*\top} \mathbf{x}_i \right)^2. \end{aligned} \quad (26)$$

We start with upper bounding the term  $(\mathbf{w}^{(\infty)\top} \mathbf{x}_i)^2$  for all  $i \in \tilde{S}_2$ . Since  $\mathbf{w}^{(\infty)\top} \mathbf{x}_i \leq 0$  by the conditions in Lemma 17, it suffices to lower bound  $\mathbf{w}^{(\infty)\top} \mathbf{x}_i$ . We have

$$\begin{aligned} \mathbf{w}^{(\infty)\top} \mathbf{x}_i &= \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top \boldsymbol{\alpha}^{(\infty)} \\ &= \mathbf{e}_i^\top \left[ \|\boldsymbol{\lambda}\|_1 \mathbf{I} + (\mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}) \right] \boldsymbol{\alpha}^{(\infty)} \\ &= \|\boldsymbol{\lambda}\|_1 \alpha_i^{(\infty)} + \mathbf{e}_i^\top (\mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}) \boldsymbol{\alpha}^{(\infty)} \\ &\geq \|\boldsymbol{\lambda}\|_1 \alpha_i^{(\infty)} - \left\| \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right\|_2 \left\| \boldsymbol{\alpha}^{(\infty)} \right\|_2 \\ &\geq \|\boldsymbol{\lambda}\|_1 \left[ \alpha_i^{(\infty)} - C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \left\| \boldsymbol{\alpha}^{(\infty)} \right\|_2 \right], \end{aligned}$$

where the last inequality applies Corollary 14. Substituting the bounds of  $\alpha_i^{(\infty)}$  and  $\|\alpha^{(\infty)}\|_2$  from Lemma 17, we have

$$\begin{aligned} \mathbf{w}^{(\infty)\top} \mathbf{x}_i &\geq \|\lambda\|_1 \left[ -\frac{3y_{\max}}{C_g \|\lambda\|_1} - C \cdot \max\left(\sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty}\right) \frac{C_\alpha \sqrt{n} y_{\max}}{\|\lambda\|_1} \right] \\ &\geq \|\lambda\|_1 \left[ -\frac{3y_{\max}}{C_g \|\lambda\|_1} - C \cdot \frac{y_{\min}}{C_0 y_{\max}} \frac{C_\alpha y_{\max}}{\|\lambda\|_1} \right] \\ &\geq -\frac{4}{C_g} y_{\max}, \end{aligned}$$

where the inequalities above substitute  $d_2 \geq C_0^2 \frac{n^2 y_{\max}^2}{y_{\min}^2}$  and  $d_\infty \geq C_0 \frac{n^{1.5} y_{\max}}{y_{\min}}$  in Assumption 2 with  $C_0 \gtrsim C_\alpha^2$  and  $C_\alpha \gtrsim \max\{C_g^2, C_y C_g\}$ . Therefore, we have  $(\mathbf{w}^{(\infty)\top} \mathbf{x}_i)^2 \leq \frac{16}{C_g^2} y_{\max}^2$  for all  $i \in \tilde{S}_2$ . Next, we upper bound the term  $(\mathbf{w}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}^{\star\top} \mathbf{x}_i)^2$  for all  $i \in S_3$ . We use the key idea that  $\alpha_i^* = 0$  for all  $i \in S_3$ . We have

$$\begin{aligned} \mathbf{w}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}^{\star\top} \mathbf{x}_i &= \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top \alpha^{(\infty)} - \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top \alpha^* \\ &= \mathbf{e}_i^\top \left[ \|\lambda\|_1 \mathbf{I} + (\mathbf{X} \mathbf{X}^\top - \|\lambda\|_1 \mathbf{I}) \right] (\alpha^{(\infty)} - \alpha^*) \\ &= \|\lambda\|_1 \alpha_i^{(\infty)} + \mathbf{e}_i^\top (\mathbf{X} \mathbf{X}^\top - \|\lambda\|_1 \mathbf{I}) (\alpha^{(\infty)} - \alpha^*) \\ &\geq \|\lambda\|_1 \alpha_i^{(\infty)} - \|\mathbf{X} \mathbf{X}^\top - \|\lambda\|_1 \mathbf{I}\|_2 \left( \|\alpha^{(\infty)}\|_2 + \|\alpha^*\|_2 \right) \\ &\geq \|\lambda\|_1 \left[ -\frac{3y_{\max}}{C_g \|\lambda\|_1} - C \cdot \max\left(\sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty}\right) \left( \frac{C_\alpha \sqrt{n} y_{\max}}{\|\lambda\|_1} + \frac{C_g \sqrt{n} y_{\max}}{\|\lambda\|_1} \right) \right] \\ &\geq \|\lambda\|_1 \left[ -\frac{3y_{\max}}{C_g \|\lambda\|_1} - C \cdot \frac{y_{\min}}{C_0 y_{\max}} \left( \frac{C_\alpha y_{\max}}{\|\lambda\|_1} + \frac{C_g y_{\max}}{\|\lambda\|_1} \right) \right] \\ &\geq -\frac{4}{C_g} y_{\max}, \end{aligned}$$

by applying the same argument and noting from Equation (25) that  $\|\alpha^*\|_2 \leq \frac{C_g \sqrt{n} y_{\max}}{\|\lambda\|_1}$ . Substituting the upper bounds into Equation (26) gives us

$$\begin{aligned} \|\mathbf{X} \mathbf{w}^{(\infty)} - \mathbf{X} \mathbf{w}^*\|_2^2 &= \sum_{i \in \tilde{S}_2} \left( \mathbf{w}^{(\infty)\top} \mathbf{x}_i \right)^2 + \sum_{i \in S_3} \left( \mathbf{w}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}^{\star\top} \mathbf{x}_i \right)^2 \\ &\leq \sum_{i \in \tilde{S}_2} \frac{16}{C_g^2} y_{\max}^2 + \sum_{i \in S_3} \frac{16}{C_g^2} y_{\max}^2 \\ &= \frac{16}{C_g^2} n - y_{\max}^2. \end{aligned} \tag{27}$$

Finally, putting together Equation (24) and (27), we have

$$\|\mathbf{w}^{(\infty)} - \mathbf{w}^*\|_2^2 \leq \frac{\|\mathbf{X} \mathbf{w}^{(\infty)} - \mathbf{X} \mathbf{w}^*\|_2^2}{\mu_n(\mathbf{X} \mathbf{X}^\top)} \leq \frac{16n - y_{\max}^2}{C_g \|\lambda\|_1},$$

which completes the proof of the upper bound. Next, we derive the lower bound of  $\|\mathbf{w}^{(\infty)} - \mathbf{w}^*\|_2$  in a similar approach. We again start with the prediction distance, given by

$$\begin{aligned} \|\mathbf{X}\mathbf{w}^{(\infty)} - \mathbf{X}\mathbf{w}^*\|_2^2 &= \|\mathbf{X}(\mathbf{w}^{(\infty)} - \mathbf{w}^*)\|_2^2 \leq \mu_1(\mathbf{X}^\top \mathbf{X}) \|\mathbf{w}^{(\infty)} - \mathbf{w}^*\|_2^2 \\ &= \mu_1(\mathbf{X}\mathbf{X}^\top) \|\mathbf{w}^{(\infty)} - \mathbf{w}^*\|_2^2. \end{aligned} \quad (28)$$

It suffices to lower bound  $\|\mathbf{X}\mathbf{w}^{(\infty)} - \mathbf{X}\mathbf{w}^*\|_2$  to get the lower bound of  $\|\mathbf{w}^{(\infty)} - \mathbf{w}^*\|_2$ . By Equation (26), we have

$$\|\mathbf{X}\mathbf{w}^{(\infty)} - \mathbf{X}\mathbf{w}^*\|_2^2 = \sum_{i \in \tilde{S}_2} (\mathbf{w}^{(\infty)\top} \mathbf{x}_i)^2 + \sum_{i \in S_3} (\mathbf{w}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}^{*\top} \mathbf{x}_i)^2.$$

Therefore, we need to lower bound  $(\mathbf{w}^{(\infty)\top} \mathbf{x}_i)^2$  for  $i \in \tilde{S}_2$ , and  $(\mathbf{w}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}^{*\top} \mathbf{x}_i)^2$  for  $i \in S_3$ . For  $\mathbf{w}^{(\infty)\top} \mathbf{x}_i$ , since  $\mathbf{w}^{(\infty)\top} \mathbf{x}_i < 0$ , we have

$$\begin{aligned} \mathbf{w}^{(\infty)\top} \mathbf{x}_i &= \mathbf{e}_i^\top \mathbf{X}\mathbf{X}^\top \boldsymbol{\alpha}^{(\infty)} \\ &= \mathbf{e}_i^\top \left[ \|\boldsymbol{\lambda}\|_1 \mathbf{I} + (\mathbf{X}\mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}) \right] \boldsymbol{\alpha}^{(\infty)} \\ &= \|\boldsymbol{\lambda}\|_1 \alpha_i^{(\infty)} + \mathbf{e}_i^\top (\mathbf{X}\mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}) \boldsymbol{\alpha}^{(\infty)} \\ &\leq \|\boldsymbol{\lambda}\|_1 \alpha_i^{(\infty)} + \left\| \mathbf{X}\mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right\|_2 \left\| \boldsymbol{\alpha}^{(\infty)} \right\|_2 \\ &\leq \|\boldsymbol{\lambda}\|_1 \left[ \alpha_i^{(\infty)} + C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \left\| \boldsymbol{\alpha}^{(\infty)} \right\|_2 \right], \end{aligned}$$

where the last inequality applies Corollary 14. Substituting the bounds of  $\alpha_i^{(\infty)}$  and  $\left\| \boldsymbol{\alpha}^{(\infty)} \right\|_2$  from Lemma 17, we have

$$\begin{aligned} \mathbf{w}^{(\infty)\top} \mathbf{x}_i &\leq \|\boldsymbol{\lambda}\|_1 \left[ -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1} + C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \frac{C_\alpha \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1} \right] \\ &\leq \|\boldsymbol{\lambda}\|_1 \left[ -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1} + C \cdot \frac{y_{\min}}{C_0 y_{\max}} \frac{C_\alpha y_{\max}}{\|\boldsymbol{\lambda}\|_1} \right] \\ &\leq -\left(1 - \frac{C \cdot C_\alpha^2}{C_0}\right) \frac{y_{\min}}{C_\alpha}, \end{aligned}$$

where the inequalities above substitute  $d_2 \geq C_0^2 \frac{n^2 y_{\max}^2}{y_{\min}^2}$  and  $d_\infty \geq C_0 \frac{n^{1.5} y_{\max}}{y_{\min}}$  in Assumption 2 with  $C_0 \gtrsim C_\alpha^2$ . Similarly, for  $\mathbf{w}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}^{\star\top} \mathbf{x}_i$ , we have

$$\begin{aligned}
\mathbf{w}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}^{\star\top} \mathbf{x}_i &= \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top \boldsymbol{\alpha}^{(\infty)} - \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top \boldsymbol{\alpha}^\star \\
&= \mathbf{e}_i^\top \left[ \|\boldsymbol{\lambda}\|_1 \mathbf{I} + (\mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}) \right] (\boldsymbol{\alpha}^{(\infty)} - \boldsymbol{\alpha}^\star) \\
&= \|\boldsymbol{\lambda}\|_1 \alpha_i^{(\infty)} + \mathbf{e}_i^\top (\mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}) (\boldsymbol{\alpha}^{(\infty)} - \boldsymbol{\alpha}^\star) \\
&\leq \|\boldsymbol{\lambda}\|_1 \alpha_i^{(\infty)} + \left\| \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right\|_2 (\|\boldsymbol{\alpha}^{(\infty)}\|_2 + \|\boldsymbol{\alpha}^\star\|_2) \\
&\leq \|\boldsymbol{\lambda}\|_1 \left[ -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1} + C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \left( \frac{C_\alpha \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1} + \frac{C_g \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1} \right) \right] \\
&\leq \|\boldsymbol{\lambda}\|_1 \left[ -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1} + C \cdot \frac{y_{\min}}{C_0 y_{\max}} \left( \frac{C_\alpha y_{\max}}{\|\boldsymbol{\lambda}\|_1} + \frac{C_g y_{\max}}{\|\boldsymbol{\lambda}\|_1} \right) \right] \\
&\leq -\left(1 - \frac{2C \cdot C_\alpha^2}{C_0}\right) \frac{y_{\min}}{C_\alpha},
\end{aligned}$$

by applying the same argument and noting from Equation (25) that  $\|\boldsymbol{\alpha}^\star\|_2 \leq \frac{C_g \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1}$ . Substituting the lower bounds into Equation (26) gives us

$$\begin{aligned}
\left\| \mathbf{X} \mathbf{w}^{(\infty)} - \mathbf{X} \mathbf{w}^\star \right\|_2^2 &= \sum_{i \in \tilde{S}_2} \left( \mathbf{w}^{(\infty)\top} \mathbf{x}_i \right)^2 + \sum_{i \in S_3} \left( \mathbf{w}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}^{\star\top} \mathbf{x}_i \right)^2 \\
&\geq \sum_{i \in \tilde{S}_2} \left(1 - \frac{C \cdot C_\alpha^2}{C_0}\right)^2 \frac{y_{\min}^2}{C_\alpha^2} + \sum_{i \in S_3} \left(1 - \frac{2C \cdot C_\alpha^2}{C_0}\right)^2 \frac{y_{\min}^2}{C_\alpha^2} \\
&\geq \left(1 - \frac{2C \cdot C_\alpha^2}{C_0}\right)^2 \frac{n - y_{\min}^2}{C_\alpha^2} \\
&= \frac{n - y_{\min}^2}{\tilde{C}}, \tag{29}
\end{aligned}$$

where we let  $\tilde{C} := \frac{C_0^2 C_\alpha^2}{(C_0 - 2C \cdot C_\alpha^2)^2} > 1$ . Finally, putting together Equations (28) and (29), we have

$$\left\| \mathbf{w}^{(\infty)} - \mathbf{w}^\star \right\|_2^2 \geq \frac{\left\| \mathbf{X} \mathbf{w}^{(\infty)} - \mathbf{X} \mathbf{w}^\star \right\|_2^2}{\mu_1(\mathbf{X} \mathbf{X}^\top)} \geq \frac{n - y_{\min}^2}{\tilde{C} C_g \|\boldsymbol{\lambda}\|_1}.$$

This completes the proof of the lower bound. ■

#### B.4. Local Minimum Convergence Under a Not-All-Positive Initialization

In this section, we present a simple and explicit counterexample showing that a single ReLU model initialized with a not-all-positive initialization can converge to a local minimum that is not global.

##### Lemma 18 (Local minimum convergence for not-all-positive initialization)

Assume Assumptions 1 and 2 hold. For a single ReLU model, a dataset  $n = 2$  and  $y_i > 0$  for all

$i = 1, 2$ , we choose the initialization  $\mathbf{w}^{(0)} = \mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \boldsymbol{\epsilon}$ , where  $\epsilon_1 = -\delta$  and  $\epsilon_2 = \delta$ , for some  $0 < \delta \leq \frac{1}{C} y_{\min}$ . If  $\mathbf{x}_1^\top \mathbf{x}_2 < 0$  and  $|\mathbf{x}_1^\top \mathbf{x}_2| < \|\mathbf{x}_2\|_2^2$ , with a step size  $\eta < \frac{1}{\mu_1(\mathbf{X}\mathbf{X}^\top)}$ , gradient descent converges to a local minimum.

**Proof** To show that gradient descent converges to a local minimum, it suffices to prove that  $\beta_1^{(t)} < 0$  for all  $t \geq 0$ : since  $y_1 > 0$ , this condition implies that the first training example is not perfectly fit. Hence, we prove by induction that if  $\beta_1^{(t)} < 0$  for all  $t \leq \tilde{t}$ , then  $\beta_1^{(t+1)} < 0$ . Note that the base case, i.e.  $\beta_1^{(0)} < 0$ , follows from the initialization condition  $\beta_1^{(0)} = \epsilon_1 = -\delta < 0$ . Therefore, we only need to show the inductive step. From the primal gradient update in Equation (4a), we obtain the gradient update of  $\beta_1^{(\tilde{t}+1)}$  as

$$\begin{aligned} \beta_1^{(\tilde{t}+1)} &= \beta_1^{(\tilde{t})} - \eta \|\mathbf{x}_1\|_2^2 \cdot \underbrace{\mathbb{1}_{\beta_1^{(\tilde{t})} > 0}}_{=0} \cdot (\beta_1^{(\tilde{t})} - y_1) - \eta \mathbf{x}_1^\top \mathbf{x}_2 \cdot \mathbb{1}_{\beta_2^{(\tilde{t})} > 0} \cdot (\beta_2^{(\tilde{t})} - y_2) \\ &= \beta_1^{(0)} - \eta \mathbf{x}_1^\top \mathbf{x}_2 \sum_{t=0}^{\tilde{t}} \mathbb{1}_{\beta_2^{(t)} > 0} (\beta_2^{(t)} - y_2). \end{aligned} \quad (30)$$

On the other hand, for  $\beta_2^{(t)}$ , since  $\beta_1^{(t)} < 0$  for all  $t \leq \tilde{t}$ , we have

$$\beta_2^{(t+1)} = \beta_2^{(t)} - \eta \|\mathbf{x}_2\|_2^2 \mathbb{1}_{\beta_2^{(t)} > 0} (\beta_2^{(t)} - y_2) = \beta_2^{(t)} (1 - \eta \|\mathbf{x}_2\|_2^2 \mathbb{1}_{\beta_2^{(t)} > 0}) + \eta \|\mathbf{x}_2\|_2^2 y_2 \mathbb{1}_{\beta_2^{(t)} > 0},$$

for all  $t \leq \tilde{t}$ . Since  $\eta \|\mathbf{x}_2\|_2^2 < \frac{\|\mathbf{x}_2\|_2^2}{\mu_1(\mathbf{X}\mathbf{X}^\top)} \leq 1$ , we have  $\beta_2^{(t+1)} > 0$  if  $\beta_2^{(t)} > 0$ . Since  $\beta_2^{(0)} = \delta > 0$ , we have shown that  $\beta_2^{(t)} > 0$  for all  $t \leq \tilde{t}$ . Therefore, Equation (30) becomes

$$\beta_1^{(\tilde{t}+1)} = \beta_1^{(0)} - \eta \mathbf{x}_1^\top \mathbf{x}_2 \sum_{t=0}^{\tilde{t}} (\beta_2^{(t)} - y_2). \quad (31)$$

Next, we show an upper bound for  $\beta_2^{(t)}$ . Since  $\beta_1^{(t)} < 0$  and  $\beta_2^{(t)} > 0$  for all  $t \leq \tilde{t}$ , we can show that the gradient update for  $\beta_2^{(t)}$  satisfies

$$\begin{aligned} \beta_2^{(t)} &= \beta_2^{(t-1)} - \eta \|\mathbf{x}_2\|_2^2 (\beta_2^{(t-1)} - y_2) \\ &= \beta_2^{(t-1)} (1 - \eta \|\mathbf{x}_2\|_2^2) + \eta \|\mathbf{x}_2\|_2^2 y_2 \\ &= \beta_2^{(0)} (1 - \eta \|\mathbf{x}_2\|_2^2)^t + \eta \|\mathbf{x}_2\|_2^2 y_2 \sum_{k=0}^{t-1} (1 - \eta \|\mathbf{x}_2\|_2^2)^k \end{aligned}$$

for all  $t \leq \tilde{t}$ . Again, since  $\eta \|\mathbf{x}_2\|_2^2 < 1$ , the geometric series yields

$$\sum_{k=0}^{t-1} (1 - \eta \|\mathbf{x}_2\|_2^2)^k \leq \frac{1}{\eta \|\mathbf{x}_2\|_2^2},$$

and therefore, we have

$$\beta_2^{(t)} \leq \beta_2^{(0)} (1 - \eta \|\mathbf{x}_2\|_2^2)^t + y_2. \quad (32)$$

Substituting Equation (32) into Equation (31) and using the assumption that  $\mathbf{x}_1^\top \mathbf{x}_2 < 0$ , we obtain

$$\begin{aligned} \beta_1^{(\tilde{t}+1)} &\leq \beta_1^{(0)} - \eta \mathbf{x}_1^\top \mathbf{x}_2 \sum_{t=0}^{\tilde{t}} \beta_2^{(0)} (1 - \eta \|\mathbf{x}_2\|_2^2)^t \\ &\leq \beta_1^{(0)} - \frac{\mathbf{x}_1^\top \mathbf{x}_2}{\|\mathbf{x}_2\|_2^2} \beta_2^{(0)} \\ &= -\delta + \frac{|\mathbf{x}_1^\top \mathbf{x}_2|}{\|\mathbf{x}_2\|_2^2} \delta < 0, \end{aligned}$$

where the second inequality again follows from the bound on the geometric series, and the last inequality uses the assumption that  $|\mathbf{x}_1^\top \mathbf{x}_2| < \|\mathbf{x}_2\|_2^2$ . This completes the proof of the inductive step and therefore the proof that  $\beta_1^{(t)} < 0$  for all  $t \geq 0$ , implying that gradient descent remains stuck in a local minimum.  $\blacksquare$

According to Lemma 18, a simple example is given by i.i.d. Gaussian vectors  $\mathbf{x}_1, \mathbf{x}_2 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  for  $n = 2$ . By symmetry, the event  $\mathbf{x}_1^\top \mathbf{x}_2 < 0$  occurs with constant probability. Moreover, in the high-dimensional regime ( $d \gg n$ ), concentration implies that  $|\mathbf{x}_1^\top \mathbf{x}_2| < \|\mathbf{x}_2\|_2^2$  with a high probability.

## Appendix C. Proofs for the Two ReLU Model ( $m = 2$ ) Trained with Gradient Descent

In this section, we present the proofs concerning the behavior of the 2-ReLU model trained with gradient descent.

### C.1. Proof of Lemma 7 (Characterization of $w^*$ )

**Proof** (Lemma 7)

We first show that the feasible set of (8) is nonempty. Define  $\tilde{w}_\oplus := \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{y}_\oplus$  and  $\tilde{w}_\ominus := \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{y}_\ominus$ , where we define  $y_{\oplus,i} := \max\{y_i, 0\}$  and  $y_{\ominus,i} := -\min\{y_i, 0\}$ . Then for all  $i \in [n]$ , we have  $\sigma(\tilde{w}_\oplus^\top \mathbf{x}_i) - \sigma(\tilde{w}_\ominus^\top \mathbf{x}_i) = \sigma(y_{\oplus,i}) - \sigma(y_{\ominus,i}) = y_i$ . Thus,  $\{\tilde{w}_\oplus, \tilde{w}_\ominus\}$  is feasible, and the feasible set is nonempty.

Next, we show that any optimal solution of (8) corresponds to an optimal solution of (9). Let  $\{w_\oplus^*, w_\ominus^*\}$  be an optimal solution of (8).

