
Implicit Bias of Gradient Descent for Two-layer ReLU and Leaky ReLU Networks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The implicit bias towards solutions with favorable properties is believed to a key
2 reason why neural networks trained by gradient-based optimization can generalize
3 well. While the implicit bias of gradient flow has been widely studied for homoge-
4 neous neural networks (including ReLU and leaky ReLU networks), the implicit
5 bias of gradient descent is currently only understood for smooth neural networks.
6 Therefore, implicit bias in non-smooth neural networks trained by gradient de-
7 scent remains an open question. In this paper, we aim to answer this question
8 by studying the implicit bias of gradient descent for training two-layer fully con-
9 nected (leaky) ReLU neural networks. We showed that when the training data are
10 nearly-orthogonal, for leaky ReLU activation function, gradient descent will find
11 a network with a stable rank that converges to 1, whereas for ReLU activation
12 function, gradient descent will find a neural network with a stable rank that is upper
13 bounded by a constant. Additionally, we conducted experiments to validate these
14 theoretical findings.

15 1 Introduction

16 Neural networks have achieved remarkable success in a variety of applications, such as image and
17 speech recognition, natural language processing, and many others. Recent studies have revealed that
18 the effectiveness of neural networks is attributed to their implicit bias towards particular solutions
19 which enjoy favorable properties. Understanding how this bias is affected by factors such as network
20 architecture, optimization algorithms and data used for training, has become an active research area
21 in the field of deep learning theory.

22 The literature on the implicit bias in neural networks has expanded rapidly in recent years (Vardi,
23 2022), with numerous studies shedding light on the implicit bias of gradient flow (GF) with a wide
24 range of neural network architecture, including deep linear networks (Ji and Telgarsky, 2018, 2020;
25 Gunasekar et al., 2018), homogeneous networks (Lyu and Li, 2019; Vardi et al., 2022a) and more
26 specific cases (Chizat and Bach, 2020; Lyu et al., 2021; Frei et al., 2022b; Safran et al., 2022). The
27 implicit bias of gradient descent (GD), on the other hand, is better understood for linear predictors
28 (Soudry et al., 2018) and smoothed neural networks (Lyu and Li, 2019; Frei et al., 2022b). Therefore,
29 an open question still remains:

30 *What is the implicit bias of leaky ReLU and ReLU networks trained by gradient descent?*

31 In this paper, we will answer this question by investigating gradient descent for both two-layer leaky
32 ReLU and ReLU neural networks on specific training data, where $\{\mathbf{x}_i\}_{i=1}^n$ are nearly-orthogonal
33 (Frei et al., 2022b), that is, $\|\mathbf{x}_i\|_2^2 \geq Cn \max_{k \neq i} |\langle \mathbf{x}_i, \mathbf{x}_k \rangle|$ with a constant C .

34 Our main results are summarized as follows:

Submitted to 37th Conference on Neural Information Processing Systems (NeurIPS 2023). Do not distribute.

- For two-layer leaky ReLU networks trained by GD, we demonstrate that the neuron activation pattern reaches a stable state beyond a specific time threshold and provide rigorous proof of the convergence of the stable rank of the weight matrix to 1, matching the results of Frei et al. (2022b) regarding gradient flow. Additionally, we conduct a convergence rate analysis of the training loss related to the extent of orthogonality in the training data. Notably, when the training data exhibit mutual orthogonality, we establish a convergence rate of $O(t^{-1})$, where t is the number of gradient descent iterations. This convergence rate improves upon the $O(t^{-1/2})$ rate provided in Frei et al. (2022b) for the case of a two-layer *smoothed* leaky ReLU network trained using gradient descent.
- For two-layer ReLU networks trained by GD, we proved that the stable rank of weight matrix can be upper bounded by a constant. Moreover, we present an illustrative example using completely orthogonal training data, showing that the stable rank of the weight matrix converges to a value approximately equal to 2. This finding suggests that ReLU networks possess superior learning ability compared to leaky ReLU networks. To the best of our knowledge, this is the first implicit bias result for two-layer ReLU networks trained by gradient descent beyond the Karush–Kuhn–Tucker (KKT) point