**Case 1:  $i \in S_+$  (positive labels)**

For  $i \in S_+$ , since  $\sigma(w_\oplus^{*\top} \mathbf{x}_i) - \sigma(w_\ominus^{*\top} \mathbf{x}_i) = y_i > 0$ , we have

$$\sigma(w_\oplus^{*\top} \mathbf{x}_i) = y_i + \sigma(w_\ominus^{*\top} \mathbf{x}_i) \geq y_i > 0.$$

Hence,  $w_\oplus^{*\top} \mathbf{x}_i > 0$  and  $\sigma(w_\oplus^{*\top} \mathbf{x}_i) = w_\oplus^{*\top} \mathbf{x}_i$ . There are two possible activation patterns:

- If  $w_\ominus^{*\top} \mathbf{x}_i \leq 0$ , then we have  $\sigma(w_\oplus^{*\top} \mathbf{x}_i) - \sigma(w_\ominus^{*\top} \mathbf{x}_i) = w_\oplus^{*\top} \mathbf{x}_i = y_i$ .
- If  $w_\ominus^{*\top} \mathbf{x}_i \geq 0$ , then we have  $\sigma(w_\oplus^{*\top} \mathbf{x}_i) - \sigma(w_\ominus^{*\top} \mathbf{x}_i) = w_\oplus^{*\top} \mathbf{x}_i - w_\ominus^{*\top} \mathbf{x}_i = y_i$ .

(Note that  $w_\ominus^{*\top} \mathbf{x}_i = 0$  is covered by both cases.)

**Case 2:  $i \in S_-$  (negative labels)**

For  $i \in S_-$ , since  $\sigma(w_\oplus^{*\top} \mathbf{x}_i) - \sigma(w_\ominus^{*\top} \mathbf{x}_i) = y_i < 0$ , we obtain

$$\sigma(w_\oplus^{*\top} \mathbf{x}_i) = -y_i + \sigma(w_\ominus^{*\top} \mathbf{x}_i) \geq -y_i > 0,$$

which implies  $w_\oplus^{*\top} \mathbf{x}_i > 0$  and  $\sigma(w_\oplus^{*\top} \mathbf{x}_i) = w_\oplus^{*\top} \mathbf{x}_i$ . Again, two activation patterns are possible:

- If  $w_\oplus^{*\top} \mathbf{x}_i \leq 0$ , then we have  $\sigma(w_\oplus^{*\top} \mathbf{x}_i) - \sigma(w_\ominus^{*\top} \mathbf{x}_i) = -w_\ominus^{*\top} \mathbf{x}_i = y_i$ .
- If  $w_\oplus^{*\top} \mathbf{x}_i \geq 0$ , then we have  $\sigma(w_\oplus^{*\top} \mathbf{x}_i) - \sigma(w_\ominus^{*\top} \mathbf{x}_i) = w_\oplus^{*\top} \mathbf{x}_i - w_\ominus^{*\top} \mathbf{x}_i = y_i$ .

(Note that  $w_\oplus^{*\top} \mathbf{x}_i = 0$  is covered by both cases.)

Combining the two cases (in total four patterns), there exist disjoint partitions

$$S_1 \cup S_2 = S_+, S_1 \cap S_2 = \emptyset, \text{ and } S_3 \cup S_4 = S_-, S_3 \cap S_4 = \emptyset,$$

such that the optimal solution  $\{w_\oplus^*, w_\ominus^*\}$  satisfies

$$\begin{aligned} w_\oplus^{*\top} \mathbf{x}_i &= y_i, & w_\ominus^{*\top} \mathbf{x}_i &\leq 0, & \text{for all } i \in S_1, \\ w_\oplus^{*\top} \mathbf{x}_i - w_\ominus^{*\top} \mathbf{x}_i &= y_i, & -w_\ominus^{*\top} \mathbf{x}_i &\leq 0, & \text{for all } i \in S_2, \\ -w_\ominus^{*\top} \mathbf{x}_i &= y_i, & w_\oplus^{*\top} \mathbf{x}_i &\leq 0, & \text{for all } i \in S_3, \\ w_\oplus^{*\top} \mathbf{x}_i - w_\ominus^{*\top} \mathbf{x}_i &= y_i, & -w_\ominus^{*\top} \mathbf{x}_i &\leq 0, & \text{for all } i \in S_4. \end{aligned}$$

These constraints are exactly those in (9). Moreover, the feasible set of (9) is a subset of the feasible set of (8), since every feasible solution of (9) also satisfies the constraints of (8) (the converse need not hold). Since  $\{\mathbf{w}_\oplus^*, \mathbf{w}_\ominus^*\}$  is feasible for both problems and is optimal for the larger feasible set (8), it must also be optimal for the restricted problem (9).  $\blacksquare$

## C.2. Proof of Theorem 8 (High-dimensional Implicit Bias)

In this section, we present the proof of Theorem 8. For the 2-ReLU model ( $m = 2$ ), the primal-dual gradient update in (4) simplifies to

$$\text{(Primal)} \quad \beta_{\oplus}^{(t+1)} = \beta_{\oplus}^{(t)} - \eta \mathbf{X} \mathbf{X}^\top \mathbf{D}(\beta_{\oplus}^{(t)})(h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}), \quad (33a)$$

$$\text{(Dual)} \quad \alpha_{\oplus}^{(t+1)} = \alpha_{\oplus}^{(t)} - \eta \mathbf{D}(\beta_{\oplus}^{(t)})(h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}), \quad (33b)$$

and

$$\text{(Primal)} \quad \beta_{\ominus}^{(t+1)} = \beta_{\ominus}^{(t)} + \eta \mathbf{X} \mathbf{X}^\top \mathbf{D}(\beta_{\ominus}^{(t)})(h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}), \quad (34a)$$

$$\text{(Dual)} \quad \alpha_{\ominus}^{(t+1)} = \alpha_{\ominus}^{(t)} + \eta \mathbf{D}(\beta_{\ominus}^{(t)})(h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}). \quad (34b)$$

Before proceeding to the proof, we again introduce a set of sufficient conditions under which the signs of the primal variables agree with the signs of the labels times the sign of the ReLU neuron at iteration  $t$ , and moreover, these conditions are preserved at iteration  $t + 1$ . We use the results of Lemma 10 and Lemma 11 again to prove Lemma 19.

**Lemma 19** *Under Assumptions 1 and 2, suppose the gradient descent step size satisfies  $\eta \leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1}$ . For a 2-ReLU model, if the following eight conditions hold at some iteration  $t \geq 0$ , then they also hold at iteration  $t + 1$ .*

- a.  $\beta_{\oplus,i}^{(t)} > 0$  for all  $i \in [n]$  with  $y_i > 0$ .
- b.  $\beta_{\ominus,j}^{(t)} > 0$  for all  $j \in [n]$  with  $y_j < 0$ .
- c.  $-\frac{3y_{\max}}{C_g \|\boldsymbol{\lambda}\|_1} \leq \alpha_{\oplus,j}^{(t)} \leq -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1}$  for all  $j \in [n]$  with  $y_j < 0$ .
- d.  $-\frac{3y_{\max}}{C_g \|\boldsymbol{\lambda}\|_1} \leq \alpha_{\ominus,i}^{(t)} \leq -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1}$  for all  $i \in [n]$  with  $y_i > 0$ .
- e.  $\|\beta_{\oplus,+}^{(t)} - \mathbf{y}_+\|_2 \leq C_y \|\mathbf{y}_+\|_2$ , and  $\|\beta_{\ominus,-}^{(t)} + \mathbf{y}_-\|_2 \leq C_y \|\mathbf{y}_-\|_2$ .
- f.  $\|\boldsymbol{\alpha}_{\oplus}^{(t)}\|_2 \leq \frac{C_\alpha \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1}$  and  $\|\boldsymbol{\alpha}_{\ominus}^{(t)}\|_2 \leq \frac{C_\alpha \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1}$ .
- g.  $\beta_{\oplus,j}^{(t)} \leq 0$  for all  $j \in [n]$  with  $y_j < 0$ .
- h.  $\beta_{\ominus,i}^{(t)} \leq 0$  for all  $i \in [n]$  with  $y_i > 0$ .

Consequently, the set of active examples consists exactly of the positively labeled examples for the positive neuron, and the activation pattern remains unchanged, i.e.,  $\mathbf{D}(\beta_{\oplus}^{(t)}) = \mathbf{D}(\beta_{\oplus}^{(t+1)})$ . The set of active examples consists exactly of the negatively labeled examples for the negative neuron, and the activation pattern remains unchanged, i.e.,  $\mathbf{D}(\beta_{\ominus}^{(t)}) = \mathbf{D}(\beta_{\ominus}^{(t+1)})$ .

**Proof** (Lemma 19) We now verify that these conditions are preserved from iteration  $t$  to  $t + 1$ .

Part (a): By conditions (a), (b), (e), (g) and (h) at iteration  $t$ , we have

$$\begin{aligned} \|h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}\|_2^2 &= \left\| \sigma(\boldsymbol{\beta}_{\oplus}^{(t)}) - \sigma(\boldsymbol{\beta}_{\ominus}^{(t)}) - \mathbf{y} \right\|_2^2 \\ &= \left\| \sigma(\boldsymbol{\beta}_{\oplus}^{(t)}) - \begin{bmatrix} \mathbf{y}_+ \\ \mathbf{0} \end{bmatrix} - (\sigma(\boldsymbol{\beta}_{\ominus}^{(t)}) + \begin{bmatrix} \mathbf{0} \\ \mathbf{y}_- \end{bmatrix}) \right\|_2^2 \\ &= \left\| \boldsymbol{\beta}_{\oplus,+}^{(t)} - \mathbf{y}_+ \right\|_2^2 + \left\| \boldsymbol{\beta}_{\ominus,-}^{(t)} + \mathbf{y}_- \right\|_2^2 \leq C_y^2 \|\mathbf{y}\|_2^2, \end{aligned}$$

and therefore,  $\|h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}\|_2 \leq C_y \|\mathbf{y}\|_2$ . Together with  $h_{\Theta^{(t)}}(\mathbf{x}_i) = \beta_{\oplus,i}^{(t)}$  and condition (a), the assumptions of Lemma 10 are satisfied for all  $i$  with  $y_i > 0$ . Consequently,  $\beta_{\oplus,i}^{(t+1)} > 0$  for all  $i \in [n]$  with  $y_i > 0$ .

Part (b): According to the proof of part (a), we have  $\|h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}\|_2 \leq C_y \|\mathbf{y}\|_2$  and  $-h_{\Theta^{(t)}}(\mathbf{x}_j) = \beta_{\ominus,j}^{(t)}$ . Together with condition (b), the assumptions of Lemma 10 are satisfied for all  $j$  with  $y_j < 0$ . Consequently, we have  $\beta_{\ominus,j}^{(t+1)} > 0$  for all  $j \in [n]$  with  $y_j < 0$ .

Part (c): According to the dual gradient update in Equation (33b), and using condition (g) at iteration  $t$ , we have:

$$\alpha_{\oplus,j}^{(t+1)} = \alpha_{\oplus,j}^{(t)} \quad \text{for all } j \in [n] \text{ with } y_j < 0.$$

Therefore, condition (c) continues to hold at iteration  $t + 1$ .

Part (d): According to the dual gradient update in Equation (34b), and using conditions (h) at iteration  $t$ , we have:

$$\alpha_{\ominus,i}^{(t+1)} = \alpha_{\ominus,i}^{(t)} \quad \text{for all } i \in [n] \text{ with } y_i > 0.$$

Therefore, condition (d) continues to hold at iteration  $t + 1$ .

Part (e): By conditions (a), (b), (g), and (h), the gradient update at iteration  $t$  for  $\beta_{\oplus}^{(t)}$  depends only on the positively labeled examples, and the update for  $\beta_{\ominus}^{(t)}$  depends only on the negatively labeled examples. Hence, the gradient update for an individual neuron is equivalent to gradient descent on a certain linear regression problem. Since the step size satisfies  $\eta \leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1}$ , the linear regression squared loss is monotonically nonincreasing (as in the proof of Lemma 2), and by condition (e) at iteration  $t$ , we obtain

$$\begin{aligned} \left\| \boldsymbol{\beta}_{\oplus,+}^{(t+1)} - \mathbf{y}_+ \right\|_2 &\leq \left\| \boldsymbol{\beta}_{\oplus,+}^{(t)} - \mathbf{y}_+ \right\|_2 \leq C_y \|\mathbf{y}_+\|_2, \\ \left\| -\boldsymbol{\beta}_{\ominus,-}^{(t+1)} - \mathbf{y}_- \right\|_2 &\leq \left\| -\boldsymbol{\beta}_{\ominus,-}^{(t)} - \mathbf{y}_- \right\|_2 \leq C_y \|\mathbf{y}_-\|_2, \end{aligned}$$

where we use  $\mathbf{y}_+$  ( $\mathbf{y}_-$ ) to denote the vector of positively labeled (negatively labeled) examples. Therefore, condition (e) holds at iteration  $t + 1$ .

Part (f): Following the same argument as in Part (d) of Lemma 17, using conditions (c), (d), and (e) at iteration  $t + 1$ , together with the eigenvalue bounds from Lemma 12, we have

$$\left\| \boldsymbol{\alpha}_{\oplus}^{(t+1)} \right\|_2 \leq \frac{C_\alpha \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1}, \quad \left\| \boldsymbol{\alpha}_{\ominus}^{(t+1)} \right\|_2 \leq \frac{C_\alpha \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1}$$

with probability at least  $1 - 2e^{-n/C_g}$ . Thus, condition (f) holds at iteration  $t + 1$ .

Part (g): By Lemma 11, since conditions (c) and (f) hold at iteration  $t + 1$ , we conclude that  $\beta_{\oplus, j}^{(t+1)} \leq 0$  for all  $j \in [n]$  with  $y_j < 0$ . Thus, condition (g) holds at iteration  $t + 1$ .

Part (h): Similarly, since conditions (d) and (f) hold at iteration  $t + 1$ , we have  $\beta_{\ominus, i}^{(t+1)} \leq 0$  for all  $i \in [n]$  with  $y_i > 0$ . Thus, condition (h) holds at iteration  $t + 1$ . ■

Equipped with Lemma 19, we are ready to prove Theorem 8.

**Proof** (Theorem 8) The proof follows a similar structure to that of Theorem 4 for single ReLU models, but now we must track the dynamics for both  $\boldsymbol{w}_{\oplus}$  and  $\boldsymbol{w}_{\ominus}$  simultaneously. Equipped with sufficient conditions under which the activation patterns are preserved in Lemma 19, we verify these conditions hold after the first gradient step, and use induction to characterize the full gradient descent dynamics.

We first verify that the iterate at  $t = 1$  satisfies all the sufficient conditions. With the initialization

$$\boldsymbol{w}_{\oplus}^{(0)} = \boldsymbol{X}^\top \left( \boldsymbol{X} \boldsymbol{X}^\top \right)^{-1} \boldsymbol{\epsilon}_{\oplus}, \quad \boldsymbol{w}_{\ominus}^{(0)} = \boldsymbol{X}^\top \left( \boldsymbol{X} \boldsymbol{X}^\top \right)^{-1} \boldsymbol{\epsilon}_{\ominus},$$

we have  $\beta_{\oplus}^{(0)} = \boldsymbol{\epsilon}_{\oplus}$  and  $\beta_{\ominus}^{(0)} = \boldsymbol{\epsilon}_{\ominus}$ . By the theorem assumptions,  $0 < \epsilon_{\oplus, i} \leq \frac{1}{2C_\alpha} y_{\min}$  and  $0 < \epsilon_{\ominus, i} \leq \frac{1}{2C_\alpha} y_{\min}$  for all  $i \in [n]$ . Using the primal gradient updates in Equations (33a) and (34a), we have

$$\begin{aligned} \beta_{\oplus}^{(1)} &= \boldsymbol{\epsilon}_{\oplus} - \eta \boldsymbol{X} \boldsymbol{X}^\top \boldsymbol{D}(\boldsymbol{\epsilon}_{\oplus}) (h_{\boldsymbol{\Theta}^{(0)}}(\boldsymbol{X}) - \boldsymbol{y}) \\ &= \boldsymbol{\epsilon}_{\oplus} - \eta \boldsymbol{X} \boldsymbol{X}^\top \boldsymbol{D}(\boldsymbol{\epsilon}_{\oplus}) (\sigma(\boldsymbol{\epsilon}_{\oplus}) - \sigma(\boldsymbol{\epsilon}_{\ominus}) - \boldsymbol{y}) \\ &= \boldsymbol{\epsilon}_{\oplus} - \eta \boldsymbol{X} \boldsymbol{X}^\top (\boldsymbol{\epsilon}_{\oplus} - \boldsymbol{\epsilon}_{\ominus} - \boldsymbol{y}), \end{aligned} \quad (35)$$

where the last equality uses the fact that  $\boldsymbol{\epsilon}_{\oplus} > \mathbf{0}$  and  $\boldsymbol{\epsilon}_{\ominus} > \mathbf{0}$  componentwise, so  $\boldsymbol{D}(\boldsymbol{\epsilon}_{\oplus}) = \boldsymbol{I}$ ,  $\sigma(\boldsymbol{\epsilon}_{\oplus}) = \boldsymbol{\epsilon}_{\oplus}$ , and  $\sigma(\boldsymbol{\epsilon}_{\ominus}) = \boldsymbol{\epsilon}_{\ominus}$ . Similarly, we have

$$\begin{aligned} \beta_{\ominus}^{(1)} &= \boldsymbol{\epsilon}_{\ominus} + \eta \boldsymbol{X} \boldsymbol{X}^\top \boldsymbol{D}(\boldsymbol{\epsilon}_{\ominus}) (\sigma(\boldsymbol{\epsilon}_{\oplus}) - \sigma(\boldsymbol{\epsilon}_{\ominus}) - \boldsymbol{y}) \\ &= \boldsymbol{\epsilon}_{\ominus} + \eta \boldsymbol{X} \boldsymbol{X}^\top (\boldsymbol{\epsilon}_{\oplus} - \boldsymbol{\epsilon}_{\ominus} - \boldsymbol{y}). \end{aligned} \quad (36)$$

For the dual variables, we have  $\boldsymbol{\alpha}_{\oplus}^{(0)} = (\boldsymbol{X} \boldsymbol{X}^\top)^{-1} \boldsymbol{\epsilon}_{\oplus}$  and  $\boldsymbol{\alpha}_{\ominus}^{(0)} = (\boldsymbol{X} \boldsymbol{X}^\top)^{-1} \boldsymbol{\epsilon}_{\ominus}$ . The dual updates give:

$$\begin{aligned} \boldsymbol{\alpha}_{\oplus}^{(1)} &= \boldsymbol{\alpha}_{\oplus}^{(0)} - \eta \boldsymbol{D}(\boldsymbol{\epsilon}_{\oplus}) (\boldsymbol{\epsilon}_{\oplus} - \boldsymbol{\epsilon}_{\ominus} - \boldsymbol{y}) \\ &= \left( \boldsymbol{X} \boldsymbol{X}^\top \right)^{-1} \boldsymbol{\epsilon}_{\oplus} - \eta (\boldsymbol{\epsilon}_{\oplus} - \boldsymbol{\epsilon}_{\ominus} - \boldsymbol{y}) \\ &= \eta \left( \boldsymbol{y} - \boldsymbol{\epsilon}_{\oplus} + \boldsymbol{\epsilon}_{\ominus} + \frac{1}{\eta} \left( \boldsymbol{X} \boldsymbol{X}^\top \right)^{-1} \boldsymbol{\epsilon}_{\oplus} \right), \end{aligned}$$

and

$$\begin{aligned}
 \alpha_{\ominus}^{(1)} &= \alpha_{\ominus}^{(0)} + \eta \mathbf{D}(\epsilon_{\ominus})(\epsilon_{\oplus} - \epsilon_{\ominus} - \mathbf{y}) \\
 &= \left( \mathbf{X} \mathbf{X}^{\top} \right)^{-1} \epsilon_{\ominus} + \eta (\epsilon_{\oplus} - \epsilon_{\ominus} - \mathbf{y}) \\
 &= \eta \left( -\mathbf{y} + \epsilon_{\oplus} - \epsilon_{\ominus} + \frac{1}{\eta} \left( \mathbf{X} \mathbf{X}^{\top} \right)^{-1} \epsilon_{\ominus} \right).
 \end{aligned}$$

We now verify each condition at  $t = 1$ .

Part (a): For all  $i \in [n]$  with  $y_i > 0$ , we apply Lemma 10. Since  $\beta_{\oplus,i}^{(0)} = \epsilon_{\oplus,i} > 0$ ,  $h_{\Theta^{(0)}}(\mathbf{x}_i) = \beta_{\oplus,i}^{(0)} - \beta_{\ominus,i}^{(0)} \leq \beta_{\oplus,i}^{(0)}$ , and

$$\begin{aligned}
 \|h_{\Theta^{(0)}}(\mathbf{X}) - \mathbf{y}\|_2 &= \|\epsilon_{\oplus} - \epsilon_{\ominus} - \mathbf{y}\|_2 \\
 &\leq \|\epsilon_{\oplus}\|_2 + \|\epsilon_{\ominus}\|_2 + \|\mathbf{y}\|_2 \leq \frac{\sqrt{n}}{C_{\alpha}} y_{\min} + \|\mathbf{y}\|_2 \leq C_y \|\mathbf{y}\|_2, \quad (37)
 \end{aligned}$$

where we have used  $C_y \geq \frac{1}{C_{\alpha}} + 1$ . We conclude that  $\beta_{\oplus,i}^{(1)} > 0$  for all  $i$  with  $y_i > 0$ .

Part (b): For all  $j \in [n]$  with  $y_j < 0$ , we apply Lemma 10. Since  $\beta_{\ominus,j}^{(0)} = \epsilon_{\ominus,j} > 0$ ,  $-h_{\Theta^{(0)}}(\mathbf{x}_j) = -\beta_{\oplus,j}^{(0)} + \beta_{\ominus,j}^{(0)} \leq \beta_{\ominus,j}^{(0)}$ , and  $\|h_{\Theta^{(0)}}(\mathbf{X}) - \mathbf{y}\|_2 \leq C_y \|\mathbf{y}\|_2$  by Equation (37), we conclude that  $\beta_{\ominus,j}^{(1)} > 0$  for all  $j \in [n]$  with  $y_j < 0$ .

Part (c): For all  $j \in [n]$  with  $y_j < 0$ , we verify that  $\alpha_{\oplus,j}^{(1)}$  satisfies the required upper and lower bounds. For the upper bound, we have

$$\begin{aligned}
 \alpha_{\oplus,j}^{(1)} &= \eta \left( y_j - \epsilon_{\oplus,j} + \epsilon_{\ominus,j} + \frac{1}{\eta} \mathbf{e}_j^{\top} \left( \mathbf{X} \mathbf{X}^{\top} \right)^{-1} \epsilon_{\oplus} \right) \\
 &= \eta \left( y_j - \epsilon_{\oplus,j} + \epsilon_{\ominus,j} + \frac{1}{\eta} \mathbf{e}_j^{\top} \left[ \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} + \left( \left( \mathbf{X} \mathbf{X}^{\top} \right)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right) \right] \epsilon_{\oplus} \right) \\
 &= \eta \left( y_j - \epsilon_{\oplus,j} + \epsilon_{\ominus,j} + \frac{\epsilon_{\oplus,j}}{\eta \|\boldsymbol{\lambda}\|_1} + \frac{1}{\eta} \mathbf{e}_j^{\top} \left( \left( \mathbf{X} \mathbf{X}^{\top} \right)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right) \epsilon_{\oplus} \right) \\
 &\leq \eta \left( y_j + \epsilon_{\ominus,j} + \frac{\epsilon_{\oplus,j}}{\eta \|\boldsymbol{\lambda}\|_1} + \frac{1}{\eta} \left\| \left( \mathbf{X} \mathbf{X}^{\top} \right)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right\|_2 \|\epsilon_{\oplus}\|_2 \right),
 \end{aligned}$$

where the inequality drops the negative term  $-\epsilon_{\oplus,j}$ . By Corollary 14, we have

$$\left\| \left( \mathbf{X} \mathbf{X}^{\top} \right)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right\|_2 \leq \frac{C_g C}{\|\boldsymbol{\lambda}\|_1} \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_{\infty}} \right),$$

with probability at least  $1 - 2 \exp(-n(Cc - \ln 9))$ . Moreover, by the theorem assumptions,  $\epsilon_{\oplus,j} \leq \frac{1}{2C_{\alpha}} y_{\min}$ ,  $\epsilon_{\ominus,j} \leq \frac{1}{2C_{\alpha}} y_{\min}$ , and  $\frac{1}{\eta} \leq CC_g \|\boldsymbol{\lambda}\|_1$ . Combining these bounds

yields

$$\begin{aligned}
 \alpha_{\oplus,j}^{(1)} &\leq \frac{1}{CC_g \|\boldsymbol{\lambda}\|_1} \left( -y_{\min} + \frac{1}{2C_\alpha} y_{\min} + \frac{CC_g}{2C_\alpha} y_{\min} \right. \\
 &\quad \left. + C^2 C_g^2 \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \cdot \frac{\sqrt{n}}{2C_\alpha} y_{\min} \right) \\
 &\leq \frac{1}{CC_g \|\boldsymbol{\lambda}\|_1} \left( -y_{\min} + \frac{1}{2C_\alpha} y_{\min} + \frac{CC_g}{2C_\alpha} y_{\min} + C^2 C_g^2 \cdot \frac{y_{\min}}{C_0 y_{\max}} \cdot \frac{1}{2C_\alpha} y_{\min} \right) \\
 &= -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1} \left( \frac{C_\alpha}{CC_g} - \frac{1}{2CC_g} - \frac{1}{2} - \frac{CC_g y_{\min}}{2C_0 y_{\max}} \right) \\
 &\leq -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1}. \tag{38}
 \end{aligned}$$

The second inequality follows from  $d_2 \geq C_0^2 \frac{n^2 y_{\max}^2}{y_{\min}^2}$  and  $d_\infty \geq C_0 \frac{n^{1.5} y_{\max}}{y_{\min}}$  in Assumption 2, and the last inequality follows the relationship between constants  $C_0 \gtrsim C_\alpha^2$  and  $C_\alpha \gtrsim \max\{C_g^2, C_y C_g\}$ . For the lower bound, we have

$$\begin{aligned}
 \alpha_{\oplus,j}^{(1)} &= \eta \left( y_j - \epsilon_{\oplus,j} + \epsilon_{\ominus,j} + \frac{\epsilon_{\oplus,j}}{\eta \|\boldsymbol{\lambda}\|_1} + \frac{1}{\eta} \mathbf{e}_j^\top \left( (\mathbf{X}\mathbf{X}^\top)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right) \boldsymbol{\epsilon}_\oplus \right) \\
 &\geq \eta \left( -y_{\max} - \epsilon_{\oplus,j} - \frac{1}{\eta} \left\| \left( \mathbf{X}\mathbf{X}^\top \right)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right\|_2 \|\boldsymbol{\epsilon}_\oplus\|_2 \right) \\
 &\geq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1} \left( -y_{\max} - \frac{1}{2C_\alpha} y_{\min} - C^2 C_g^2 \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \cdot \frac{\sqrt{n}}{2C_\alpha} y_{\min} \right) \\
 &\geq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1} \left( -y_{\max} - \frac{1}{2C_\alpha} y_{\min} - C^2 C_g^2 \cdot \frac{y_{\min}}{C_0 y_{\max}} \cdot \frac{1}{2C_\alpha} y_{\min} \right) \\
 &\geq -\frac{3y_{\max}}{C_g \|\boldsymbol{\lambda}\|_1}, \tag{39}
 \end{aligned}$$

by the same arguments. Thus,  $\alpha_{\oplus,j}^{(1)}$  satisfies both the required upper and lower bounds for all  $j$  with  $y_j < 0$ .

Part (d): For all  $i \in [n]$  with  $y_i > 0$ , we verify that  $\alpha_{\ominus,i}^{(1)}$  satisfies the required bounds in the approach analogous to Part (c). For the upper bound, we have

$$\begin{aligned}
 \alpha_{\ominus,i}^{(1)} &= \eta \left( -y_i + \epsilon_{\oplus,i} - \epsilon_{\ominus,i} + \frac{1}{\eta} \mathbf{e}_i^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \boldsymbol{\epsilon}_\ominus \right) \\
 &= \eta \left( -y_i + \epsilon_{\oplus,i} - \epsilon_{\ominus,i} + \frac{1}{\eta} \mathbf{e}_i^\top \left[ \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} + \left( (\mathbf{X}\mathbf{X}^\top)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right) \right] \boldsymbol{\epsilon}_\ominus \right) \\
 &= \eta \left( -y_i + \epsilon_{\oplus,i} - \epsilon_{\ominus,i} + \frac{\epsilon_{\ominus,i}}{\eta \|\boldsymbol{\lambda}\|_1} + \frac{1}{\eta} \mathbf{e}_i^\top \left( (\mathbf{X}\mathbf{X}^\top)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right) \boldsymbol{\epsilon}_\ominus \right) \\
 &\leq \eta \left( -y_i + \epsilon_{\oplus,i} + \frac{\epsilon_{\ominus,i}}{\eta \|\boldsymbol{\lambda}\|_1} + \frac{1}{\eta} \left\| \left( \mathbf{X}\mathbf{X}^\top \right)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right\|_2 \|\boldsymbol{\epsilon}_\ominus\|_2 \right),
 \end{aligned}$$

where the inequality drops the negative term  $-\epsilon_{\ominus,i}$ . Applying the upper bound in Corollary 14 and the theorem assumptions  $\epsilon_{\oplus,i} \leq \frac{1}{2C_\alpha} y_{\min}$ ,  $\epsilon_{\ominus,i} \leq \frac{1}{2C_\alpha} y_{\min}$ , and  $\frac{1}{\eta} \leq CC_g \|\boldsymbol{\lambda}\|_1$ , we have

$$\begin{aligned} \alpha_{\ominus,i}^{(1)} &\leq \frac{1}{CC_g \|\boldsymbol{\lambda}\|_1} \left( -y_{\min} + \frac{1}{2C_\alpha} y_{\min} + \frac{CC_g}{2C_\alpha} y_{\min} \right. \\ &\quad \left. + C^2 C_g^2 \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \cdot \frac{\sqrt{n}}{2C_\alpha} y_{\min} \right) \\ &\leq -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1}, \end{aligned}$$

where the second inequality follows the argument used in Equation (38). For the lower bound, following the same argument as in Part (c), we have

$$\begin{aligned} \alpha_{\ominus,i}^{(1)} &= \eta \left( -y_i + \epsilon_{\oplus,i} - \epsilon_{\ominus,i} + \frac{\epsilon_{\ominus,i}}{\eta \|\boldsymbol{\lambda}\|_1} + \frac{1}{\eta} \mathbf{e}_i^\top \left( (\mathbf{X}\mathbf{X}^\top)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right) \boldsymbol{\epsilon}_\ominus \right) \\ &\geq \eta \left( -y_{\max} - \epsilon_{\ominus,i} - \frac{1}{\eta} \left\| \left( \mathbf{X}\mathbf{X}^\top \right)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right\|_2 \|\boldsymbol{\epsilon}_\ominus\|_2 \right) \\ &\geq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1} \left( -y_{\max} - \frac{1}{2C_\alpha} y_{\min} - C^2 C_g^2 \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \cdot \frac{\sqrt{n}}{2C_\alpha} y_{\min} \right) \\ &\geq -\frac{3y_{\max}}{C_g \|\boldsymbol{\lambda}\|_1}, \end{aligned}$$

where the last inequality follows follows the argument used in Equation (39). Thus,  $\alpha_{\ominus,i}^{(1)}$  satisfies both bounds for all  $i$  with  $y_i > 0$ .

Part (e): We verify that the primal variables  $\boldsymbol{\beta}_\oplus^{(1)}$  corresponding to positively labeled examples minus  $\mathbf{y}_+$  satisfy the norm bound. Specifically, we show that  $\left\| \boldsymbol{\beta}_{\oplus,+}^{(1)} - \mathbf{y}_+ \right\|_2^2 \leq C_y^2 \|\mathbf{y}_+\|_2^2$ . According to Equation (35), we have

$$\begin{aligned} \left\| \boldsymbol{\beta}_{\oplus,+}^{(1)} - \mathbf{y}_+ \right\|_2^2 &= \sum_{i:y_i>0} \left( \beta_{\oplus,i}^{(1)} - y_i \right)^2 \\ &= \sum_{i:y_i>0} \left( \underbrace{\epsilon_{\oplus,i} - \eta \mathbf{e}_i^\top \mathbf{X}\mathbf{X}^\top (\boldsymbol{\epsilon}_\oplus - \boldsymbol{\epsilon}_\ominus - \mathbf{y}) - y_i}_{=:T_i} \right)^2. \end{aligned} \quad (40)$$

Next, we bound the term  $T_i := \epsilon_{\oplus,i} - \eta \mathbf{e}_i^\top \mathbf{X}\mathbf{X}^\top (\boldsymbol{\epsilon}_\oplus - \boldsymbol{\epsilon}_\ominus - \mathbf{y}) - y_i$  for all  $i \in [n]$  with  $y_i > 0$ . We have

$$\begin{aligned} T_i &= \epsilon_{\oplus,i} - \eta \mathbf{e}_i^\top \mathbf{X}\mathbf{X}^\top (\boldsymbol{\epsilon}_\oplus - \boldsymbol{\epsilon}_\ominus - \mathbf{y}) - y_i \\ &= (\epsilon_{\oplus,i} - y_i) - \eta \mathbf{e}_i^\top \left[ \|\boldsymbol{\lambda}\|_1 \mathbf{I} + \left( \mathbf{X}\mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right) \right] (\boldsymbol{\epsilon}_\oplus - \boldsymbol{\epsilon}_\ominus - \mathbf{y}) \\ &= (1 - \eta \|\boldsymbol{\lambda}\|_1) \epsilon_{\oplus,i} + \eta \|\boldsymbol{\lambda}\|_1 \epsilon_{\ominus,i} - (1 - \eta \|\boldsymbol{\lambda}\|_1) y_i \\ &\quad - \eta \mathbf{e}_i^\top \left( \mathbf{X}\mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right) (\boldsymbol{\epsilon}_\oplus - \boldsymbol{\epsilon}_\ominus - \mathbf{y}). \end{aligned}$$

Since the step size assumption guarantees that  $\frac{1}{CC_g\|\boldsymbol{\lambda}\|_1} \leq \eta \leq \frac{1}{C_g\|\boldsymbol{\lambda}\|_1}$ , and  $\epsilon_{\oplus,i}, \epsilon_{\ominus,i} \leq \frac{1}{2C_\alpha} y_{\min}$ , we have

$$\begin{aligned} (1 - \eta \|\boldsymbol{\lambda}\|_1) \epsilon_{\oplus,i} + \eta \|\boldsymbol{\lambda}\|_1 \epsilon_{\ominus,i} - (1 - \eta \|\boldsymbol{\lambda}\|_1) y_i & \\ & \leq \epsilon_{\oplus,i} + \eta \|\boldsymbol{\lambda}\|_1 \epsilon_{\ominus,i} - (1 - \eta \|\boldsymbol{\lambda}\|_1) y_i \\ & \leq \left(1 + \frac{1}{C_g}\right) \frac{1}{2C_\alpha} y_{\min} - \left(1 - \frac{1}{C_g}\right) y_{\min} \\ & < 0, \end{aligned}$$

with  $C_\alpha \gtrsim C_g^2$ . Hence, in order to upper bound  $T_i^2$ , it suffices to find the lower bound for  $T_i$ . We have

$$\begin{aligned} T_i &= (1 - \eta \|\boldsymbol{\lambda}\|_1) \epsilon_{\oplus,i} + \eta \|\boldsymbol{\lambda}\|_1 \epsilon_{\ominus,i} - (1 - \eta \|\boldsymbol{\lambda}\|_1) y_i \\ & \quad - \eta \mathbf{e}_i^\top (\mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}) (\boldsymbol{\epsilon}_\oplus - \boldsymbol{\epsilon}_\ominus - \mathbf{y}) \\ & \geq -y_i - \eta \left\| \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right\|_2 \|\boldsymbol{\epsilon}_\oplus - \boldsymbol{\epsilon}_\ominus - \mathbf{y}\|_2, \end{aligned}$$

where the inequality drops the positive terms  $(1 - \eta \|\boldsymbol{\lambda}\|_1) \epsilon_{\oplus,i}$ ,  $\eta \|\boldsymbol{\lambda}\|_1 \epsilon_{\ominus,i}$ , and  $\eta \|\boldsymbol{\lambda}\|_1 y_i$ . We again upper bound  $\left\| \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right\|_2$  by Corollary 14. With probability at least  $1 - 2 \exp(-n(Cc - \ln 9))$ , we have

$$\begin{aligned} T_i &\geq -y_i - \eta \cdot C \|\boldsymbol{\lambda}\|_1 \cdot \max\left(\sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty}\right) \|\boldsymbol{\epsilon}_\oplus - \boldsymbol{\epsilon}_\ominus - \mathbf{y}\|_2 \\ &\geq -y_i - \frac{C}{C_g} \cdot \max\left(\sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty}\right) \|\boldsymbol{\epsilon}_\oplus - \boldsymbol{\epsilon}_\ominus - \mathbf{y}\|_2, \end{aligned}$$

by applying  $\eta \leq \frac{1}{C_g\|\boldsymbol{\lambda}\|_1}$ . Finally, we apply the upper bounds for  $\|\boldsymbol{\epsilon}_\oplus\|_2$ ,  $\|\boldsymbol{\epsilon}_\ominus\|_2$  and  $\|\mathbf{y}\|_2$ , and Assumption 2 ensures that  $d_2 \geq C_0^2 \frac{n^2 y_{\max}^2}{y_{\min}^2}$  and  $d_\infty \geq C_0 \frac{n^{1.5} y_{\max}}{y_{\min}}$ . We have

$$\begin{aligned} T_i &\geq -y_i - \frac{C}{C_g} \max\left(\sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty}\right) \left(\frac{\sqrt{n}}{C_\alpha} y_{\min} + \sqrt{n} y_{\max}\right) \\ &\geq -y_i - \frac{C y_{\min}}{C_g C_0 y_{\max}} \left(\frac{1}{C_\alpha} y_{\min} + y_{\max}\right) \\ &\geq -y_i \left(1 + \frac{2C}{C_g C_0}\right) \\ &\geq -C_y y_i, \end{aligned}$$

with the choice of  $C_y \geq 2$ . Substituting  $T_i^2 \leq C_y^2 y_i^2$  into Equation (40), we have

$$\left\| \boldsymbol{\beta}_{\oplus,+}^{(1)} - \mathbf{y}_+ \right\|_2^2 \leq \sum_{i:y_i>0} C_y^2 y_i^2 = C_y^2 \|\mathbf{y}_+\|_2^2.$$

As a result, we conclude that  $\left\| \boldsymbol{\beta}_{\oplus,+}^{(1)} - \mathbf{y}_+ \right\|_2 \leq C_y \|\mathbf{y}_+\|_2$  as required. The same derivation holds for  $\left\| \boldsymbol{\beta}_{\ominus,-}^{(1)} + \mathbf{y}_- \right\|_2 \leq C_y \|\mathbf{y}_-\|_2$  by an analogous argument. Therefore, condition (e) holds at  $t = 1$ .

Part (f): We verify the norm bounds on the dual variables. By the triangle inequality, we have

$$\begin{aligned} \left\| \boldsymbol{\alpha}_{\oplus}^{(1)} \right\|_2 &= \left\| \eta \left( \mathbf{y} - \boldsymbol{\epsilon}_{\oplus} + \boldsymbol{\epsilon}_{\ominus} + \frac{1}{\eta} \left( \mathbf{X} \mathbf{X}^{\top} \right)^{-1} \boldsymbol{\epsilon}_{\oplus} \right) \right\|_2 \\ &\leq \eta \left( \|\mathbf{y}\|_2 + \|\boldsymbol{\epsilon}_{\oplus}\|_2 + \|\boldsymbol{\epsilon}_{\ominus}\|_2 + \frac{1}{\eta} \left\| \left( \mathbf{X} \mathbf{X}^{\top} \right)^{-1} \right\|_2 \|\boldsymbol{\epsilon}_{\oplus}\|_2 \right). \end{aligned}$$

Using  $\|\mathbf{y}\|_2 \leq \sqrt{n}y_{\max}$ ,  $\|\boldsymbol{\epsilon}_{\oplus}\|_2, \|\boldsymbol{\epsilon}_{\ominus}\|_2 \leq \frac{\sqrt{n}}{2C_{\alpha}}y_{\min}$ ,  $\left\| \left( \mathbf{X} \mathbf{X}^{\top} \right)^{-1} \right\|_2 \leq \frac{C_g}{\|\boldsymbol{\lambda}\|_1}$ , and  $\frac{1}{CC_g\|\boldsymbol{\lambda}\|_1} \leq \eta \leq \frac{1}{C_g\|\boldsymbol{\lambda}\|_1}$ , we have

$$\begin{aligned} \left\| \boldsymbol{\alpha}_{\oplus}^{(1)} \right\|_2 &\leq \frac{1}{C_g\|\boldsymbol{\lambda}\|_1} \left( \sqrt{n}y_{\max} + \frac{\sqrt{n}}{C_{\alpha}}y_{\min} + CC_g\|\boldsymbol{\lambda}\|_1 \cdot \frac{C_g}{\|\boldsymbol{\lambda}\|_1} \cdot \frac{\sqrt{n}}{C_{\alpha}}y_{\min} \right) \\ &\leq \frac{1}{C_g\|\boldsymbol{\lambda}\|_1} (3\sqrt{n}y_{\max}) \\ &\leq \frac{C_{\alpha}\sqrt{n}y_{\max}}{\|\boldsymbol{\lambda}\|_1}, \end{aligned}$$

with  $C_{\alpha} \gtrsim \max\{C_g^2, C_y C_g\}$ . The same bound holds for  $\left\| \boldsymbol{\alpha}_{\ominus}^{(1)} \right\|_2$ . Thus, condition (f) holds at  $t = 1$ .

Part (g): Since we have shown that  $\alpha_{\oplus,j}^{(1)} \leq -\frac{y_{\min}}{C_{\alpha}\|\boldsymbol{\lambda}\|_1}$  and  $\left\| \boldsymbol{\alpha}_{\oplus}^{(1)} \right\|_2 \leq \frac{C_{\alpha}\sqrt{n}y_{\max}}{\|\boldsymbol{\lambda}\|_1}$  for all  $j \in [n]$  with  $y_j < 0$ , it follows from Lemma 11 that  $\beta_{\oplus,j}^{(1)} \leq 0$  for all  $j \in [n]$  with  $y_j < 0$ .

Part (h): Similarly, since we have shown that  $\alpha_{\ominus,i}^{(1)} \leq -\frac{y_{\min}}{C_{\alpha}\|\boldsymbol{\lambda}\|_1}$  and  $\left\| \boldsymbol{\alpha}_{\ominus}^{(1)} \right\|_2 \leq \frac{C_{\alpha}\sqrt{n}y_{\max}}{\|\boldsymbol{\lambda}\|_1}$  for all  $i \in [n]$  with  $y_i > 0$ , it follows from Lemma 11 that  $\beta_{\ominus,i}^{(1)} \leq 0$  for all  $i \in [n]$  with  $y_i > 0$ .

We have shown that at iteration  $t = 1$  the conditions in Lemma 19 are satisfied, and by induction, these conditions will also hold for  $t \geq 1$ . As a result, the positive neuron  $\mathbf{w}_{\oplus}$  is trained with only positive examples starting from the iteration  $t = 1$ , and it is equivalent to linear regression using only positive examples with initialization  $\mathbf{w}_{\oplus}^{(1)} = \eta \mathbf{X}^{\top} \left( \mathbf{y} - \boldsymbol{\epsilon}_{\oplus} + \boldsymbol{\epsilon}_{\ominus} + \frac{1}{\eta} \left( \mathbf{X} \mathbf{X}^{\top} \right)^{-1} \boldsymbol{\epsilon}_{\oplus} \right)$ . Finally, since  $\mathbf{w}_{\oplus}$  and  $\mathbf{w}_{\ominus}$  are trained on disjoint subsets of examples,  $\mathbf{w}_{\oplus}^{(\infty)}$  satisfies

$$\mathbf{w}_{\oplus}^{(\infty)} = \arg \min_{\mathbf{w} \in \{\mathbf{w} : \mathbf{X}_+ \mathbf{w} = \mathbf{y}_+\}} \left\| \mathbf{w} - \mathbf{w}_{\oplus}^{(1)} \right\|_2,$$

by Lemma 2. The same arguments apply to the negative neuron  $\mathbf{w}_{\ominus}$  as well. This completes the proof of Theorem 8. ■

### C.3. Proof of Theorem 9 (Implicit Bias Approximation to $\mathbf{w}^*$ )

**Proof** (Theorem 9) We divide the proof into four steps, and formally show the result for  $\mathbf{w}_{\oplus}^*$ . The result for  $\mathbf{w}_{\ominus}^*$  follows an identical series of steps. In the first step, we derive an upper bound for  $\left\| \boldsymbol{\alpha}_{\oplus}^* \right\|_2$  and  $\left\| \boldsymbol{\alpha}_{\ominus}^* \right\|_2$  where  $\mathbf{w}_{\oplus}^* := \mathbf{X}^{\top} \boldsymbol{\alpha}_{\oplus}^*$  and  $\mathbf{w}_{\ominus}^* := \mathbf{X}^{\top} \boldsymbol{\alpha}_{\ominus}^*$ , by using the optimality of the objective function in (8). In the second step, we use the KKT conditions of (9) to find a representation of  $\{\mathbf{w}_{\oplus}^*, \mathbf{w}_{\ominus}^*\}$ . In Steps 3 and 4, we derive the corresponding upper bound and lower bound.

**Step 1: Upper bounds for  $\|\alpha_\oplus^*\|_2$  and  $\|\alpha_\ominus^*\|_2$ .**

$\{\mathbf{w}_\oplus^*, \mathbf{w}_\ominus^*\}$  is the optimal solution to (8) and it achieves the minimum objective in (8). In the proof of Lemma 7, we introduce  $\{\tilde{\mathbf{w}}_\oplus, \tilde{\mathbf{w}}_\ominus\}$  which is also a feasible solution to (8), where  $\tilde{\mathbf{w}}_\oplus := \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{y}_\oplus$  and  $\tilde{\mathbf{w}}_\ominus := \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{y}_\ominus$ , with  $y_{\oplus,i} = \max\{y_i, 0\}$  and  $y_{\ominus,i} = -\min\{y_i, 0\}$ . Therefore, by the optimality of  $\{\mathbf{w}_\oplus^*, \mathbf{w}_\ominus^*\}$  in the objective, we have

$$\begin{aligned} \|\mathbf{w}_\oplus^*\|_2^2 + \|\mathbf{w}_\ominus^*\|_2^2 &= \alpha_\oplus^{*\top} \mathbf{X} \mathbf{X}^\top \alpha_\oplus^* + \alpha_\ominus^{*\top} \mathbf{X} \mathbf{X}^\top \alpha_\ominus^* \\ &\leq \|\tilde{\mathbf{w}}_\oplus\|_2^2 + \|\tilde{\mathbf{w}}_\ominus\|_2^2 \\ &= \mathbf{y}_\oplus^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{y}_\oplus + \mathbf{y}_\ominus^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{y}_\ominus \\ &\leq \left\| (\mathbf{X} \mathbf{X}^\top)^{-1} \right\|_2 \|\mathbf{y}_\oplus\|_2^2 + \left\| (\mathbf{X} \mathbf{X}^\top)^{-1} \right\|_2 \|\mathbf{y}_\ominus\|_2^2 \\ &\leq \frac{2C_g n y_{\max}^2}{\|\boldsymbol{\lambda}\|_1}, \end{aligned}$$

where the last inequality uses Lemma 12 with probability at least  $1 - 2e^{-n/C_g}$ , and we have  $\max\{\|\mathbf{y}_\oplus\|_2^2, \|\mathbf{y}_\ominus\|_2^2\} \leq \|\mathbf{y}\|_2^2 \leq n y_{\max}^2$ . Therefore, we have

$$\begin{aligned} \lambda_n(\mathbf{X} \mathbf{X}^\top) \|\alpha_\oplus^*\|_2^2 &\leq \alpha_\oplus^{*\top} \mathbf{X} \mathbf{X}^\top \alpha_\oplus^* \\ &\leq \alpha_\oplus^{*\top} \mathbf{X} \mathbf{X}^\top \alpha_\oplus^* + \alpha_\ominus^{*\top} \mathbf{X} \mathbf{X}^\top \alpha_\ominus^* \\ &\leq \frac{2C_g n y_{\max}^2}{\|\boldsymbol{\lambda}\|_1}. \end{aligned}$$

As a result, we have  $\|\alpha_\oplus^*\|_2 \leq \frac{\sqrt{2}C_g \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1}$ . This bound applies to  $\|\alpha_\ominus^*\|_2$  as well via an identical argument.