2 Related Work

Implicit bias in neural networks. Recent years have witnessed significant progress on implicit bias in neural networks trained by gradient flow (GF). Lyu and Li (2019) and Ji and Telgarsky (2020) demonstrated that homogeneous neural networks trained with exponentially-tailed classification losses converge in direction to the KKT point of a maximum-margin problem. Lyu et al. (2021) studied the implicit bias in two-layer leaky ReLU networks trained on linearly separable and symmetric data, showing that GF converges to a linear classifier maximizing the ℓ_2 margin. Frei et al. (2022b) showed that two-layer leaky ReLU networks trained by GF on nearly-orthogonal data produce a ℓ_2 -max-margin solution with a linear decision boundary and rank at most two. Other works studying the implicit bias of classification using GF in nonlinear two-layer networks include Chizat and Bach (2020); Phuong and Lampert (2021); Sarussi et al. (2021); Safran et al. (2022); Vardi et al. (2022a,b); Timor et al. (2023). Although implicit bias in neural networks trained by GF has been extensively studied, research on implicit bias in networks trained by gradient descent (GD) remains limited. Lyu and Li (2019) examined smoothed homogeneous neural network trained by GD with exponentially-tailed losses and proved a convergence to KKT points of a max-margin problem. Frei et al. (2022b) studied two-layer smoothed leaky ReLU trained by GD and revealed the implicit bias towards low-rank networks. Other works studying implicit bias towards rank minimization include Ji and Telgarsky (2018, 2020); Timor et al. (2023); Arora et al. (2019); Razin and Cohen (2020); Li et al. (2021). Lastly, Vardi (2022) provided a comprehensive literature survey on implicit bias.

Benign overfitting and double descent in neural networks. A parallel line of research aims to understand the benign overfitting phenomenon (Bartlett et al., 2020) of neural networks by considering a variety of models. For example, Allen-Zhu and Li (2020); Jelassi and Li (2022); Shen et al. (2022); Cao et al. (2022); Kou et al. (2023) studied the generalization performance of two-layer convolutional networks on patch-based data models. Several other papers studied high-dimensional mixture models (Chatterji and Long, 2021; Wang and Thrampoulidis, 2022; Cao et al., 2021; Frei et al., 2022a). Another thread of work Belkin et al. (2020); Hastie et al. (2022); Wu and Xu (2020); Mei and Montanari (2019); Liao et al. (2020) focuses on understanding the double descent phenomenon first empirically observed by Belkin et al. (2019).

3 Preliminaries

In this section, we introduce the notation, fully connected neural networks, the gradient descent-based training algorithm, and a signal-noise decomposition technique.

Notation. We use lower case letters, lower case bold face letters, and upper case bold face letters to denote scalars, vectors, and matrices respectively. For a vector $\mathbf{v} = (v_1, \dots, v_d)^\top$, we denote by $\|\mathbf{v}\|_2 := (\sum_{j=1}^d v_j^2)^{1/2}$ its ℓ_2 norm. For a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, we use $\|\mathbf{A}\|_F$ to denote its Frobenius norm and $\|\mathbf{A}\|_2$ its spectral norm. We use $\text{sign}(z)$ as the function that is 1 when $z > 0$ and -1

otherwise. For a vector $\mathbf{v} \in \mathbb{R}^d$, we use $[\mathbf{v}]_i \in \mathbb{R}$ to denote the i -th component of the vector. For two sequence $\{a_k\}$ and $\{b_k\}$, we denote $a_k = O(b_k)$ if $|a_k| \leq C|b_k|$ for some absolute constant C , denote $a_k = \Omega(b_k)$ if $b_k = O(a_k)$, and denote $a_k = \Theta(b_k)$ if $a_k = O(b_k)$ and $a_k = \Omega(b_k)$. We also denote $a_k = o(b_k)$ if $\lim |a_k/b_k| = 0$.

Two-layer fully connected neural network. We consider a two-layer neural network described as follows: its first layer consists of m positive neurons and m negative neurons; its second layer parameters are fixed as $+1/m$ and $-1/m$ respectively for positive and negative neurons. Then the network can be written as $f(\mathbf{W}, \mathbf{x}) = F_{+1}(\mathbf{W}_{+1}, \mathbf{x}) - F_{-1}(\mathbf{W}_{-1}, \mathbf{x})$, where the partial network function of positive and negative neurons, i.e., $F_{+1}(\mathbf{W}_{+1}, \mathbf{x})$, $F_{-1}(\mathbf{W}_{-1}, \mathbf{x})$, are defined as:

$$F_j(\mathbf{W}_j, \mathbf{x}) = \frac{1}{m} \sum_{r=1}^m \sigma(\langle \mathbf{w}_{j,r}, \mathbf{x} \rangle) \quad (3.1)$$

for $j \in \{\pm 1\}$. Here, $\sigma(z)$ represents the activation function. For ReLU, $\sigma(z) = \max\{0, z\}$, and for leaky ReLU, $\sigma(z) = \max\{\gamma z, z\}$, where $\gamma \in (0, 1)$. $\mathbf{W}_j \in \mathbb{R}^{m \times d}$ is the collection of model weights associated with F_j , and $\mathbf{w}_{j,r} \in \mathbb{R}^d$ denotes the weight vector for the r -th neuron in \mathbf{W}_j . We use \mathbf{W} to denote the collection of all model weights.

Gradient Descent. Instead of considering the gradient flow (GF) that commonly studied in prior work on implicit bias, we use gradient descent (GD) to optimize the empirical loss on the training data

$$L_S(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \ell(y_i \cdot f(\mathbf{W}, \mathbf{x}_i)),$$

where $\ell(z) = \log(1 + \exp(-z))$ is the logistic loss, and $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ is the training data set. The gradient descent update rule of each neuron in the two-layer neural network can be written as

$$\mathbf{w}_{j,r}^{(t+1)} = \mathbf{w}_{j,r}^{(t)} - \eta \cdot \nabla_{\mathbf{w}_{j,r}} L_S(\mathbf{W}^{(t)}) = \mathbf{w}_{j,r}^{(t)} - \frac{\eta}{nm} \sum_{i=1}^n \ell_i'^{(t)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle) \cdot y_i \mathbf{x}_i \quad (3.2)$$

for all $j \in \{\pm 1\}$ and $r \in [m]$, where we introduce a shorthand notation $\ell_i'^{(t)} = \ell'[y_i \cdot f(\mathbf{W}^{(t)}, \mathbf{x}_i)]$ and assume the derivative of the ReLU activation function at 0 is $\sigma'(0) = 1$ without loss of generality. We initialize the gradient descent by Gaussian initialization, where all the entries of $\mathbf{W}^{(0)}$ are sampled from i.i.d. Gaussian distributions $\mathcal{N}(0, \sigma_0^2)$ with σ_0^2 being the variance.

4 Main Results

In this section, we present our main theoretical results. For training data $\mathcal{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \subseteq \mathbb{R}^d \times \{\pm 1\}$, let $R_{\min} = \min_i \|\mathbf{x}_i\|_2$, $R_{\max} = \max_i \|\mathbf{x}_i\|_2$, $p = \max_{i \neq k} |\langle \mathbf{x}_i, \mathbf{x}_k \rangle|$, and suppose $R = R_{\max}/R_{\min}$ is at most an absolute constant.

Theorem 4.1 (Leaky ReLU Networks). For two-layer neural network defined in (3.1) with leaky ReLU activation $\sigma(z) = \max\{\gamma z, z\}$, $\gamma \in (0, 1)$. Assume the training data satisfy $R_{\min}^2 \geq CR^2 \gamma^{-4} np$ for some sufficiently large constant C . For any $\delta \in (0, 1)$, if the learning rate $\eta \leq (CR_{\max}^2/nm)^{-1}$ and the initialization scale $\sigma_0 \leq \gamma(CR_{\max} \sqrt{\log(mn/\delta)})^{-1}$, then with probability at least $1 - \delta$ over the random initialization of gradient descent, the trained network satisfies:

- The ℓ_2 norm of each neuron increases to infinity at a logarithmic rate: $\|\mathbf{w}_{j,r}^{(t)}\|_2 = \Theta(\log(t))$ for all $j \in \{\pm 1\}$ and $r \in [m]$.
- Throughout the gradient descent trajectory, the stable rank of the weights $\mathbf{W}_j^{(t)}$ for all $j \in \{\pm 1\}$ satisfies,

$$\lim_{t \rightarrow \infty} \|\mathbf{W}_j^{(t)}\|_F^2 / \|\mathbf{W}_j^{(t)}\|_2^2 = 1,$$

with a convergence rate of $O(1/\log(t))$.

- The empirical risk under the logistic loss converges to zero at the following rate:

$$L_S(\mathbf{W}^{(t)}) = O\left(t^{-(1-c\gamma^{-2}R_{\min}^{-2}pn)/(1+c\gamma^{-2}R_{\min}^{-2}pn)}\right),$$

where c is a constant.