**Step 2: KKT conditions of  $\{\mathbf{w}_\oplus^*, \mathbf{w}_\ominus^*\}$  by Lemma 7.**

Based on Lemma 7, the optimal solution  $\{\mathbf{w}_\oplus^*, \mathbf{w}_\ominus^*\}$  of (8) is also the optimal solution of a convex program (9). Hence, we restate the convex program in (9) below

$$\begin{aligned} \mathbf{w}_\oplus^*, \mathbf{w}_\ominus^* &= \arg \min_{\mathbf{w}_\oplus, \mathbf{w}_\ominus} \frac{1}{2} \|\mathbf{w}_\oplus\|_2^2 + \frac{1}{2} \|\mathbf{w}_\ominus\|_2^2 \\ \text{s.t. } \mathbf{w}_\oplus^\top \mathbf{x}_i &= y_i, & \mathbf{w}_\ominus^\top \mathbf{x}_i &\leq 0, & \text{for all } i \in S_1, \\ \mathbf{w}_\oplus^\top \mathbf{x}_i - \mathbf{w}_\ominus^\top \mathbf{x}_i &= y_i, & -\mathbf{w}_\ominus^\top \mathbf{x}_i &\leq 0, & \text{for all } i \in S_2, \\ & & -\mathbf{w}_\ominus^\top \mathbf{x}_i &= y_i, & \mathbf{w}_\oplus^\top \mathbf{x}_i &\leq 0, & \text{for all } i \in S_3, \\ \mathbf{w}_\oplus^\top \mathbf{x}_i - \mathbf{w}_\ominus^\top \mathbf{x}_i &= y_i, & -\mathbf{w}_\oplus^\top \mathbf{x}_i &\leq 0, & & \text{for all } i \in S_4. \end{aligned}$$

The Lagrange function in terms of  $\delta \in \mathbb{R}^n$  and non-negative  $\mu \in \mathbb{R}_+^n$  is given by

$$\begin{aligned} \mathcal{L} = & \frac{1}{2} \|\mathbf{w}_\oplus\|_2^2 + \frac{1}{2} \|\mathbf{w}_\ominus\|_2^2 + \sum_{i \in S_1} \delta_i \left( \mathbf{w}_\oplus^\top \mathbf{x}_i - y_i \right) + \sum_{i \in S_1} \mu_i \left( \mathbf{w}_\ominus^\top \mathbf{x}_i \right) \\ & + \sum_{i \in S_2} \delta_i \left( \mathbf{w}_\oplus^\top \mathbf{x}_i - \mathbf{w}_\ominus^\top \mathbf{x}_i - y_i \right) - \sum_{i \in S_2} \mu_i \left( \mathbf{w}_\ominus^\top \mathbf{x}_i \right) \\ & + \sum_{i \in S_3} \delta_i \left( -\mathbf{w}_\ominus^\top \mathbf{x}_i - y_i \right) + \sum_{i \in S_3} \mu_i \left( \mathbf{w}_\oplus^\top \mathbf{x}_i \right) \\ & + \sum_{i \in S_4} \delta_i \left( \mathbf{w}_\oplus^\top \mathbf{x}_i - \mathbf{w}_\ominus^\top \mathbf{x}_i - y_i \right) - \sum_{i \in S_4} \mu_i \left( \mathbf{w}_\oplus^\top \mathbf{x}_i \right). \end{aligned}$$

The KKT conditions are given below.

**Stationarity:**

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{w}_\oplus} &= \mathbf{w}_\oplus^* + \sum_{i \in S_1} \delta_i^* \mathbf{x}_i + \sum_{i \in S_2} \delta_i^* \mathbf{x}_i + \sum_{i \in S_3} \mu_i^* \mathbf{x}_i + \sum_{i \in S_4} (\delta_i^* - \mu_i^*) \mathbf{x}_i = 0, \\ \frac{\partial \mathcal{L}}{\partial \mathbf{w}_\ominus} &= \mathbf{w}_\ominus^* + \sum_{i \in S_1} \mu_i^* \mathbf{x}_i - \sum_{i \in S_2} (\delta_i^* + \mu_i^*) \mathbf{x}_i - \sum_{i \in S_3} \delta_i^* \mathbf{x}_i - \sum_{i \in S_4} \delta_i^* \mathbf{x}_i = 0, \\ \Leftrightarrow \mathbf{w}_\oplus^* &= - \sum_{i \in S_1} \delta_i^* \mathbf{x}_i - \sum_{i \in S_2} \delta_i^* \mathbf{x}_i - \sum_{i \in S_3} \mu_i^* \mathbf{x}_i + \sum_{i \in S_4} (-\delta_i^* + \mu_i^*) \mathbf{x}_i, \end{aligned} \quad (41)$$

$$\mathbf{w}_\ominus^* = - \sum_{i \in S_1} \mu_i^* \mathbf{x}_i + \sum_{i \in S_2} (\delta_i^* + \mu_i^*) \mathbf{x}_i + \sum_{i \in S_3} \delta_i^* \mathbf{x}_i + \sum_{i \in S_4} \delta_i^* \mathbf{x}_i. \quad (42)$$

**Primal feasibility:**

$$\begin{aligned} \mathbf{w}_\oplus^{*\top} \mathbf{x}_i &= y_i, & \mathbf{w}_\ominus^{*\top} \mathbf{x}_i &\leq 0, & \text{for all } i \in S_1, \\ \mathbf{w}_\oplus^{*\top} \mathbf{x}_i - \mathbf{w}_\ominus^{*\top} \mathbf{x}_i &= y_i, & -\mathbf{w}_\ominus^{*\top} \mathbf{x}_i &\leq 0, & \text{for all } i \in S_2, \\ -\mathbf{w}_\ominus^{*\top} \mathbf{x}_i &= y_i, & \mathbf{w}_\oplus^{*\top} \mathbf{x}_i &\leq 0, & \text{for all } i \in S_3, \\ \mathbf{w}_\oplus^{*\top} \mathbf{x}_i - \mathbf{w}_\ominus^{*\top} \mathbf{x}_i &= y_i, & -\mathbf{w}_\oplus^{*\top} \mathbf{x}_i &\leq 0, & \text{for all } i \in S_4. \end{aligned}$$

**Dual feasibility:**

$$\begin{aligned} \delta_i^* &\in \mathbb{R}, \text{ for all } i \in [n], \\ \mu_i^* &\geq 0, \text{ for all } i \in [n]. \end{aligned}$$

**Complementary slackness:**

$$\sum_{i \in S_1} \mu_i^* \left( \mathbf{w}_\ominus^{*\top} \mathbf{x}_i \right) + \sum_{i \in S_2} \mu_i^* \left( -\mathbf{w}_\ominus^{*\top} \mathbf{x}_i \right) + \sum_{i \in S_3} \mu_i^* \left( \mathbf{w}_\oplus^{*\top} \mathbf{x}_i \right) + \sum_{i \in S_4} \mu_i^* \left( -\mathbf{w}_\oplus^{*\top} \mathbf{x}_i \right) = 0.$$

Note that the representation of  $\mathbf{w}_\oplus^*$  and  $\mathbf{w}_\ominus^*$  shares the parameters  $\{\delta_i^* : i \in S_2 \cup S_4\}$ . As a result, since we define  $\mathbf{w}_\oplus^* = \mathbf{X}^\top \alpha_\oplus^*$  and  $\mathbf{w}_\ominus^* = \mathbf{X}^\top \alpha_\ominus^*$ , we can write  $\alpha_{\oplus,i}^*$  and  $\alpha_{\ominus,i}^*$  in terms of  $\delta_i$  and  $\mu_i$  for all  $i \in [n]$  by Equations (41) and (42) as

$$\alpha_{\oplus,i}^* = \begin{cases} -\delta_i^* & \forall i \in S_1 \\ -\delta_i^* & \forall i \in S_2 \\ -\mu_i^* & \forall i \in S_3 \\ -\delta_i^* + \mu_i^* & \forall i \in S_4 \end{cases}, \text{ and } \alpha_{\ominus,i}^* = \begin{cases} -\mu_i^* & \forall i \in S_1 \\ \delta_i^* + \mu_i^* & \forall i \in S_2 \\ \delta_i^* & \forall i \in S_3 \\ \delta_i^* & \forall i \in S_4 \end{cases}.$$

**Step 3: Upper bound for**  $\left\| \mathbf{w}_{\oplus}^{(\infty)} - \mathbf{w}_{\oplus}^* \right\|_2$ .

We now generalize the proof of Theorem 6. We first relate the distance between the predictors  $\mathbf{w}_{\oplus}^{(\infty)}$  and  $\mathbf{w}_{\oplus}^*$  to the distance in their predictions, i.e.,  $\left\| \mathbf{X} \mathbf{w}_{\oplus}^{(\infty)} - \mathbf{X} \mathbf{w}_{\oplus}^* \right\|_2$ . Since both vectors lie in the span of the data  $\{\mathbf{x}_i\}_{i=1}^n$ , their difference has no component in the null space corresponding to the smallest  $d - n$  eigenvalues of  $\mathbf{X}^\top \mathbf{X}$ . Therefore, we have

$$\begin{aligned} \left\| \mathbf{X} \mathbf{w}_{\oplus}^{(\infty)} - \mathbf{X} \mathbf{w}_{\oplus}^* \right\|_2^2 &= \left\| \mathbf{X} \left( \mathbf{w}_{\oplus}^{(\infty)} - \mathbf{w}_{\oplus}^* \right) \right\|_2^2 \geq \mu_n(\mathbf{X}^\top \mathbf{X}) \left\| \mathbf{w}_{\oplus}^{(\infty)} - \mathbf{w}_{\oplus}^* \right\|_2^2 \\ &= \mu_n(\mathbf{X} \mathbf{X}^\top) \left\| \mathbf{w}_{\oplus}^{(\infty)} - \mathbf{w}_{\oplus}^* \right\|_2^2. \end{aligned} \quad (43)$$

As a result, to derive an upper bound for  $\left\| \mathbf{w}_{\oplus}^{(\infty)} - \mathbf{w}_{\oplus}^* \right\|_2$ , it suffices to upper bound the distance between their predictions  $\left\| \mathbf{X} \mathbf{w}_{\oplus}^{(\infty)} - \mathbf{X} \mathbf{w}_{\oplus}^* \right\|_2$ . We begin with analyzing  $\mathbf{w}_{\oplus}^{(\infty)}$ . By Theorem 8,  $\mathbf{w}_{\oplus}^{(\infty)}$  satisfies the following

$$\mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i = y_i \quad \text{for all } y_i > 0, \quad (44a)$$

$$\alpha_{\oplus,j}^{(\infty)} = \alpha_{\oplus,j}^{(1)} = \eta(y_j - \epsilon_{\oplus,j} + \epsilon_{\ominus,j} + \frac{1}{\eta} \mathbf{e}_j^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \boldsymbol{\epsilon}_{\oplus}) \quad \text{for all } y_j < 0, \quad (44b)$$

and also all the conditions in Lemma 19.

We know that  $\mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i = \mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i = y_i$  for all  $i \in S_1$ , and  $\mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i = y_i$  and  $\mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i = y_i + \mathbf{w}_{\ominus}^{*\top} \mathbf{x}_i$  for all  $i \in S_2$ . Therefore, we can write

$$\begin{aligned} &\left\| \mathbf{X} \mathbf{w}_{\oplus}^{(\infty)} - \mathbf{X} \mathbf{w}_{\oplus}^* \right\|_2^2 \\ &= \sum_{i=1}^n \left( \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i \right)^2 \\ &= \sum_{i \in S_2} \left( -\mathbf{w}_{\ominus}^{*\top} \mathbf{x}_i \right)^2 + \sum_{i \in S_3} \left( \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i \right)^2 + \sum_{i \in S_4} \left( \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i \right)^2. \end{aligned} \quad (45)$$

We start with upper bounding the term  $(-\mathbf{w}_{\ominus}^{*\top} \mathbf{x}_i)^2$  for all  $i \in S_2$ . For  $i \in S_2$ , by the complementary slackness, we either have  $-\mathbf{w}_{\ominus}^{*\top} \mathbf{x}_i = 0$  with  $\mu_i^* \geq 0$  or  $-\mathbf{w}_{\ominus}^{*\top} \mathbf{x}_i \leq 0$  with  $\mu_i^* = 0$ . In the first case, we have  $(-\mathbf{w}_{\ominus}^{*\top} \mathbf{x}_i)^2 = 0$ . In the second case, we define  $\tilde{S}_2 \subseteq S_2$  such that  $\mu_i^* = 0$  and  $-\mathbf{w}_{\ominus}^{*\top} \mathbf{x}_i \leq 0$  for all  $i \in \tilde{S}_2$ , and we will show that  $\tilde{S}_2$  is empty with probability at least  $1 - 2 \exp(-n(Cc - \ln 9))$ . Since  $\mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i - \mathbf{w}_{\ominus}^{*\top} \mathbf{x}_i = y_i$  for all  $i \in \tilde{S}_2 \subseteq S_2$ , we have

$$\begin{aligned} y_i &= \mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i - \mathbf{w}_{\ominus}^{*\top} \mathbf{x}_i = \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top (\boldsymbol{\alpha}_{\oplus}^* - \boldsymbol{\alpha}_{\ominus}^*) \\ &= \|\boldsymbol{\lambda}\|_1 (-2\delta_i^* - \underbrace{\mu_i^*}_{=0}) + \mathbf{e}_i^\top (\mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}) (\boldsymbol{\alpha}_{\oplus}^* - \boldsymbol{\alpha}_{\ominus}^*). \end{aligned}$$

The, for all  $i \in \tilde{S}_2$ , we have

$$\delta_i^* = \frac{y_i - \mathbf{e}_i^\top (\mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}) (\boldsymbol{\alpha}_{\oplus}^* - \boldsymbol{\alpha}_{\ominus}^*)}{-2 \|\boldsymbol{\lambda}\|_1}.$$

Based on this representation of  $\delta_i^*$ , for all  $i \in \tilde{S}_2$ , we have

$$\begin{aligned}
 \mathbf{w}_\ominus^{*\top} \mathbf{x}_i &= \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top \boldsymbol{\alpha}_\ominus^* \\
 &= \|\boldsymbol{\lambda}\|_1 (\delta_i^*) + \mathbf{e}_i^\top \left( \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right) \boldsymbol{\alpha}_\ominus^* \\
 &= -\frac{y_i}{2} + \frac{1}{2} \mathbf{e}_i^\top \left( \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right) (\boldsymbol{\alpha}_\oplus^* + \boldsymbol{\alpha}_\ominus^*) \\
 &\leq -\frac{y_i}{2} + \frac{1}{2} \left\| \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right\|_2 (\|\boldsymbol{\alpha}_\oplus^*\|_2 + \|\boldsymbol{\alpha}_\ominus^*\|_2) \\
 &\leq \|\boldsymbol{\lambda}\|_1 \left[ -\frac{y_i}{2 \|\boldsymbol{\lambda}\|_1} + \frac{1}{2} C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \cdot \frac{2\sqrt{2} C_g \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1} \right] \\
 &\leq \|\boldsymbol{\lambda}\|_1 \left[ -\frac{y_{\min}}{2 \|\boldsymbol{\lambda}\|_1} + \frac{1}{2} C \cdot \frac{y_{\min}}{C_0 y_{\max}} \cdot \frac{2\sqrt{2} C_g y_{\max}}{\|\boldsymbol{\lambda}\|_1} \right] \\
 &< 0,
 \end{aligned}$$

where the inequalities above apply Corollary 14 and the upper bound of  $\|\boldsymbol{\alpha}_\oplus^*\|_2, \|\boldsymbol{\alpha}_\ominus^*\|_2$  from Step 1, and substitute  $d_2 \geq C_0^2 \frac{n^2 y_{\max}^2}{y_{\min}^2}$  and  $d_\infty \geq C_0 \frac{n^{1.5} y_{\max}}{y_{\min}}$  in Assumption 2 with  $C_0 \gtrsim C_\alpha^2$  and  $C_\alpha \gtrsim \max\{C_g^2, C_y C_g\}$ . However,  $\mathbf{w}_\ominus^{*\top} \mathbf{x}_i < 0$  contradicts with the condition  $-\mathbf{w}_\ominus^{*\top} \mathbf{x}_i \leq 0$  for  $i \in \tilde{S}_2$ . Therefore,  $\tilde{S}_2 = \emptyset$ . By combining these two cases, we conclude that  $\sum_{i \in S_2} (-\mathbf{w}_\ominus^{*\top} \mathbf{x}_i)^2 = 0$ .

Next, we upper bound the term  $(\mathbf{w}_\oplus^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_\oplus^{*\top} \mathbf{x}_i)^2$  for all  $i \in S_3$ . We know that  $\mathbf{w}_\oplus^{(\infty)\top} \mathbf{x}_i < 0$  in Theorem 8. For  $\mathbf{w}_\oplus^{*\top} \mathbf{x}_i$  with  $i \in S_3$ , by the complementary slackness, we either have  $\mathbf{w}_\oplus^{*\top} \mathbf{x}_i = 0$  with  $\mu_i^* \geq 0$  or  $\mu_i^* = 0$  with  $\mathbf{w}_\oplus^{*\top} \mathbf{x}_i \leq 0$ . In the first case,  $\mathbf{w}_\oplus^{*\top} \mathbf{x}_i = 0$ , we have

$$\begin{aligned}
 \mathbf{w}_\oplus^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_\oplus^{*\top} \mathbf{x}_i &= \mathbf{w}_\oplus^{(\infty)\top} \mathbf{x}_i \\
 &= \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top \boldsymbol{\alpha}_\oplus^{(\infty)} \\
 &= \mathbf{e}_i^\top \left[ \|\boldsymbol{\lambda}\|_1 \mathbf{I} + \left( \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right) \right] \boldsymbol{\alpha}_\oplus^{(\infty)} \\
 &= \|\boldsymbol{\lambda}\|_1 \alpha_{\oplus,i}^{(\infty)} + \mathbf{e}_i^\top \left( \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right) \boldsymbol{\alpha}_\oplus^{(\infty)} \\
 &\geq \|\boldsymbol{\lambda}\|_1 \alpha_{\oplus,i}^{(\infty)} - \left\| \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right\|_2 \left\| \boldsymbol{\alpha}_\oplus^{(\infty)} \right\|_2 \\
 &\geq \|\boldsymbol{\lambda}\|_1 \left[ \alpha_{\oplus,i}^{(\infty)} - C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \left\| \boldsymbol{\alpha}_\oplus^{(\infty)} \right\|_2 \right],
 \end{aligned}$$

where the last inequality applies Corollary 14. Substituting the bounds of  $\alpha_{\oplus,i}^{(\infty)}$  and  $\left\| \boldsymbol{\alpha}_\oplus^{(\infty)} \right\|_2$  from Lemma 19, we have

$$\begin{aligned}
 \mathbf{w}_\oplus^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_\oplus^{*\top} \mathbf{x}_i &\geq \|\boldsymbol{\lambda}\|_1 \left[ -\frac{3y_{\max}}{C_g \|\boldsymbol{\lambda}\|_1} - C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \frac{C_\alpha y_{\max}}{\|\boldsymbol{\lambda}\|_1} \right] \\
 &\geq \|\boldsymbol{\lambda}\|_1 \left[ -\frac{3y_{\max}}{C_g \|\boldsymbol{\lambda}\|_1} - C \cdot \frac{y_{\min}}{C_0 y_{\max}} \frac{C_\alpha y_{\max}}{\|\boldsymbol{\lambda}\|_1} \right] \\
 &\geq -\frac{4}{C_g} y_{\max},
 \end{aligned}$$

where the inequalities substitute  $d_2 \geq C_0^2 \frac{n^2 y_{\max}^2}{y_{\min}^2}$  and  $d_\infty \geq C_0 \frac{n^{1.5} y_{\max}}{y_{\min}}$  in Assumption 2 with  $C_0 \gtrsim C_\alpha^2$  and  $C_\alpha \gtrsim \max\{C_g^2, C_y C_g\}$ . In the second case,  $\alpha_{\oplus, i}^* = -\mu_i^* = 0$  for  $i \in S_3$ , we have

$$\begin{aligned}
 \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{\star\top} \mathbf{x}_i &= \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top (\boldsymbol{\alpha}_{\oplus}^{(\infty)} - \boldsymbol{\alpha}_{\oplus}^*) \\
 &= \mathbf{e}_i^\top \left[ \|\boldsymbol{\lambda}\|_1 \mathbf{I} + (\mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}) \right] (\boldsymbol{\alpha}_{\oplus}^{(\infty)} - \boldsymbol{\alpha}_{\oplus}^*) \\
 &= \|\boldsymbol{\lambda}\|_1 \boldsymbol{\alpha}_{\oplus, i}^{(\infty)} + \mathbf{e}_i^\top (\mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}) (\boldsymbol{\alpha}_{\oplus}^{(\infty)} - \boldsymbol{\alpha}_{\oplus}^*) \\
 &\geq \|\boldsymbol{\lambda}\|_1 \boldsymbol{\alpha}_{\oplus, i}^{(\infty)} - \|\mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}\|_2 \left( \|\boldsymbol{\alpha}_{\oplus}^{(\infty)}\|_2 + \|\boldsymbol{\alpha}_{\oplus}^*\|_2 \right) \\
 &\geq \|\boldsymbol{\lambda}\|_1 \left[ -\frac{3y_{\max}}{C_g \|\boldsymbol{\lambda}\|_1} - C \cdot \max\left(\sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty}\right) \left( \frac{C_\alpha \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1} + \frac{\sqrt{2} C_g \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1} \right) \right] \\
 &\geq \|\boldsymbol{\lambda}\|_1 \left[ -\frac{3y_{\max}}{C_g \|\boldsymbol{\lambda}\|_1} - C \cdot \frac{y_{\min}}{C_0 y_{\max}} \left( \frac{C_\alpha y_{\max}}{\|\boldsymbol{\lambda}\|_1} + \frac{\sqrt{2} C_g y_{\max}}{\|\boldsymbol{\lambda}\|_1} \right) \right] \\
 &\geq -\frac{4}{C_g} y_{\max},
 \end{aligned}$$

by applying the same argument and the upper bound from Step 1 that  $\|\boldsymbol{\alpha}_{\oplus}^*\|_2 \leq \frac{\sqrt{2} C_g \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1}$ .

Therefore, we have  $(\mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{\star\top} \mathbf{x}_i)^2 \leq \frac{16}{C_g^2} y_{\max}^2$  for all  $i \in S_3$ .

Next, we upper bound the term  $(\mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{\star\top} \mathbf{x}_i)^2$  for all  $i \in S_4$  in a similar way compared to  $S_2$ . For  $i \in S_4$ , by the complementary slackness, we either have  $-\mathbf{w}_{\oplus}^{\star\top} \mathbf{x}_i = 0$  with  $\mu_i^* \geq 0$  or  $-\mathbf{w}_{\oplus}^{\star\top} \mathbf{x}_i \leq 0$  with  $\mu_i^* = 0$ . In the first case,  $(-\mathbf{w}_{\oplus}^{\star\top} \mathbf{x}_i)^2 = 0$ , and we have  $(\mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{\star\top} \mathbf{x}_i)^2 = (\mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i)^2$ . Therefore, we can reuse the upper bound we derived in  $S_3$  such that  $0 \geq \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i \geq -\frac{4}{C_g} y_{\max}$ . In the second case, we define  $\tilde{S}_4 \subseteq S_4$  such that  $\mu_i^* = 0$  and  $-\mathbf{w}_{\oplus}^{\star\top} \mathbf{x}_i \leq 0$  for all  $i \in \tilde{S}_4$ , and we will show that  $\tilde{S}_4$  is empty with probability at least  $1 - 2 \exp(-n(Cc - \ln 9))$ . Since  $\mathbf{w}_{\oplus}^{\star\top} \mathbf{x}_i - \mathbf{w}_{\ominus}^{\star\top} \mathbf{x}_i = y_i$  for all  $i \in \tilde{S}_4 \subseteq S_4$ , we have

$$\begin{aligned}
 y_i = \mathbf{w}_{\oplus}^{\star\top} \mathbf{x}_i - \mathbf{w}_{\ominus}^{\star\top} \mathbf{x}_i &= \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top (\boldsymbol{\alpha}_{\oplus}^* - \boldsymbol{\alpha}_{\ominus}^*) \\
 &= \|\boldsymbol{\lambda}\|_1 (-2\delta_i^* + \underbrace{\mu_i^*}_{=0}) + \mathbf{e}_i^\top (\mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}) (\boldsymbol{\alpha}_{\oplus}^* - \boldsymbol{\alpha}_{\ominus}^*).
 \end{aligned}$$

For all  $i \in \tilde{S}_4$ , we have

$$\delta_i^* = \frac{y_i - \mathbf{e}_i^\top (\mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}) (\boldsymbol{\alpha}_{\oplus}^* - \boldsymbol{\alpha}_{\ominus}^*)}{-2 \|\boldsymbol{\lambda}\|_1}.$$

Based on this representation of  $\delta_i^*$ , for all  $i \in \tilde{S}_4$ , we have

$$\begin{aligned}
 \mathbf{w}_\oplus^{\star\top} \mathbf{x}_i &= \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top \boldsymbol{\alpha}_\oplus^* \\
 &= \|\boldsymbol{\lambda}\|_1 (-\delta_i^*) + \mathbf{e}_i^\top \left( \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right) \boldsymbol{\alpha}_\oplus^* \\
 &= \frac{y_i}{2} + \frac{1}{2} \mathbf{e}_i^\top \left( \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right) (\boldsymbol{\alpha}_\oplus^* + \boldsymbol{\alpha}_\ominus^*) \\
 &\leq \frac{y_i}{2} + \frac{1}{2} \left\| \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right\|_2 (\|\boldsymbol{\alpha}_\oplus^*\|_2 + \|\boldsymbol{\alpha}_\ominus^*\|_2) \\
 &\leq \|\boldsymbol{\lambda}\|_1 \left[ \frac{y_i}{2 \|\boldsymbol{\lambda}\|_1} + \frac{1}{2} C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \cdot \frac{2\sqrt{2} C_g \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1} \right] \\
 &\leq \|\boldsymbol{\lambda}\|_1 \left[ -\frac{y_{\min}}{2 \|\boldsymbol{\lambda}\|_1} + \frac{1}{2} C \cdot \frac{y_{\min}}{C_0 y_{\max}} \cdot \frac{2\sqrt{2} C_g y_{\max}}{\|\boldsymbol{\lambda}\|_1} \right] \\
 &< 0,
 \end{aligned}$$

where the inequalities apply Corollary 14 and the upper bound of  $\|\boldsymbol{\alpha}_\oplus^*\|_2$  in Step 1, and substitute  $d_2 \geq C_0^2 \frac{n^2 y_{\max}^2}{y_{\min}^2}$  and  $d_\infty \geq C_0 \frac{n^{1.5} y_{\max}}{y_{\min}}$  in Assumption 2 with  $C_0 \gtrsim C_\alpha^2$  and  $C_\alpha \gtrsim \max\{C_g^2, C_y C_g\}$ . However,  $\mathbf{w}_\oplus^{\star\top} \mathbf{x}_i < 0$  contradicts with the condition  $-\mathbf{w}_\oplus^{\star\top} \mathbf{x}_i \leq 0$  for  $i \in \tilde{S}_4$ . Therefore,  $\tilde{S}_4 = \emptyset$ . By combining these two cases, we have  $\sum_{i \in S_4} \left( \mathbf{w}_\oplus^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_\oplus^{\star\top} \mathbf{x}_i \right)^2 = \sum_{i \in S_4} \left( \mathbf{w}_\oplus^{(\infty)\top} \mathbf{x}_i \right)^2 \leq \sum_{i \in S_4} \frac{16}{C_g^2} y_{\max}^2$ .

Substituting the upper bounds into Equation (45) gives us

$$\begin{aligned}
 &\left\| \mathbf{X} \mathbf{w}^{(\infty)} - \mathbf{X} \mathbf{w}^* \right\|_2^2 \\
 &= \sum_{i \in S_2} \left( -\mathbf{w}_\ominus^{\star\top} \mathbf{x}_i \right)^2 + \sum_{i \in S_3} \left( \mathbf{w}_\oplus^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_\oplus^{\star\top} \mathbf{x}_i \right)^2 + \sum_{i \in S_4} \left( \mathbf{w}_\oplus^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_\oplus^{\star\top} \mathbf{x}_i \right)^2 \\
 &\leq \sum_{i \in S_3} \frac{16}{C_g^2} y_{\max}^2 + \sum_{i \in S_4} \frac{16}{C_g^2} y_{\max}^2 \\
 &= \frac{16}{C_g^2} n - y_{\max}^2.
 \end{aligned} \tag{46}$$

Finally, putting together Equation (43) and (46), we have

$$\left\| \mathbf{w}^{(\infty)} - \mathbf{w}^* \right\|_2^2 \leq \frac{\left\| \mathbf{X} \mathbf{w}^{(\infty)} - \mathbf{X} \mathbf{w}^* \right\|_2^2}{\mu_n(\mathbf{X} \mathbf{X}^\top)} \leq \frac{16n - y_{\max}^2}{C_g \|\boldsymbol{\lambda}\|_1},$$

which completes the proof of the upper bound.

**Step 4: Lower bound for  $\left\| \mathbf{w}_\oplus^{(\infty)} - \mathbf{w}_\oplus^* \right\|_2$ .**

Now, we derive the lower bound of  $\left\| \mathbf{w}_\oplus^{(\infty)} - \mathbf{w}_\oplus^* \right\|_2$  in a similar approach. We again start with the

prediction distance

$$\begin{aligned} \left\| \mathbf{X} \mathbf{w}_{\oplus}^{(\infty)} - \mathbf{X} \mathbf{w}_{\oplus}^* \right\|_2^2 &= \left\| \mathbf{X} \left( \mathbf{w}_{\oplus}^{(\infty)} - \mathbf{w}_{\oplus}^* \right) \right\|_2^2 \leq \mu_1(\mathbf{X}^\top \mathbf{X}) \left\| \mathbf{w}_{\oplus}^{(\infty)} - \mathbf{w}_{\oplus}^* \right\|_2^2 \\ &= \mu_1(\mathbf{X} \mathbf{X}^\top) \left\| \mathbf{w}_{\oplus}^{(\infty)} - \mathbf{w}_{\oplus}^* \right\|_2^2. \end{aligned} \quad (47)$$

Therefore, it suffices to lower bound  $\left\| \mathbf{X} \mathbf{w}_{\oplus}^{(\infty)} - \mathbf{X} \mathbf{w}_{\oplus}^* \right\|_2$  to get the lower bound of  $\left\| \mathbf{w}_{\oplus}^{(\infty)} - \mathbf{w}_{\oplus}^* \right\|_2$ . By Equation (45), we have

$$\begin{aligned} &\left\| \mathbf{X} \mathbf{w}_{\oplus}^{(\infty)} - \mathbf{X} \mathbf{w}_{\oplus}^* \right\|_2^2 \\ &= \sum_{i \in S_2} \left( -\mathbf{w}_{\ominus}^{*\top} \mathbf{x}_i \right)^2 + \sum_{i \in S_3} \left( \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i \right)^2 + \sum_{i \in S_4} \left( \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i \right)^2 \\ &\geq \sum_{i \in S_3} \left( \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i \right)^2 + \sum_{i \in S_4} \left( \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i \right)^2. \end{aligned} \quad (48)$$

We omit the partition in  $S_2$  because we have shown that  $\sum_{i \in S_2} \left( -\mathbf{w}_{\ominus}^{*\top} \mathbf{x}_i \right)^2 = 0$  with probability at least  $1 - 2 \exp(-n(Cc - \ln 9))$ . Therefore, we need to lower bound  $\left( \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i \right)^2$  for  $i \in S_3$  and  $\left( \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i \right)^2$  for  $i \in S_4$ .

We start with lower bounding  $\left( \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i \right)^2$  for  $i \in S_3$ . We know that  $\mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i < 0$  in Theorem 8. For  $\mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i$  with  $i \in S_3$ , by the complementary slackness, we either have  $\mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i = 0$  with  $\mu_i^* \geq 0$  or  $\mu_i^* = 0$  with  $\mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i \leq 0$ . In the first case,  $\mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i = 0$ , we have

$$\begin{aligned} \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i &= \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i \\ &= \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top \boldsymbol{\alpha}_{\oplus}^{(\infty)} \\ &= \mathbf{e}_i^\top \left[ \|\boldsymbol{\lambda}\|_1 \mathbf{I} + \left( \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right) \right] \boldsymbol{\alpha}_{\oplus}^{(\infty)} \\ &= \|\boldsymbol{\lambda}\|_1 \alpha_{\oplus, i}^{(\infty)} + \mathbf{e}_i^\top \left( \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right) \boldsymbol{\alpha}_{\oplus}^{(\infty)} \\ &\leq \|\boldsymbol{\lambda}\|_1 \alpha_{\oplus, i}^{(\infty)} + \left\| \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right\|_2 \left\| \boldsymbol{\alpha}_{\oplus}^{(\infty)} \right\|_2 \\ &\leq \|\boldsymbol{\lambda}\|_1 \left[ \alpha_{\oplus, i}^{(\infty)} + C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \left\| \boldsymbol{\alpha}_{\oplus}^{(\infty)} \right\|_2 \right], \end{aligned}$$

where the last inequality applies Corollary 14. Substituting the bounds of  $\alpha_{\oplus, i}^{(\infty)}$  and  $\left\| \boldsymbol{\alpha}_{\oplus}^{(\infty)} \right\|_2$  from Lemma 19, we have

$$\begin{aligned} \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{*\top} \mathbf{x}_i &\leq \|\boldsymbol{\lambda}\|_1 \left[ -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1} + C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \frac{C_\alpha \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1} \right] \\ &\leq \|\boldsymbol{\lambda}\|_1 \left[ -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1} + C \cdot \frac{y_{\min}}{C_0 y_{\max}} \frac{C_\alpha y_{\max}}{\|\boldsymbol{\lambda}\|_1} \right] \\ &\leq -\left(1 - \frac{C \cdot C_\alpha^2}{C_0}\right) \frac{y_{\min}}{C_\alpha}, \end{aligned}$$

where the inequalities substitute  $d_2 \geq C_0^2 \frac{n^2 y_{\max}^2}{y_{\min}^2}$  and  $d_\infty \geq C_0 \frac{n^{1.5} y_{\max}}{y_{\min}}$  in Assumption 2 with  $C_0 \gtrsim C_\alpha^2$ . In the second case,  $\alpha_{\oplus, i}^* = -\mu_i^* = 0$  for  $i \in S_3$ , we have

$$\begin{aligned}
 & \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{\star\top} \mathbf{x}_i \\
 &= \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top \left( \boldsymbol{\alpha}_{\oplus}^{(\infty)} - \boldsymbol{\alpha}_{\oplus}^* \right) \\
 &= \mathbf{e}_i^\top \left[ \|\boldsymbol{\lambda}\|_1 \mathbf{I} + \left( \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right) \right] \left( \boldsymbol{\alpha}_{\oplus}^{(\infty)} - \boldsymbol{\alpha}_{\oplus}^* \right) \\
 &= \|\boldsymbol{\lambda}\|_1 \alpha_{\oplus, i}^{(\infty)} + \mathbf{e}_i^\top \left( \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right) \left( \boldsymbol{\alpha}_{\oplus}^{(\infty)} - \boldsymbol{\alpha}_{\oplus}^* \right) \\
 &\leq \|\boldsymbol{\lambda}\|_1 \alpha_{\oplus, i}^{(\infty)} + \left\| \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right\|_2 \left( \left\| \boldsymbol{\alpha}_{\oplus}^{(\infty)} \right\|_2 + \left\| \boldsymbol{\alpha}_{\oplus}^* \right\|_2 \right) \\
 &\leq \|\boldsymbol{\lambda}\|_1 \left[ -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1} + C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \left( \frac{C_\alpha \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1} + \frac{\sqrt{2} C_g \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1} \right) \right] \\
 &\leq \|\boldsymbol{\lambda}\|_1 \left[ -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1} + C \cdot \frac{y_{\min}}{C_0 y_{\max}} \left( \frac{C_\alpha y_{\max}}{\|\boldsymbol{\lambda}\|_1} + \frac{\sqrt{2} C_g y_{\max}}{\|\boldsymbol{\lambda}\|_1} \right) \right] \\
 &\leq -\left( 1 - \frac{2C \cdot C_\alpha^2}{C_0} \right) \frac{y_{\min}}{C_\alpha},
 \end{aligned}$$

by applying the same argument and the upper bound from Step 1 that  $\|\boldsymbol{\alpha}_{\oplus}^*\|_2 \leq \frac{\sqrt{2} C_g \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1}$ .

Therefore, we have  $\left( \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{\star\top} \mathbf{x}_i \right)^2 \geq \left( 1 - \frac{2C \cdot C_\alpha^2}{C_0} \right)^2 \frac{y_{\min}^2}{C_\alpha^2}$ , for all  $i \in S_3$ .

Next, we lower bound the term  $\left( \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{\star\top} \mathbf{x}_i \right)^2$  for all  $i \in S_4$ . In Step 3, we already showed the two cases in  $\mathbf{w}_{\oplus}^{\star\top} \mathbf{x}_i$  with  $i \in S_4$  by the complementary slackness. In the first case,  $(-\mathbf{w}_{\oplus}^{\star\top} \mathbf{x}_i)^2 = 0$ , and we have  $\left( \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{\star\top} \mathbf{x}_i \right)^2 = \left( \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i \right)^2$ . Therefore, we can reuse the lower bound we derived in  $S_3$  such that  $\mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i \leq -\left( 1 - \frac{C \cdot C_\alpha^2}{C_0} \right) \frac{y_{\min}}{C_\alpha}$ . In the second case, we have shown that  $\tilde{S}_4 = \emptyset$ . By concluding two cases, we have  $\sum_{i \in S_4} \left( \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{\star\top} \mathbf{x}_i \right)^2 = \sum_{i \in S_4} \left( \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i \right)^2 \geq \sum_{i \in S_4} \left( 1 - \frac{2C \cdot C_\alpha^2}{C_0} \right)^2 \frac{y_{\min}^2}{C_\alpha^2}$ .

Substituting the lower bounds into Equation (48) gives us

$$\begin{aligned}
 \left\| \mathbf{X} \mathbf{w}_{\oplus}^{(\infty)} - \mathbf{X} \mathbf{w}_{\oplus}^* \right\|_2^2 &\geq \sum_{i \in S_3} \left( \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{\star\top} \mathbf{x}_i \right)^2 + \sum_{i \in S_4} \left( \mathbf{w}_{\oplus}^{(\infty)\top} \mathbf{x}_i - \mathbf{w}_{\oplus}^{\star\top} \mathbf{x}_i \right)^2 \\
 &\geq \sum_{i \in S_3} \left( 1 - \frac{2C \cdot C_\alpha^2}{C_0} \right)^2 \frac{y_{\min}^2}{C_\alpha^2} + \sum_{i \in S_4} \left( 1 - \frac{2C \cdot C_\alpha^2}{C_0} \right)^2 \frac{y_{\min}^2}{C_\alpha^2} \\
 &= \frac{n_- y_{\min}^2}{\tilde{C}},
 \end{aligned} \tag{49}$$

where we let  $\tilde{C} := \frac{C_0^2 C_\alpha^2}{(C_0 - 2C \cdot C_\alpha^2)^2} > 1$ . Finally, putting together Equation (47) and (49), we have

$$\left\| \mathbf{w}_{\oplus}^{(\infty)} - \mathbf{w}_{\oplus}^* \right\|_2^2 \geq \frac{\left\| \mathbf{X} \mathbf{w}_{\oplus}^{(\infty)} - \mathbf{X} \mathbf{w}_{\oplus}^* \right\|_2^2}{\mu_1(\mathbf{X} \mathbf{X}^\top)} \geq \frac{n_- y_{\min}^2}{\tilde{C} C_g \|\boldsymbol{\lambda}\|_1}.$$

This completes the proof of the lower bound. ■

## Appendix D. Implicit Bias of Multiple ReLU Models ( $m > 2$ ) Under Gradient Descent

In this section, we extend our analysis to multiple ReLU models trained with  $m > 2$  neurons under stronger assumptions on the initialization. We consider models of the form:  $h_{\Theta}(\mathbf{x}) := h_{\{\mathbf{w}_k\}_{k=1}^m}(\mathbf{x}) = \sum_{k=1}^m s_k \sigma(\mathbf{w}_k^\top \mathbf{x})$ , where  $\mathbf{w}_k \in \mathbb{R}^d$  are the model weights and there are at least one positive neuron and one negative neuron. The parameter set is hence denoted by  $\Theta = \{\mathbf{w}_k\}_{k=1}^m$ . The empirical risk is defined in (1) as

$$\mathcal{R}(\Theta) = \frac{1}{2} \sum_{i=1}^n (h_{\Theta}(\mathbf{x}_i) - y_i)^2 = \frac{1}{2} \sum_{i=1}^n \left( \sum_{k=1}^m s_k \sigma(\mathbf{w}_k^\top \mathbf{x}_i) - y_i \right)^2.$$

Here, we fix  $s_k \in \{\pm 1\}$  and only train the hidden weights  $\{\mathbf{w}_k\}_{k=1}^m$ .

### D.1. Gradient Descent Updates and Convergence

The gradient of the empirical risk in (1) with respect to  $\mathbf{w}_k$  is given in (2) as

$$\mathbf{w}_k^{(t+1)} = \mathbf{w}_k^{(t)} - \eta \nabla_{\mathbf{w}_k} \mathcal{R}(\Theta^{(t)}) = \mathbf{w}_k^{(t)} - \eta s_k \mathbf{X}^\top \mathbf{D}(\mathbf{X} \mathbf{w}_k^{(t)}) (h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}). \quad (50)$$

The primal-dual gradient update in (4) is given by

$$\text{(Primal)} \quad \beta_k^{(t+1)} = \beta_k^{(t)} - \eta s_k \mathbf{X} \mathbf{X}^\top \mathbf{D}(\beta_k^{(t)}) (h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}), \quad (51a)$$

$$\text{(Dual)} \quad \alpha_k^{(t+1)} = \alpha_k^{(t)} - \eta s_k \mathbf{D}(\beta_k^{(t)}) (h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}). \quad (51b)$$

Next, we consider a regime in which, after some time  $t_0$ , each neuron activates on a fixed subset of training examples, and this activation pattern remains unchanged throughout the subsequent dynamics. Moreover, these active subsets are disjoint across different neurons. That is, for every training example, at most one neuron is active, while each neuron may be active on a subset of examples. In this regime, each neuron effectively reduces to a linear model trained only on its own active examples.

**Lemma 20** *Consider a multiple ReLU model  $h_{\Theta}$ . For each neuron  $k \in [m]$ , suppose there exists iteration  $t_0 \geq 0$  such that*

1. *At time  $t_0$ , the subset of examples on which the  $k$ -th neuron is active is disjoint from the subsets activated by all other neurons, i.e.,  $\mathbf{D}(\mathbf{X} \mathbf{w}_k^{(t_0)}) \mathbf{D}(\mathbf{X} \mathbf{w}_\ell^{(t_0)}) = \mathbf{0}_{n \times n}$  for any  $\ell \neq k$ .*
2. *The activation pattern of the  $k$ -th remains unchanged after time  $t_0$ , i.e.,  $\mathbf{D}(\mathbf{X} \mathbf{w}_k^{(t_0)}) = \mathbf{D}(\mathbf{X} \mathbf{w}_k^{(t)})$  for all  $t \geq t_0$ .*

*Then, for all  $t \geq t_0$ , and each  $k \in [m]$ , the gradient descent dynamics of the  $k$ -th neuron are equivalent to gradient descent applied to a linear model, initialized at  $\mathbf{w}_k^{(t_0)}$ , and trained using only the subset of samples satisfying  $\mathbf{x}_i^\top \mathbf{w}_k^{(t_0)} > 0$ .*

The proof of Lemma 20 is provided in Appendix E.1.

## D.2. Minimum- $\ell_2$ -norm Solution of Multiple ReLU Models

The minimum- $\ell_2$ -norm solution for the multiple ReLU regression in (5) is given by

$$\begin{aligned} \{\mathbf{w}_k^*\}_{k=1}^m &= \arg \min_{\{\mathbf{w}_k\}_{k=1}^m} \frac{1}{2} \sum_{k=1}^m \|\mathbf{w}_k\|_2^2 \\ \text{s.t. } \sum_{k=1}^m s_k \sigma(\mathbf{w}_k^\top \mathbf{x}_i) &= y_i, \text{ for all } i \in [n]. \end{aligned} \quad (52)$$

## D.3. High-dimensional Implicit Bias of Multiple ReLU Models

In this section, we characterize the implicit bias of multiple ReLU models trained by gradient descent in the high-dimensional regime. We identify a setup in which each neuron is only active toward a fixed and disjoint subset of training examples, where the labels  $y_i$  of these examples have the same sign as the neuron's sign  $s_k$ .

To formalize this setup, we introduce an assignment vector  $\mathbf{a} \in [m]^n$ , where each entry  $a_i \in [m]$  indicates which neuron is responsible for example  $i$ .

**Assumption 4** *For a multiple ReLU model, we assume that there exists an assignment vector  $\mathbf{a} \in [m]^n$  such that for each example  $i \in [n]$ ,  $a_i = k$ , for some neuron  $k$  satisfying  $s_k \cdot y_i > 0$ . For each neuron  $k \in [m]$ , define a diagonal matrix  $\mathbf{A}_k \in \mathbb{R}^{n \times n}$  with diagonal entries*

$$(\mathbf{A}_k)_{ii} = \begin{cases} 0, & \text{if } a_i = k, \text{ or } s_k \cdot y_i < 0 \\ -\text{sign}(y_i), & \text{otherwise} \end{cases}.$$

Assumption 4 is used to design a proper initialization that ensures that the gradient descent can converge to the desired regime. In this regime, we show that if a neuron's primal variable  $\beta_{k,i}$  is positive and the sign of the neuron agrees with the label (i.e.,  $s_k \cdot y_i > 0$ ), then the corresponding example remains active throughout training. Conversely, if the associated dual variable  $\alpha_{k,j}$  stays sufficiently negative, it remains frozen and is no longer updated.

**Theorem 21** *Under Assumptions 1, 2 and 4, suppose we choose initialization*

$\mathbf{w}_k^{(0)} = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \left( \frac{1}{C_g} \mathbf{A}_k \mathbf{y} + \epsilon_k \right)$ , where  $0 < \epsilon_{k,i} \leq \frac{1}{C_{\alpha m}} y_{\min}$  for all  $k \in [m]$  and  $i \in [n]$ , and the gradient descent step size satisfies  $\frac{1}{C_g \|\lambda\|_1} \leq \eta \leq \frac{1}{C_g \|\lambda\|_1}$ . Then, the gradient descent limit  $\mathbf{w}_k^{(\infty)}$  for multiple ReLU models coincides with the solution obtained by training a linear model on disjoint subsets of examples, initialized at  $\mathbf{w}_k^{(1)}$  with probability at least  $1 - 2 \exp(-cn)$ . Formally, we have  $\mathbf{w}_k^{(\infty)} = \arg \min_{\mathbf{w} \in \{\mathbf{w} : \mathbf{X}_{S_k} \mathbf{w} = \mathbf{y}_{S_k}\}} \|\mathbf{w} - \mathbf{w}_k^{(1)}\|_2$  and  $\mathbf{X}_{S_k} \mathbf{w}_k^{(\infty)} \preceq \mathbf{0}$ , where  $S_k := \{i \in [n] : a_i = k\}$ .

The full proof is provided in Appendix E.2. Note that the initialization, constructed by the matrices  $\mathbf{A}_k$ , ensures that each training example  $i$  is activated by exactly one neuron that matches its sign—namely, the  $a_i$ -th neuron. All other neurons with the same sign remain inactive on this example.

#### D.4. Approximation to Minimum- $\ell_2$ -norm Solution in High Dimensions

In this section, we show that in high dimensions, the implicit bias solution for multiple ReLU models derived in Theorem 21 is close to the corresponding minimum- $\ell_2$ -norm solution  $\{\mathbf{w}_k^*\}_{k=1}^m$  defined in (52).

**Theorem 22** *Under Assumptions 1, 2 and 4, suppose we choose initialization*

$\mathbf{w}_k^{(0)} = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \left( \frac{1}{C_g} \mathbf{A}_k \mathbf{y} + \epsilon_k \right)$ , *where  $0 < \epsilon_{k,i} \leq \frac{1}{C_{\alpha m}} y_{\min}$  for all  $k \in [m]$  and  $i \in [n]$ , and the gradient descent step size satisfies  $\frac{1}{C C_g \|\boldsymbol{\lambda}\|_1} \leq \eta \leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1}$ . Then, we have*

$$\sqrt{\sum_{k=1}^m \left\| \mathbf{w}_k^{(\infty)} - \mathbf{w}_k^* \right\|_2^2} \leq \sqrt{\frac{4 C_g C_{\alpha}^2 m n y_{\max}^2}{\|\boldsymbol{\lambda}\|_1}} \text{ with probability at least } 1 - 2 \exp(-cn).$$

The proof is deferred to Appendix E.3. Note that since the minimum- $\ell_2$ -norm solution  $\{\mathbf{w}_k^*\}_{k=1}^m$  is more involved to characterize, Theorem 22 only provides an upper bound for the approximation of the implicit bias to  $\{\mathbf{w}_k^*\}_{k=1}^m$ . A more fine-grained characterization, as well as a deeper understanding of the role of overparameterization, is left for future work.

## Appendix E. Proofs for Multiple ReLU Models ( $m > 2$ ) Trained with Gradient Descent

In this section, we present the proofs concerning the behavior of the multiple ReLU model trained with gradient descent.

### E.1. Proof of Lemma 20 (Gradient Descent Convergence)

**Proof** (Lemma 20) This proof is analogous to Lemma 1. The key idea is to show that once the activation pattern becomes fixed after some iteration  $t_0 \geq 0$ , the gradient descent dynamics of each neuron are equivalent to those of a linear model trained on a fixed subset of examples.

Fix a neuron  $k \in [m]$ . Consider the linear model

$$h(\mathbf{x}) = s_k \mathbf{w}^\top \mathbf{x},$$

where  $\mathbf{w} \in \mathbb{R}^d$  is the linear model parameter (also called weight). Let  $S_k^{(t_0)} \subseteq [n]$  denote the active set of the  $k$ -th neuron at iteration  $t_0$ , defined by  $S_k^{(t_0)} := \{i \in [n] : \mathbf{x}_i^\top \mathbf{w}_k^{(t_0)} > 0\}$ . We define the empirical risk with the linear model using only the examples in  $S_k^{(t_0)}$  as

$$\mathcal{R}_{S_k^{(t_0)}}(\mathbf{w}) = \frac{1}{2} \sum_{i \in S_k^{(t_0)}} (s_k \mathbf{w}^\top \mathbf{x}_i - y_i)^2.$$

The gradient descent update for this linear model is then given by

$$\begin{aligned} \mathbf{w}^{(t+1)} &= \mathbf{w}^{(t)} - \eta \nabla \mathcal{R}_{S_k^{(t_0)}}(\mathbf{w}^{(t)}) \\ &= \mathbf{w}^{(t)} - \eta s_k \sum_{i \in S_k^{(t_0)}} (s_k \mathbf{w}^{(t)\top} \mathbf{x}_i - y_i) \mathbf{x}_i. \end{aligned} \quad (53)$$

On the other hand, the original gradient descent update of the multiple ReLU model in Equation (50) tells us that

$$\mathbf{w}_k^{(t+1)} = \mathbf{w}_k^{(t)} - \eta s_k \mathbf{X}^\top \mathbf{D}(\mathbf{X} \mathbf{w}_k^{(t)}) (h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}).$$

Under the first assumption in the lemma,  $\mathbf{D}(\mathbf{X} \mathbf{w}_k^{(t_0)}) \mathbf{D}(\mathbf{X} \mathbf{w}_\ell^{(t_0)}) = \mathbf{0}_{n \times n}$  for any  $\ell \neq k$ , the activation patterns of different neurons are disjoint at iteration  $t_0$ . Consequently, for any  $i \in S_k^{(t_0)}$ , only the  $k$ -th neuron is active, and therefore we have

$$h_{\Theta^{(t_0)}}(\mathbf{x}_i) = \sum_{k=1}^m s_k \sigma(\mathbf{w}_k^\top \mathbf{x}_i) = s_k \mathbf{w}_k^\top \mathbf{x}_i.$$

Moreover, by the second assumption of the lemma, the activation pattern of the  $k$ -th neuron remains unchanged after iteration  $t_0$ , i.e.,  $\mathbf{D}(\mathbf{X} \mathbf{w}_k^{(t_0)}) = \mathbf{D}(\mathbf{X} \mathbf{w}_k^{(t)})$  for all  $t \geq t_0$ . Hence, for all  $t \geq t_0$ , the diagonal entries of  $\mathbf{D}(\mathbf{X} \mathbf{w}_k^{(t)})$  satisfy  $D_{ii} = \mathbb{1}_{i \in S_k^{(t_0)}}$  for all  $i \in [n]$ . Therefore, for all  $t \geq t_0$ , the gradient update of the  $k$ -th neuron in the multiple ReLU model is given by

$$\begin{aligned} \mathbf{w}_k^{(t+1)} &= \mathbf{w}_k^{(t)} - \eta s_k \mathbf{X}^\top \mathbf{D}(\mathbf{X} \mathbf{w}_k^{(t)}) (h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}) \\ &= \mathbf{w}_k^{(t)} - \eta s_k \sum_{i \in S_k^{(t_0)}} (s_k \mathbf{w}_k^{(t)\top} \mathbf{x}_i - y_i) \mathbf{x}_i. \end{aligned}$$

This update is identical to the gradient descent update of the linear model in Equation (53). Hence, for all  $t \geq t_0$ , the gradient descent dynamics of the  $k$ -th neuron in the multiple ReLU model are equivalent to those of a linear model trained using only the examples in  $S_k^{(t_0)}$ . This completes the proof of the lemma. ■

## E.2. Proof of Theorem 21 (High-dimensional Implicit Bias)

In this section, we present the proof of Theorem 21. Before proceeding to the proof, we again introduce a set of sufficient conditions under which the active pattern for a neuron at iteration  $t$  will be preserved at iteration  $t + 1$ . Similar to the single ReLU model and 2-ReLU model cases, our analysis relies on Lemma 10 and Lemma 11 to characterize the dynamics of primal and dual variables. Using these results, we establish Lemma 23, which characterizes that the active sets of all neurons remain unchanged across gradient descent iterations.

**Lemma 23** *Under Assumption 1, 2 and 4, suppose the gradient descent step size satisfies  $\eta \leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1}$ . For a multiple ReLU model, if the following five conditions hold at some iteration  $t \geq 0$ , then they also hold at iteration  $t + 1$ .*

- a.  $\beta_{a_i, i}^{(t)} > 0$ , for all  $i \in [n]$ .
- b.  $-\frac{3y_{\max}}{C_g \|\boldsymbol{\lambda}\|_1} \leq \alpha_{k, j}^{(t)} \leq -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1}$ , for all  $j \in [n]$  with  $k \neq a_j$ .
- c.  $\left\| \boldsymbol{\beta}_{k, S_k}^{(t)} - s_k \mathbf{y}_{S_k} \right\|_2 \leq C_y \|\mathbf{y}_{S_k}\|_2$ , for all  $k \in [m]$ .
- d.  $\left\| \boldsymbol{\alpha}_k^{(t)} \right\|_2 \leq \frac{C_\alpha \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1}$ , for all  $k \in [m]$ .
- e.  $\beta_{k, j}^{(t)} \leq 0$ , for all  $j \in [n]$  with  $k \neq a_j$ .

Consequently, the activation pattern of each neuron remains unchanged from iteration  $t$  to  $t + 1$ . In the above, we define  $S_k := \{i \in [n] : a_i = k\}$ , and for any vector  $\mathbf{v} \in \mathbb{R}^n$ , we use  $\mathbf{v}_{S_k}$  to denote the subvector of entries indexed by  $S_k$ .

**Proof** (Lemma 23) We now verify that these conditions are preserved from iteration  $t$  to  $t + 1$ .

Part (a): By conditions (a), (c) and (e) at iteration  $t$ , we have

$$\|h_{\Theta^{(t)}}(\mathbf{X}) - \mathbf{y}\|_2^2 = \left\| \sum_{k=1}^m s_k \sigma(\boldsymbol{\beta}_k^{(t)}) - \mathbf{y} \right\|_2^2 = \sum_{k=1}^m \left\| s_k \left( \boldsymbol{\beta}_{k, S_k}^{(t)} - s_k \mathbf{y}_{S_k} \right) \right\|_2^2 \leq C_y^2 \|\mathbf{y}\|_2^2,$$

where the last inequality uses the fact that the sets  $\{S_k\}_{k=1}^m$  are disjoint. Also, we have  $s_{a_i} h_{\Theta^{(t)}}(\mathbf{x}_i) = \beta_{a_i, i}^{(t)}$  from conditions (a) and (e). Together with condition (a), the assumptions of Lemma 10 are satisfied for all  $i \in [n]$ . Consequently,  $\beta_{a_i, i}^{(t+1)} > 0$  for all  $i \in [n]$ .

Part (b): According to the dual gradient update in Equation (51b), and using condition (e) at iteration  $t$ , we have:

$$\alpha_{k, j}^{(t+1)} = \alpha_{k, j}^{(t)} \quad \text{for all } j \in [n] \text{ with } k \neq a_j.$$

Therefore, condition (b) continues to hold at iteration  $t + 1$ .

Part (c): By conditions (a) and (e), the gradient update at iteration  $t$  for  $\beta_k^{(t)}$  depends only on the examples in the subset  $S_k$ . Hence, the gradient update for an individual neuron is equivalent to a linear regression gradient descent. As similarly argued in the proof of Lemma 2, since the step size satisfies  $\eta \leq \frac{1}{C_g \|\lambda\|_1}$ , the linear regression squared loss is monotonically nonincreasing, and by condition (c) at iteration  $t$ , we obtain

$$\left\| \beta_{k,S_k}^{(t+1)} - s_k \mathbf{y}_{S_k} \right\|_2 \leq \left\| \beta_{k,S_k}^{(t)} - s_k \mathbf{y}_{S_k} \right\|_2 \leq C_y \|\mathbf{y}_{S_k}\|_2.$$

Therefore, condition (c) holds at iteration  $t + 1$ .

Part (d): Following the same argument as in Part (d) of Lemma 17, using conditions (b) and (c) at iteration  $t + 1$ , together with the eigenvalue bounds from Lemma 12, we can establish that

$$\left\| \alpha_k^{(t+1)} \right\|_2 \leq \frac{C_\alpha \sqrt{ny_{\max}}}{\|\lambda\|_1} \text{ for all } k \in [m],$$

with probability at least  $1 - 2e^{-n/C_g}$ . Thus, condition (d) holds at iteration  $t + 1$ .

Part (e): By Lemma 11, since conditions (b) and (d) hold at iteration  $t + 1$ , we conclude that  $\beta_{k,j}^{(t+1)} \leq 0$  for all  $j \in [n]$  with  $k \neq a_j$ . Thus, condition (e) holds at iteration  $t + 1$ . ■

Equipped with Lemma 23, we are ready to prove Theorem 21.

**Proof** (Theorem 21) The proof follows a similar structure to that of Theorem 4 for single ReLU models, but now we must track the dynamics of all the neurons  $\{\mathbf{w}_k\}_{k=1}^m$  simultaneously. Equipped with sufficient conditions under which the activation patterns are preserved in Lemma 23, we verify these conditions hold after the first gradient step, and use induction to characterize the full gradient descent dynamics.

We first verify that the iterate at  $t = 1$  satisfies all the sufficient conditions. With the initialization

$$\mathbf{w}_k^{(0)} = \mathbf{X}^\top \left( \mathbf{X} \mathbf{X}^\top \right)^{-1} \left( \frac{1}{C_g} \mathbf{A}_k \mathbf{y} + \epsilon_k \right),$$

we have  $\beta_k^{(0)} = \frac{1}{C_g} \mathbf{A}_k \mathbf{y} + \epsilon_k$ . Recalling the definition of  $\mathbf{A}_k$  in Assumption 4, we have

$$\beta_{k,i}^{(0)} = \begin{cases} \epsilon_{a_i,i}, & \text{if } a_i = k, \text{ or } s_k \cdot y_i < 0 \\ -\frac{|y_i|}{C_g} + \epsilon_{k,i}, & \text{otherwise} \end{cases}, \quad (54)$$

for all  $k \in [m]$  and  $i \in [n]$ . Since the theorem assumption ensures  $\epsilon_{k,i} \leq \frac{1}{C_\alpha m} y_{\min}$  and  $C_\alpha \gtrsim C_g^2$ , we have  $-\frac{|y_i|}{C_g} + \epsilon_{k,i} < 0$ . Therefore, we obtain

$$h_{\Theta^{(0)}}(\mathbf{x}_i) = \sum_{k=1}^m s_k \sigma \left( \beta_{k,i}^{(0)} \right) = s_{a_i} \epsilon_{a_i,i} - s_{a_i} \sum_{k: s_k \cdot y_i < 0} \epsilon_{k,i}, \quad (55)$$

for all  $i \in [n]$ . Therefore, using the primal gradient update in Equation (51a), we obtain

$$\begin{aligned} \beta_k^{(1)} &= \beta_k^{(0)} - \eta s_k \mathbf{X} \mathbf{X}^\top \mathbf{D}(\beta_k^{(0)})(h_{\Theta^{(0)}}(\mathbf{X}) - \mathbf{y}) \\ &= \mathbf{X} \mathbf{X}^\top \left[ \underbrace{\eta \left( s_k \mathbf{D}(\beta_k^{(0)}) (\mathbf{y} - h_{\Theta^{(0)}}(\mathbf{X})) + \frac{1}{\eta} (\mathbf{X} \mathbf{X}^\top)^{-1} \beta_k^{(0)} \right)}_{=: \alpha_k^{(1)}} \right], \end{aligned} \quad (56)$$

according to the primal-dual formulation  $\beta_k^{(1)} = \mathbf{X} \mathbf{X}^\top \alpha_k^{(1)}$  in Equation (3). In the below, we show that at iteration  $t = 1$ , the variables  $\beta_k^{(1)}$  and  $\alpha_k^{(1)}$  satisfy all the conditions in Lemma 23.

Part (a): For all  $i \in [n]$ , we show that  $\beta_{a_i,i}^{(1)} > 0$  by applying Lemma 10. According to Equation (54) and Equation (55), we have  $\beta_{a_i,i}^{(0)} = \epsilon_{a_i,i} > 0$  and  $s_{a_i} \cdot h_{\Theta^{(0)}}(\mathbf{x}_i) = \epsilon_{a_i,i} - \sum_{k: s_k \cdot y_i < 0} \epsilon_{k,i} \leq \beta_{a_i,i}^{(0)}$ . Moreover, we have

$$\begin{aligned} \|h_{\Theta^{(0)}}(\mathbf{X}) - \mathbf{y}\|_2 &\leq \|h_{\Theta^{(0)}}(\mathbf{X})\|_2 + \|\mathbf{y}\|_2 \leq \sum_{k=1}^m \|\epsilon_k\|_2 + \|\mathbf{y}\|_2 \leq \frac{\sqrt{n}}{C_\alpha} y_{\min} + \|\mathbf{y}\|_2 \\ &\leq C_y \|\mathbf{y}\|_2, \end{aligned}$$

with  $C_y \geq 1 + \frac{1}{C_\alpha}$ . All the conditions of Lemma 10 are satisfied, and therefore,  $\beta_{a_i,i}^{(1)} > 0$  for all  $i \in [n]$ .

Part (b): For all  $j \in [n]$  with  $k \neq a_j$ , we verify that  $\alpha_{k,j}^{(1)}$  satisfies the required upper and lower bounds. We need to discuss two cases: 1)  $\beta_{k,j}^{(0)} = \epsilon_{k,j} > 0$  with  $s_k \cdot y_j < 0$ , and 2)  $\beta_{k,j}^{(0)} = -\frac{|y_j|}{C_g} + \epsilon_{k,j} < 0$  with  $s_k \cdot y_j > 0$ .

**Case 1):** For  $\beta_{k,j}^{(0)} = \epsilon_{k,j} > 0$  with  $s_k \cdot y_j < 0$ , we work from Equation (56) to get

$$\begin{aligned} \alpha_{k,j}^{(1)} &= \eta e_j^\top \left( s_k \mathbf{D}(\beta_k^{(0)}) (\mathbf{y} - h_{\Theta^{(0)}}(\mathbf{X})) + \frac{1}{\eta} (\mathbf{X} \mathbf{X}^\top)^{-1} \beta_k^{(0)} \right) \\ &\stackrel{(i)}{=} \eta \left( -|y_j| - s_k \left( s_{a_j} \epsilon_{a_j,j} - s_{a_j} \sum_{k: s_k \cdot y_j < 0} \epsilon_{k,j} \right) \right. \\ &\quad \left. + \frac{1}{\eta} e_j^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \left( \frac{1}{C_g} \mathbf{A}_k \mathbf{y} + \epsilon_k \right) \right) \\ &= \eta \left( -|y_j| + \epsilon_{a_j,j} - \sum_{k: s_k \cdot y_j < 0} \epsilon_{k,j} + \frac{1}{\eta} e_j^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \left( \frac{1}{C_g} \mathbf{A}_k \mathbf{y} + \epsilon_k \right) \right) \end{aligned}$$

$$\begin{aligned}
 &= \eta \left( -|y_j| + \epsilon_{a_j,j} - \sum_{k:s_k \cdot y_j < 0} \epsilon_{k,j} \right. \\
 &\quad \left. + \frac{1}{\eta} \mathbf{e}_j^\top \left[ \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} + \left( (\mathbf{X}\mathbf{X}^\top)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right) \right] \left( \frac{1}{C_g} \mathbf{A}_k \mathbf{y} + \boldsymbol{\epsilon}_k \right) \right) \\
 &\stackrel{\text{(ii)}}{=} \eta \left( -|y_j| + \epsilon_{a_j,j} - \sum_{k:s_k \cdot y_j < 0} \epsilon_{k,j} + \frac{\epsilon_{k,j}}{\eta \|\boldsymbol{\lambda}\|_1} \right. \\
 &\quad \left. + \frac{1}{\eta} \mathbf{e}_j^\top \left( (\mathbf{X}\mathbf{X}^\top)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right) \left( \frac{1}{C_g} \mathbf{A}_k \mathbf{y} + \boldsymbol{\epsilon}_k \right) \right) \quad (57)
 \end{aligned}$$

where equality (i) substitutes  $h_{\Theta^{(0)}}(\mathbf{x}_j) = s_{a_j} \epsilon_{a_j,j} - s_{a_j} \sum_{k:s_k \cdot y_j < 0} \epsilon_{k,j}$  from Equation (54) and  $\beta_k^{(0)} = \frac{1}{C_g} \mathbf{A}_k \mathbf{y} + \boldsymbol{\epsilon}_k$ , and equality (ii) applies  $(\mathbf{A}_k)_{jj} = 0$  for  $k \neq a_j$  with  $s_k \cdot y_j < 0$ . For the upper bound, we have

$$\alpha_{k,j}^{(1)} \leq \eta \left( -|y_j| + \epsilon_{a_j,j} + \frac{\epsilon_{k,j}}{\eta \|\boldsymbol{\lambda}\|_1} + \frac{1}{\eta} \left\| \left( \mathbf{X}\mathbf{X}^\top \right)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right\|_2 \left\| \frac{1}{C_g} \mathbf{A}_k \mathbf{y} + \boldsymbol{\epsilon}_k \right\|_2 \right), \quad (58)$$

by dropping negative terms  $-\sum_{k:s_k \cdot y_j < 0} \epsilon_{k,j}$ . Next, by applying the upper bound in Corollary 14 and the upper bounds for  $\|\mathbf{y}\|_2$  and  $\|\boldsymbol{\epsilon}_k\|_2$ , we have

$$\begin{aligned}
 \alpha_{k,j}^{(1)} &\leq \eta \left( -|y_j| + \epsilon_{a_j,j} + \frac{\epsilon_{k,j}}{\eta \|\boldsymbol{\lambda}\|_1} \right. \\
 &\quad \left. + \frac{C_g}{\eta \|\boldsymbol{\lambda}\|_1} C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \left( \frac{\sqrt{n} y_{\max}}{C_g} + \frac{\sqrt{n}}{C_\alpha m} y_{\min} \right) \right) \\
 &\leq \eta \left( -|y_j| + \epsilon_{a_j,j} + \frac{\epsilon_{k,j}}{\eta \|\boldsymbol{\lambda}\|_1} + \frac{C_g}{\eta \|\boldsymbol{\lambda}\|_1} C \cdot \frac{y_{\min}}{C_0 y_{\max}} \left( \frac{y_{\max}}{C_g} + \frac{y_{\min}}{C_\alpha m} \right) \right),
 \end{aligned}$$

where the second inequality substitutes  $d_2 \geq C_0^2 \frac{n^2 y_{\max}^2}{y_{\min}^2}$  and  $d_\infty \geq C_0 \frac{n^{1.5} y_{\max}}{y_{\min}}$  in Assumption 2. Finally, by using the step size assumption  $\frac{1}{\eta} \leq CC_g \|\boldsymbol{\lambda}\|_1$  and the theorem assumption  $\epsilon_{k,j} \leq \frac{1}{C_\alpha m} y_{\min}$ , we have

$$\begin{aligned}
 \alpha_{k,j}^{(1)} &\leq \frac{1}{CC_g \|\boldsymbol{\lambda}\|_1} \left( -y_{\min} + \frac{y_{\min}}{C_\alpha m} + \frac{CC_g y_{\min}}{C_\alpha m} + \frac{2C^2 C_g^2}{C_0} y_{\min} \right) \\
 &\leq -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1},
 \end{aligned}$$

with  $C_0 \gtrsim C_\alpha^2$  and  $C_\alpha \gtrsim \max\{C_g^2, C_y C_g\}$ . For the lower bound, starting from Equation (57), we have

$$\begin{aligned}
 \alpha_{k,j}^{(1)} &\geq \eta \left( -|y_j| - \sum_{k:s_k \cdot y_j < 0} \epsilon_{k,j} - \frac{1}{\eta} \left\| \left( \mathbf{X} \mathbf{X}^\top \right)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right\|_2 \left\| \frac{1}{C_g} \mathbf{A}_k \mathbf{y} + \boldsymbol{\epsilon}_k \right\|_2 \right) \\
 &\stackrel{(i)}{\geq} \eta \left( -|y_j| - \sum_{k:s_k \cdot y_j < 0} \epsilon_{k,j} \right. \\
 &\quad \left. - \frac{C_g}{\eta \|\boldsymbol{\lambda}\|_1} C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \left( \frac{\sqrt{n} y_{\max}}{C_g} + \frac{\sqrt{n}}{C_\alpha m} y_{\min} \right) \right) \\
 &\stackrel{(ii)}{\geq} \frac{1}{C_g \|\boldsymbol{\lambda}\|_1} \left( -y_{\max} - \frac{y_{\min}}{C_\alpha} - C^2 C_g^2 \cdot \frac{y_{\min}}{C_0 y_{\max}} \left( \frac{y_{\max}}{C_g} + \frac{y_{\min}}{C_\alpha m} \right) \right) \\
 &\stackrel{(iii)}{\geq} -\frac{3y_{\max}}{C_g \|\boldsymbol{\lambda}\|_1},
 \end{aligned}$$

where inequality (i) applies the upper bound in Corollary 14 and the upper bounds for  $\|\mathbf{y}\|_2$  and  $\|\boldsymbol{\epsilon}_k\|_2$ , inequality (ii) applies  $d_2 \geq C_0^2 \frac{n^2 y_{\max}^2}{y_{\min}}$  and  $d_\infty \geq C_0 \frac{n^{1.5} y_{\max}}{y_{\min}}$  in Assumption 2, and inequality (iii) follows by the constant relationship that  $C_0 \gtrsim C_\alpha^2$  and  $C_\alpha \gtrsim \max\{C_g^2, C_y C_g\}$ . Thus, for  $\beta_{k,j}^{(0)} = \epsilon_{k,j} > 0$  with  $s_k \cdot y_j < 0$ ,  $\alpha_{k,j}^{(1)}$  satisfies both the required upper and lower bounds.

**Case 2):** For  $\beta_{k,j}^{(0)} = -\frac{|y_j|}{C_g} + \epsilon_{k,j} < 0$  with  $s_k \cdot y_j > 0$ , we work from Equation (56) to get

$$\begin{aligned}
 \alpha_{k,j}^{(1)} &= \eta \mathbf{e}_j^\top \left( s_k \mathbf{D}(\boldsymbol{\beta}_k^{(0)}) (\mathbf{y} - h_{\Theta^{(0)}}(\mathbf{X})) + \frac{1}{\eta} \left( \mathbf{X} \mathbf{X}^\top \right)^{-1} \boldsymbol{\beta}_k^{(0)} \right) \\
 &= \mathbf{e}_j^\top \left( \mathbf{X} \mathbf{X}^\top \right)^{-1} \left( \frac{1}{C_g} \mathbf{A}_k \mathbf{y} + \boldsymbol{\epsilon}_k \right),
 \end{aligned}$$

where we substitute  $\boldsymbol{\beta}_k^{(0)} = \frac{1}{C_g} \mathbf{A}_k \mathbf{y} + \boldsymbol{\epsilon}_k$  and  $\beta_{k,j}^{(0)} = -\frac{1}{C_g} |y_j| + \epsilon_{k,j} < 0$ , and this eliminates the first term, since  $D_{jj} = 0$ . Then,  $\alpha_{k,j}^{(1)}$  can further be written as

$$\begin{aligned}
 \alpha_{k,j}^{(1)} &= \mathbf{e}_j^\top \left[ \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} + \left( \left( \mathbf{X} \mathbf{X}^\top \right)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right) \right] \left( \frac{1}{C_g} \mathbf{A}_k \mathbf{y} + \boldsymbol{\epsilon}_k \right) \\
 &= \frac{1}{\|\boldsymbol{\lambda}\|_1} \left( -\frac{|y_j|}{C_g} + \epsilon_{k,j} \right) + \mathbf{e}_j^\top \left( \left( \mathbf{X} \mathbf{X}^\top \right)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right) \left( \frac{1}{C_g} \mathbf{A}_k \mathbf{y} + \boldsymbol{\epsilon}_k \right).
 \end{aligned} \tag{59}$$

For the upper bound, we have

$$\begin{aligned}
 \alpha_{k,j}^{(1)} &\leq \frac{1}{\|\boldsymbol{\lambda}\|_1} \left( -\frac{1}{C_g} |y_j| + \epsilon_{k,j} \right) + \left\| \left( \mathbf{X} \mathbf{X}^\top \right)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right\|_2 \left\| \frac{1}{C_g} \mathbf{A}_k \mathbf{y} + \boldsymbol{\epsilon}_k \right\|_2 \\
 &\stackrel{(i)}{\leq} \frac{1}{\|\boldsymbol{\lambda}\|_1} \left( -\frac{|y_j|}{C_g} + \epsilon_{k,j} + C_g C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \left( \frac{\sqrt{n} y_{\max}}{C_g} + \frac{\sqrt{n}}{C_\alpha} y_{\min} \right) \right) \\
 &\stackrel{(ii)}{\leq} \frac{1}{\|\boldsymbol{\lambda}\|_1} \left( -\frac{y_{\min}}{C_g} + \frac{y_{\min}}{C_\alpha m} + C_g C \cdot \frac{y_{\min}}{C_0 y_{\max}} \left( \frac{y_{\max}}{C_g} + \frac{1}{C_\alpha} y_{\min} \right) \right) \\
 &\leq -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1},
 \end{aligned}$$

where inequality (i) applies the upper bound in Corollary 14 and the upper bounds for  $\|\mathbf{y}\|_2$  and  $\|\boldsymbol{\epsilon}_k\|_2$ , inequalities (ii) substitutes  $d_2 \geq C_0^2 \frac{n^2 y_{\max}^2}{y_{\min}^2}$  and  $d_\infty \geq C_0 \frac{n^{1.5} y_{\max}}{y_{\min}}$  in Assumption 2. The last inequality follows by  $C_0 \gtrsim C_\alpha^2$  and  $C_\alpha \gtrsim \max\{C_g^2, C_y C_g\}$ . For the lower bound, we work from Equation (59) to get

$$\begin{aligned}
 \alpha_{k,j}^{(1)} &\geq \frac{1}{\|\boldsymbol{\lambda}\|_1} \left( -\frac{|y_j|}{C_g} \right) - \left\| \left( \mathbf{X} \mathbf{X}^\top \right)^{-1} - \frac{1}{\|\boldsymbol{\lambda}\|_1} \mathbf{I} \right\|_2 \left\| \frac{1}{C_g} \mathbf{A}_k \mathbf{y} + \boldsymbol{\epsilon}_k \right\|_2 \\
 &\geq \frac{1}{\|\boldsymbol{\lambda}\|_1} \left( -\frac{|y_j|}{C_g} - C_g C \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \left( \frac{\sqrt{n} y_{\max}}{C_g} + \frac{\sqrt{n}}{C_\alpha} y_{\min} \right) \right) \\
 &\geq \frac{1}{\|\boldsymbol{\lambda}\|_1} \left( -\frac{y_{\max}}{C_g} - C_g C \cdot \frac{y_{\min}}{C_0 y_{\max}} \left( \frac{y_{\max}}{C_g} + \frac{y_{\min}}{C_\alpha} \right) \right) \\
 &\geq -\frac{3y_{\max}}{C_g \|\boldsymbol{\lambda}\|_1},
 \end{aligned}$$

by the same argument. Thus, for  $\beta_{k,j}^{(0)} = -\frac{|y_j|}{C_g} + \epsilon_{k,j} < 0$  with  $s_k \cdot y_j > 0$ ,  $\alpha_{k,j}^{(1)}$  satisfies both the required upper and lower bounds. This completes the proof of this part.

Part (c): We verify that the primal variables  $\beta_{k,S_k}^{(1)}$  corresponding to active examples minus  $\mathbf{y}_{S_k}$  satisfy the norm bound. Specifically, we show that  $\|\beta_{k,S_k}^{(1)} - s_k \mathbf{y}_{S_k}\|_2^2 \leq C_y^2 \|\mathbf{y}\|_2^2$ . According to Equation (56), we have

$$\begin{aligned}
 \|\beta_{k,S_k}^{(1)} - s_k \mathbf{y}_{S_k}\|_2^2 &= \sum_{i:a_i=k} \left( \beta_{k,i}^{(1)} - s_k y_i \right)^2 \\
 &= \sum_{i:a_i=k} \left( \beta_{k,i}^{(0)} - \eta s_k \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top \mathbf{D}(\beta_k^{(0)}) (h_{\Theta^{(0)}}(\mathbf{X}) - \mathbf{y}) - s_k y_i \right)^2 \\
 &= \sum_{i:a_i=k} \left( \underbrace{\epsilon_{a_i,i} - \eta s_{a_i} \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top \mathbf{D}(\beta_{a_i}^{(0)}) (h_{\Theta^{(0)}}(\mathbf{X}) - \mathbf{y}) - |y_i|}_{=: T_i} \right)^2,
 \end{aligned} \tag{60}$$

where we substitute  $\beta_{k,i}^{(0)} = \beta_{a_i,i}^{(0)} = \epsilon_{a_i,i}$  for  $k = a_i$ , and  $s_{a_i} \cdot y_i > 0$ . Next, we bound the term

$T_i := \epsilon_{a_i,i} - \eta s_{a_i} \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top \mathbf{D}(\beta_{a_i}^{(0)})(h_{\Theta^{(0)}}(\mathbf{X}) - \mathbf{y}) - |y_i|$  for all  $i \in [n]$ . We have

$$\begin{aligned}
 T_i &= \epsilon_{a_i,i} - \eta s_{a_i} \mathbf{e}_i^\top \mathbf{X} \mathbf{X}^\top \mathbf{D}(\beta_{a_i}^{(0)})(h_{\Theta^{(0)}}(\mathbf{X}) - \mathbf{y}) - |y_i| \\
 &= (\epsilon_{a_i,i} - |y_i|) - \eta s_{a_i} \mathbf{e}_i^\top \left[ \|\boldsymbol{\lambda}\|_1 \mathbf{I} + (\mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}) \right] \mathbf{D}(\beta_{a_i}^{(0)})(h_{\Theta^{(0)}}(\mathbf{X}) - \mathbf{y}) \\
 &= (\epsilon_{a_i,i} - |y_i|) - \eta s_{a_i} \|\boldsymbol{\lambda}\|_1 \left( s_{a_i} \epsilon_{a_i,i} - s_{a_i} \sum_{k:s_k \cdot y_i < 0} \epsilon_{k,i} - y_i \right) \\
 &\quad - \eta s_k \mathbf{e}_i^\top (\mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}) \mathbf{D}(\beta_k^{(0)})(h_{\Theta^{(0)}}(\mathbf{X}) - \mathbf{y}) \\
 &= (1 - \eta \|\boldsymbol{\lambda}\|_1) \epsilon_{a_i,i} - (1 - \eta \|\boldsymbol{\lambda}\|_1) |y_i| + \eta \|\boldsymbol{\lambda}\|_1 \sum_{k:s_k \cdot y_i < 0} \epsilon_{k,i} \\
 &\quad - \eta s_k \mathbf{e}_i^\top (\mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}) \mathbf{D}(\beta_k^{(0)})(h_{\Theta^{(0)}}(\mathbf{X}) - \mathbf{y}),
 \end{aligned}$$

by applying  $h_{\Theta^{(0)}}(\mathbf{x}_i) = s_{a_i} \epsilon_{a_i,i} - s_{a_i} \sum_{k:s_k \cdot y_i < 0} \epsilon_{k,i}$  from Equation (55). Since the step size assumption guarantees that  $\frac{1}{C_g \|\boldsymbol{\lambda}\|_1} \leq \eta \leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1}$ , and  $\epsilon_{k,i} \leq \frac{1}{C_\alpha m} y_{\min}$ , we have

$$\begin{aligned}
 (1 - \eta \|\boldsymbol{\lambda}\|_1) \epsilon_{a_i,i} - (1 - \eta \|\boldsymbol{\lambda}\|_1) |y_i| + \eta \|\boldsymbol{\lambda}\|_1 \sum_{k:s_k \cdot y_i < 0} \epsilon_{k,i} \\
 \leq \epsilon_{a_i,i} - (1 - \eta \|\boldsymbol{\lambda}\|_1) |y_i| + \eta \|\boldsymbol{\lambda}\|_1 \sum_{k:s_k \cdot y_i < 0} \epsilon_{k,i} \\
 \leq \left( \frac{1}{m} + \frac{1}{C_g} \right) \frac{1}{C_\alpha} y_{\min} - \left( 1 - \frac{1}{C_g} \right) y_{\min} \\
 < 0,
 \end{aligned}$$

with  $C_\alpha \gtrsim C_g^2$ . Hence, in order to upper bound  $T_i^2$ , it suffices to find the lower bound for  $T_i$ . We have

$$\begin{aligned}
 T_i &= (1 - \eta \|\boldsymbol{\lambda}\|_1) \epsilon_{a_i,i} - (1 - \eta \|\boldsymbol{\lambda}\|_1) |y_i| + \eta \|\boldsymbol{\lambda}\|_1 \sum_{k:s_k \cdot y_i < 0} \epsilon_{k,i} \\
 &\quad - \eta s_k \mathbf{e}_i^\top (\mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}) \mathbf{D}(\beta_k^{(0)})(h_{\Theta^{(0)}}(\mathbf{X}) - \mathbf{y}) \\
 &\geq -|y_i| - \eta \left\| \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right\|_2 \|h_{\Theta^{(0)}}(\mathbf{X}) - \mathbf{y}\|_2,
 \end{aligned}$$

where the inequality drops the positive terms  $(1 - \eta \|\boldsymbol{\lambda}\|_1) \epsilon_{a_i,i}$ ,  $\eta \|\boldsymbol{\lambda}\|_1 |y_i|$ , and  $\eta \|\boldsymbol{\lambda}\|_1 \sum_{k:s_k \cdot y_i < 0} \epsilon_{k,i}$ . We again upper bound  $\left\| \mathbf{X} \mathbf{X}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I} \right\|_2$  by Corollary 14. With probability at least  $1 - 2 \exp(-n(Cc - \ln 9))$ , we have

$$\begin{aligned}
 T_i &\geq -|y_i| - \eta \cdot C \|\boldsymbol{\lambda}\|_1 \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \|h_{\Theta^{(0)}}(\mathbf{X}) - \mathbf{y}\|_2 \\
 &\geq -|y_i| - \frac{C}{C_g} \cdot \max \left( \sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty} \right) \|h_{\Theta^{(0)}}(\mathbf{X}) - \mathbf{y}\|_2,
 \end{aligned}$$

by applying  $\eta \leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1}$ . Finally, we apply the upper bound that  $\|h_{\Theta^{(0)}}(\mathbf{X})\|_2 \leq \sum_{k=1}^m \|\boldsymbol{\epsilon}_k\|_2 \leq \frac{\sqrt{n}}{C_\alpha} y_{\min}$  and  $\|\mathbf{y}\|_2 \leq \sqrt{n} y_{\max}$ , and Assumption 2 ensures that  $d_2 \geq C_0^2 \frac{n^2 y_{\max}^2}{y_{\min}^2}$  and  $d_\infty \geq C_0 \frac{n^{1.5} y_{\max}}{y_{\min}}$ . We have

$$\begin{aligned} T_i &\geq -|y_i| - \frac{C}{C_g} \max\left(\sqrt{\frac{n}{d_2}}, \frac{n}{d_\infty}\right) \left(\frac{\sqrt{n}}{C_\alpha} y_{\min} + \sqrt{n} y_{\max}\right) \\ &\geq -|y_i| - \frac{C y_{\min}}{C_g C_0 y_{\max}} \left(\frac{1}{C_\alpha} y_{\min} + y_{\max}\right) \\ &\geq -|y_i| \left(1 + \frac{2C}{C_g C_0}\right) \\ &\geq -C_y |y_i|, \end{aligned}$$

with the choice of  $C_y \geq 2$ . Substituting  $T_i^2 \leq C_y^2 y_i^2$  into Equation (60), we have

$$\left\| \boldsymbol{\beta}_{k, S_k}^{(1)} - s_k \mathbf{y}_{S_k} \right\|_2^2 \leq \sum_{i: a_i=k} C_y^2 y_i^2 = C_y^2 \|\mathbf{y}_{S_k}\|_2^2.$$

As a result, we conclude that  $\left\| \boldsymbol{\beta}_{k, S_k}^{(1)} - s_k \mathbf{y}_{S_k} \right\|_2 \leq C_y \|\mathbf{y}_{S_k}\|_2$  as required.

Part (d): We verify the norm bounds on the dual variables. By the triangle inequality, we work from Equation (56) to get

$$\begin{aligned} \left\| \boldsymbol{\alpha}_k^{(1)} \right\|_2 &= \left\| \eta \left( s_k \mathbf{D}(\boldsymbol{\beta}_k^{(0)}) (\mathbf{y} - h_{\Theta^{(0)}}(\mathbf{X})) + \frac{1}{\eta} (\mathbf{X} \mathbf{X}^\top)^{-1} \boldsymbol{\beta}_k^{(0)} \right) \right\|_2 \\ &\leq \eta \left[ \|\mathbf{y}\|_2 + \|h_{\Theta^{(0)}}(\mathbf{X})\|_2 + \frac{1}{\eta} \left\| (\mathbf{X} \mathbf{X}^\top)^{-1} \right\|_2 \left\| \boldsymbol{\beta}_k^{(0)} \right\|_2 \right] \\ &= \eta \left[ \|\mathbf{y}\|_2 + \|h_{\Theta^{(0)}}(\mathbf{X})\|_2 + \frac{1}{\eta} \left\| (\mathbf{X} \mathbf{X}^\top)^{-1} \right\|_2 \left\| \frac{1}{C_g} \mathbf{A}_k \mathbf{y} + \boldsymbol{\epsilon}_k \right\|_2 \right], \end{aligned}$$

by substituting  $\boldsymbol{\beta}_k^{(0)} = \frac{1}{C_g} \mathbf{A}_k \mathbf{y} + \boldsymbol{\epsilon}_k$ . Using  $\|\mathbf{y}\|_2 \leq \sqrt{n} y_{\max}$ ,  $\|h_{\Theta^{(0)}}(\mathbf{X})\|_2 \leq \sum_{k=1}^m \|\boldsymbol{\epsilon}_k\|_2 \leq \frac{\sqrt{n}}{C_\alpha} y_{\min}$ ,  $\left\| (\mathbf{X} \mathbf{X}^\top)^{-1} \right\|_2 \leq \frac{C_g}{\|\boldsymbol{\lambda}\|_1}$ ,  $\epsilon_{k,i} \leq \frac{1}{C_\alpha m} y_{\min}$ , and  $\frac{1}{C C_g \|\boldsymbol{\lambda}\|_1} \leq \eta \leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1}$ , we have

$$\begin{aligned} &\left\| \boldsymbol{\alpha}_k^{(1)} \right\|_2 \\ &\leq \frac{1}{C_g \|\boldsymbol{\lambda}\|_1} \left[ \sqrt{n} y_{\max} + \frac{\sqrt{n}}{C_\alpha} y_{\min} + C C_g \|\boldsymbol{\lambda}\|_1 \cdot \frac{C_g}{\|\boldsymbol{\lambda}\|_1} \cdot \left( \frac{\sqrt{n} y_{\max}}{C_g} + \frac{\sqrt{n}}{C_\alpha m} y_{\min} \right) \right] \\ &\leq \frac{1}{\|\boldsymbol{\lambda}\|_1} (3\sqrt{n} y_{\max}) \\ &\leq \frac{C_\alpha \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1}, \end{aligned}$$

with  $C_\alpha \gtrsim \max\{C_g^2, C_y C_g\}$ . Thus, condition (d) holds at  $t = 1$ .

Part (e): Since we have shown that  $\alpha_{k,j}^{(1)} \leq -\frac{y_{\min}}{C_\alpha \|\boldsymbol{\lambda}\|_1}$  and  $\|\boldsymbol{\alpha}_k^{(1)}\|_2 \leq \frac{C_\alpha \sqrt{n} y_{\max}}{\|\boldsymbol{\lambda}\|_1}$  for all  $j \in [n]$  with  $k \neq a_j$ , by Lemma 11, it follows that  $\beta_{k,j}^{(1)} \leq 0$  for all  $j \in [n]$  with  $k \neq a_j$ .

We have shown that at iteration  $t = 1$  the conditions in Lemma 23 are satisfied, and by induction, these conditions will also hold for  $t \geq 1$ . As a result,  $\mathbf{w}_k$  is trained with only predefined active examples starting from the iteration  $t = 0$ , and it is equivalent to linear regression using only active examples with initialization  $\mathbf{w}_k^{(1)} = \eta \mathbf{X}^\top \left( s_k \mathbf{D}(\boldsymbol{\beta}_k^{(0)}) (\mathbf{y} - h_{\Theta^{(0)}}(\mathbf{X})) + \frac{1}{\eta} (\mathbf{X} \mathbf{X}^\top)^{-1} \boldsymbol{\beta}_k^{(0)} \right)$ . Finally, since  $\mathbf{w}_k$  is trained on disjoint subset of examples by Assumption 4, by Lemma 2,  $\mathbf{w}_k^{(\infty)}$  satisfies

$$\mathbf{w}_k^{(\infty)} = \arg \min_{\mathbf{w} \in \{\mathbf{w} : \mathbf{X}_{S_k} \mathbf{w} = \mathbf{y}_{S_k}\}} \|\mathbf{w} - \mathbf{w}_k^{(1)}\|_2.$$

This completes the proof of Theorem 21. ■

### E.3. Proof of Theorem 22 (Implicit Bias Approximation to $\mathbf{w}^*$ )

**Proof** (Theorem 22) We restate the definition of  $\mathbf{w}^*$  in Equation (52).

$$\begin{aligned} \{\mathbf{w}_k^*\}_{k=1}^m &= \arg \min_{\{\mathbf{w}_k\}_{k=1}^m} \frac{1}{2} \sum_{k=1}^m \|\mathbf{w}_k\|_2^2 \\ \text{s.t. } \sum_{k=1}^m s_k \sigma(\mathbf{w}_k^\top \mathbf{x}_i) &= y_i, \text{ for all } i \in [n]. \end{aligned} \quad (61)$$

Recall that the gradient descent limit  $\{\mathbf{w}_k^{(\infty)}\}_{k=1}^m$  satisfies the same set of constraints: it interpolates all examples. Consequently, both  $\{\mathbf{w}_k^{(\infty)}\}_{k=1}^m$  and  $\{\mathbf{w}_k^*\}_{k=1}^m$  are feasible solutions to (61). We show that the norm difference between  $\mathbf{w}_k^{(\infty)}$  and  $\mathbf{w}_k^*$  can be upper bounded by 2 times the norm of  $\mathbf{w}_k^{(\infty)}$ .

$$\sum_{k=1}^m \|\mathbf{w}_k^{(\infty)} - \mathbf{w}_k^*\|_2^2 \leq 2 \sum_{k=1}^m \|\mathbf{w}_k^{(\infty)}\|_2^2 + 2 \sum_{k=1}^m \|\mathbf{w}_k^*\|_2^2 \leq 4 \sum_{k=1}^m \|\mathbf{w}_k^{(\infty)}\|_2^2,$$

where it follows the definition of (61). By Lemma 23, we have the upper bound for  $\|\mathbf{w}_k^{(\infty)}\|_2^2$  as

$$\|\mathbf{w}_k^{(\infty)}\|_2^2 = \boldsymbol{\alpha}_k^{(\infty)\top} \mathbf{X} \mathbf{X}^\top \boldsymbol{\alpha}_k^{(\infty)} \leq \mu_1(\mathbf{X} \mathbf{X}^\top) \|\boldsymbol{\alpha}_k^{(\infty)}\|_2^2 \leq C_g \|\boldsymbol{\lambda}\|_1 \cdot \frac{C_\alpha^2 n y_{\max}^2}{\|\boldsymbol{\lambda}\|_1^2} = \frac{C_g C_\alpha^2 n y_{\max}^2}{\|\boldsymbol{\lambda}\|_1}.$$

As a result, we have

$$\sum_{k=1}^m \|\mathbf{w}_k^{(\infty)} - \mathbf{w}_k^*\|_2^2 \leq \frac{4 C_g C_\alpha^2 m n y_{\max}^2}{\|\boldsymbol{\lambda}\|_1}.$$

■

## Appendix F. Simulations

In this section, we present exploratory visualizations of the evolution of the primal variables at iteration checkpoints in settings that violate the assumptions underlying our theoretical results.

### F.1. Moderate-Dimensional Data and Single ReLU Model

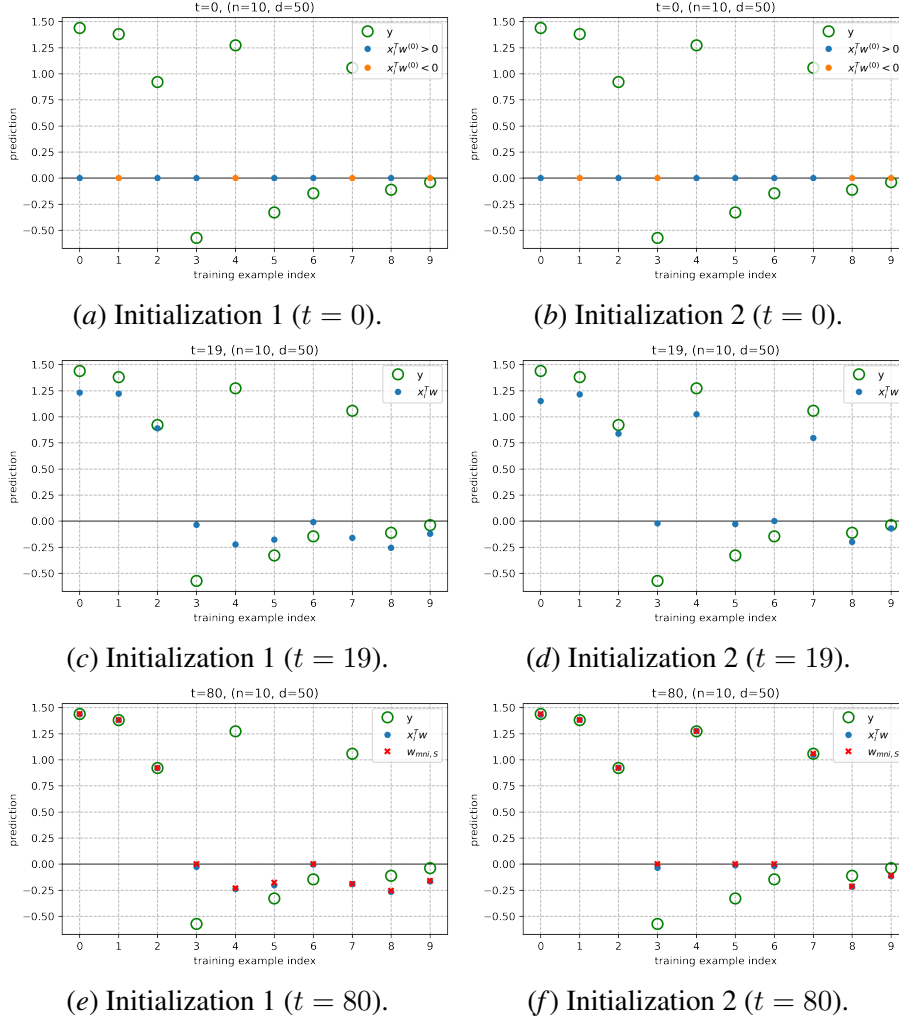


Figure 3: We illustrate the prediction dynamics of gradient descent for a single ReLU model under different random initializations when  $d$  is comparable with  $n$ . In both cases, with sufficiently small step size, the final solution converges to a linear minimum- $\ell_2$ -norm interpolator on some subset of the training examples, i.e. of the form  $\mathbf{w}_{\text{linear-MNI},S} = \mathbf{X}_S^\top (\mathbf{X}_S \mathbf{X}_S^\top)^{-1} \tilde{\mathbf{y}}_S$ , where  $\tilde{y}_{S,i} = \max\{y_i, 0\}$ . In contrast to the high-dimensional regime, *different initializations lead to different subsets  $S$* , indicating that ReLU training implicitly performs an example “selection” process, that is initialization-dependent, rather than fitting all positively-labeled samples. The experiment uses  $n = 10$ ,  $d = 50$ ,  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $y \sim \mathcal{N}(0, 1)$ ,  $\mathbf{w}^{(0)} \sim \mathcal{N}(\mathbf{0}, 2 \times 10^{-6} \mathbf{I})$ , and  $\eta = 10^{-4}$ .

## F.2. Gradient Descent Dynamics of Two ReLU Models

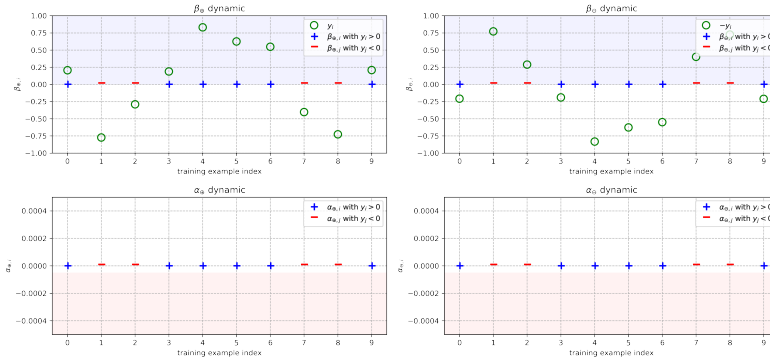
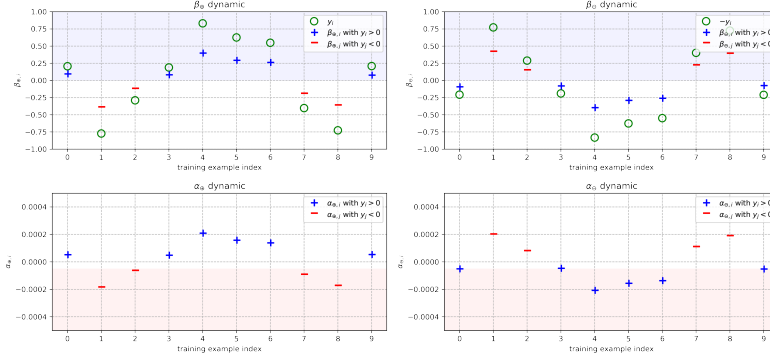
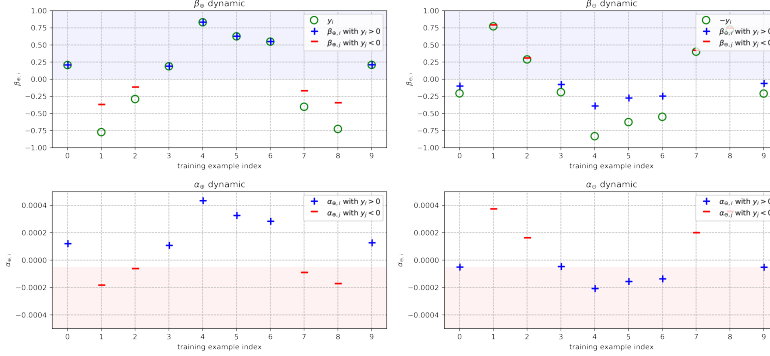
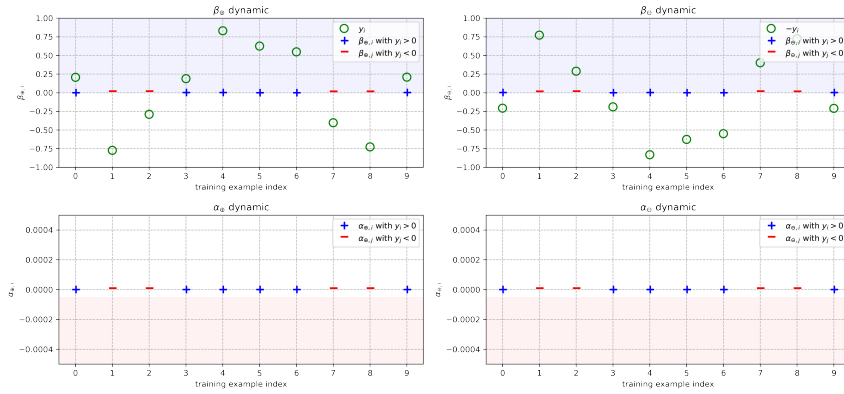
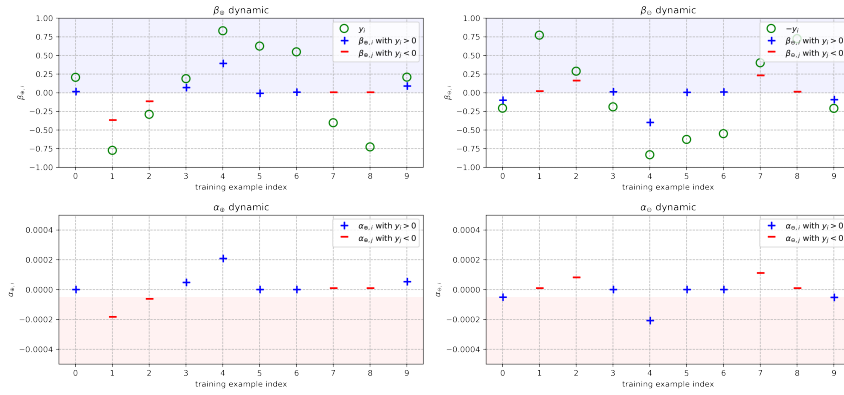
(a) Two ReLU model gradient descent dynamic ( $t = 0$ ).(b) Two ReLU model gradient descent dynamic ( $t = 1$ ).(c) Two ReLU model gradient descent dynamic ( $t = 17$ ).

Figure 4: Simulation illustrating Theorem 8. In the high-dimensional regime and under our “all-positive” initialization, after the first gradient step, examples with positive labels remain active while examples with negative labels become inactive, consistent with Lemma 19. The blue region shows primal variables that remain positive over training, whereas the red region corresponds to dual variables that are sufficiently negative and remain unchanged. As training proceeds,  $w_{\oplus}$  fits all positively labeled examples and  $w_{\ominus}$  fits all negatively labeled examples. The experiment uses  $n = 10$ ,  $d = 2000$ , features  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and labels satisfying  $|y| \sim \mathcal{U}(0.1, 1)$  with  $\text{sign}(y)$  uniformly distributed over  $\{\pm 1\}$ .

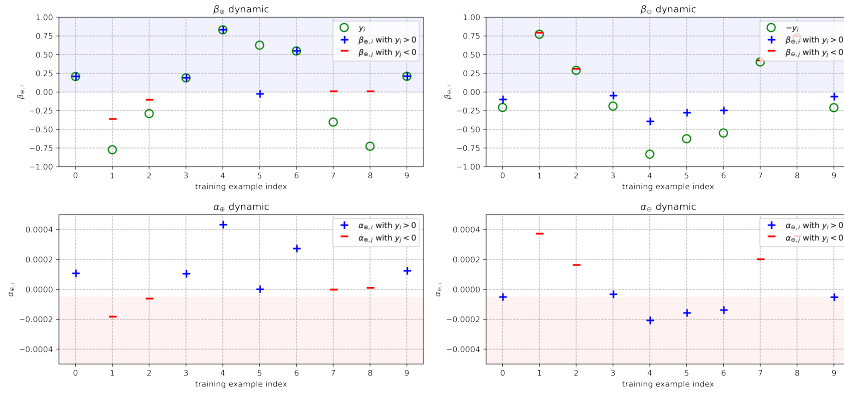
# HIGH-DIMENSIONAL IMPLICIT BIAS OF SQUARED LOSS RELU



(a) Two ReLU model gradient descent dynamic ( $t = 0$ ).



(b) Two ReLU model gradient descent dynamic ( $t = 1$ ).



(c) Two ReLU model gradient descent dynamic ( $t = 40$ ).

Figure 5: Simulation with random initialization in the high-dimensional regime, which violates our initialization assumption in Theorem 8. Under random initialization, the sufficient conditions of Lemma 19 are violated at the first gradient step. As a result, positively labeled examples do not all remain in the active (blue) regime (e.g., example no. 5), nor do negatively labeled examples consistently enter the inactive (red) regime (e.g., example no. 7). *Consequently, during training, this model fails to converge to a global minimum.* The experiment uses  $n = 10$ ,  $d = 2000$ , features  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and labels satisfying  $|y| \sim \mathcal{U}(0.1, 1)$  with  $\text{sign}(y)$  uniformly distributed over  $\{\pm 1\}$ .

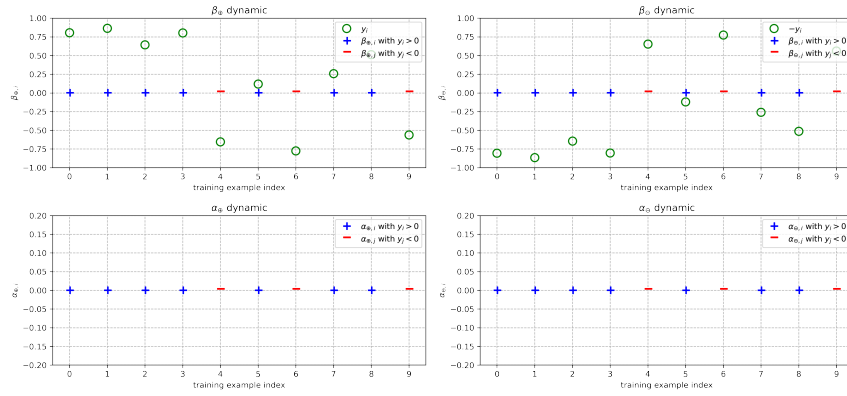
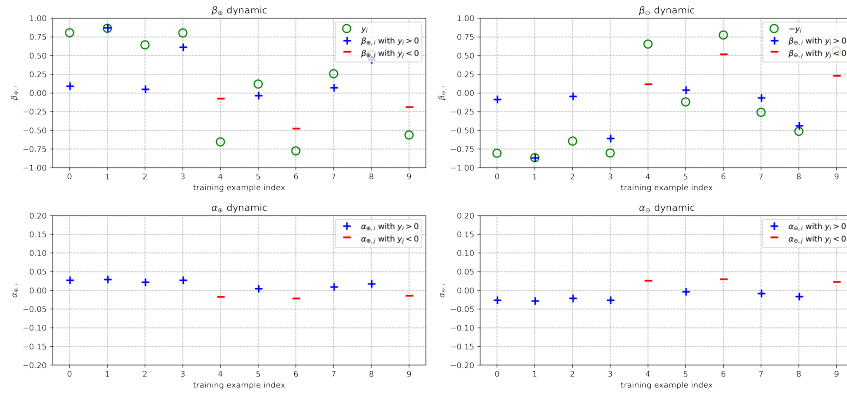
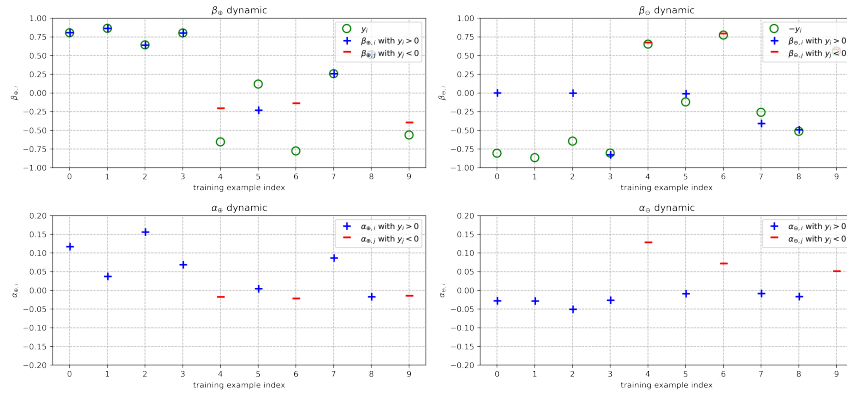
(a) Two ReLU model gradient descent dynamic ( $t = 0$ ).(b) Two ReLU model gradient descent dynamic ( $t = 1$ ).(c) Two ReLU model gradient descent dynamic ( $t = 40$ ).

Figure 6: Simulation with all-positive initialization outside the high-dimensional regime. When the data dimension is not sufficiently large, the feature vectors are no longer approximately orthogonal. As a result, the clear separation into active (blue) and inactive (red) regimes observed in Figures 4 and 5 disappears. Consequently, the gradient dynamics become highly coupled across examples and are no longer analytically tractable using our high-dimensional arguments. The experiment uses  $n = 10$ ,  $d = 15$ , features  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and labels satisfying  $|y| \sim \mathcal{U}(0.1, 1)$  with  $\text{sign}(y)$  uniformly distributed over  $\{\pm 1\}$ .

## E.3. Gradient Descent Dynamics of Multiple ReLU Models

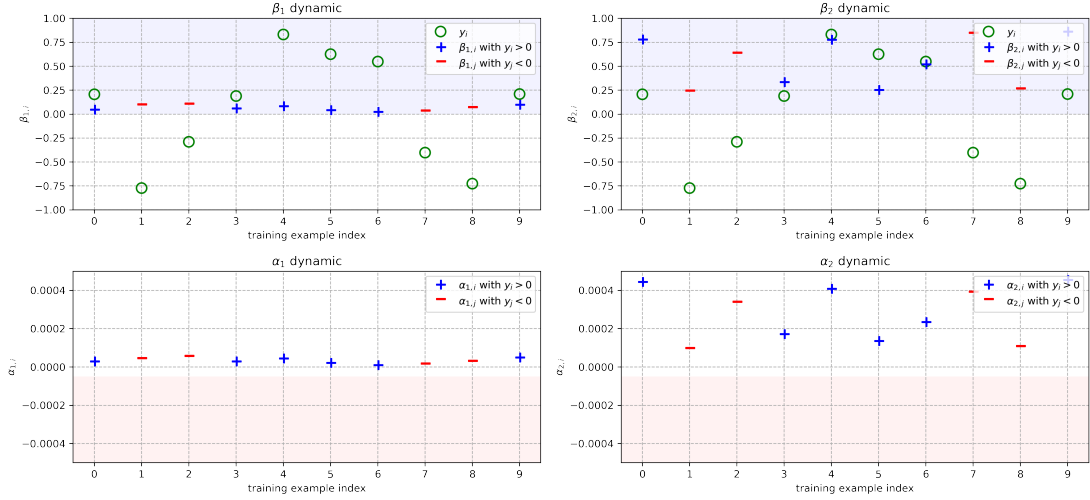
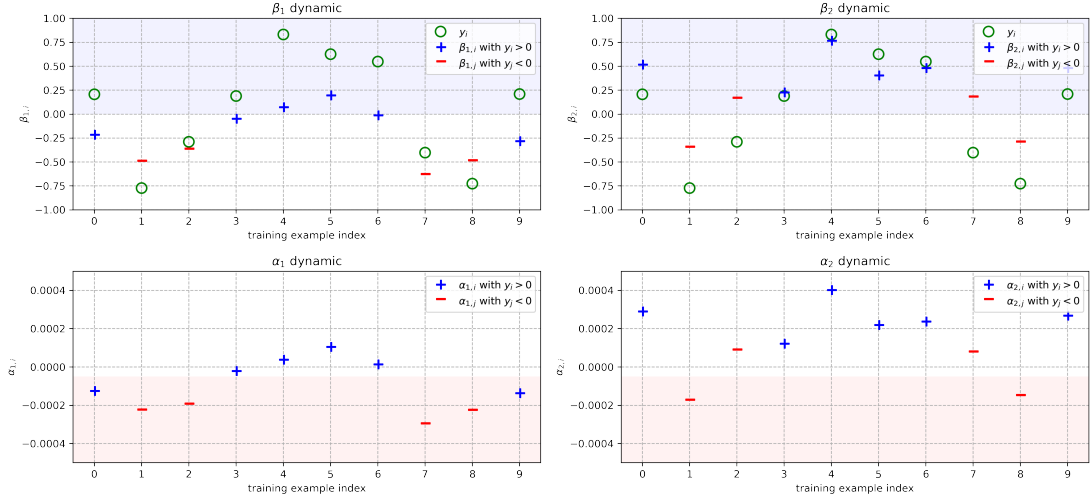

 (a) Multiple ReLU model gradient descent dynamic ( $t = 0$ ).

 (b) Multiple ReLU model gradient descent dynamic ( $t = 1$ ).

Figure 7: Failure of stable activation patterns in multiple ReLU models. We illustrate the training dynamics of a multiple ReLU model when multiple neurons share the same sign. In this setting, the sufficient conditions of Lemma 23 are violated, and positive primal variables do not necessarily remain in the active (blue) regime throughout training (e.g. training example no. 0). As a result, the activation pattern becomes unstable, and the resulting primal dynamics are no longer tractable. The experiment uses  $n = 10$ ,  $d = 2000$ ,  $m = 4$ , with neuron signs  $s_1 = s_2 = 1$  and  $s_3 = s_4 = -1$ , features  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and labels satisfying  $|y| \sim \mathcal{U}(0.1, 1)$  with  $\text{sign}(y)$  uniformly distributed over  $\{\pm 1\}$ .