Remark 4.2. In Theorem 4.1, we show that when using the leaky ReLU activation function on nearly orthogonal training data, gradient descent asymptotically finds a network with a stable rank of \mathbf{W}_j equal to 1. Additionally, we provide a convergence rate for the training loss, which depends on the level of orthogonality in the training data. The more orthogonal the training data is, the faster convergence rate we can achieve. Specifically, when the training data is completely orthogonal ($p = 0$), we attain an $O(t^{-1})$ convergence rate. Furthermore, we analyze the rate of weight norm increase and the convergence rate of the stable rank for gradient descent, both of which exhibit a logarithmic order.

Theorem 4.3 (ReLU Networks). For two-layer neural network defined in (3.1) with ReLU activation $\sigma(z) = \max\{0, z\}$. Assume the training data satisfy $R_{\min}^2 \geq CR^2np$ for some sufficiently large constant C . For any $\delta \in (0, 1)$, if the neural network width $m \geq C \log(n/\delta)$, learning rate $\eta \leq (CR_{\max}^2/nm)^{-1}$ and initialization scale $\sigma_0 \leq (CR_{\max} \sqrt{\log(mn/\delta)})^{-1}$, then with probability at least $1 - \delta$ over the random initialization of gradient descent, the trained network satisfies:

- Throughout the gradient descent trajectory, the stable rank of the weights $\mathbf{W}_j^{(t)}$ for all $j \in \{\pm 1\}$ satisfies,

$$\limsup_{t \rightarrow \infty} \|\mathbf{W}_j^{(t)}\|_F^2 / \|\mathbf{W}_j^{(t)}\|_2^2 \leq c,$$

where c is a constant.

- The empirical risk under the logistic loss converges to zero at the following rate:

$$L_S(\mathbf{W}^{(t)}) = O\left(t^{-c_1(1-c_2R_{\min}^{-2}pn)/(1+c_2R_{\min}^{-2}pn)}\right),$$

where c_1, c_2 are constants.

Remark 4.4. For ReLU networks, we provide an example in the appendix concerning fully orthogonal training data and prove that the activation pattern during training depends solely on the initial activation state. Specifically, when training a two-layer ReLU network with gradient descent using such data, the stable rank of the network's weight matrix \mathbf{W}_j eventually converges to approximately 2. It is worth noting that this stable rank value is higher than the stable rank achieved by leaky ReLU networks, which is 1.

Comparison with previous work. One notable related work is Lyu et al. (2021), which also investigates the implicit bias of two-layer leaky ReLU networks. The main distinction between our work and Lyu et al. (2021) is the optimization method employed. We utilize gradient descent, whereas they utilize gradient flow. Additionally, our assumption is that the training data is nearly-orthogonal, while they assume the training data is symmetric. Our findings are more closely related to the work by Frei et al. (2022b), which investigates both gradient flow and gradient descent. In both our study and Frei et al. (2022b), we examine two-layer neural networks with leaky ReLU activations. However, they focus on networks trained via gradient flow, while we investigate networks trained using gradient descent. For the gradient descent approach, Frei et al. (2022b) provide a constant stable rank upper bound for smoothed leaky ReLU. In contrast, we prove that the stable rank of leaky ReLU networks converges to 1, aligning with the implicit bias of gradient flow proved in Frei et al. (2022b). Furthermore, they presented an $O(t^{-1/2})$ convergence rate for the empirical loss, whereas our convergence rate is dependent on the level of orthogonality and becomes $O(t^{-1})$ when the training data is completely orthogonal. It is worth noting that Lyu et al. (2021); Frei et al. (2022b) demonstrated that neural networks trained by gradient flow converge to a Karush-Kuhn-Tucker (KKT) point of the max-margin problem. We do not have such a result and have only proven that gradient descent can find a neural network with a low stable rank.

5 Overview of Proof Techniques

In this section, we discuss the key techniques we invent in our proofs to analyze the implicit bias of ReLU and leaky ReLU networks.

5.1 Refined Analysis of Decomposition Coefficient

Signal-noise decomposition, a technique initially introduced by Cao et al. (2022), is used to analyze the learning dynamics of two-layer convolutional networks. This method decomposes the convolutional filters into a linear combination of initial filters, signal vectors, and noise vectors, converting the neural network learning into a dynamical system of coefficients derived from the decomposition. In this work, we extend the signal-noise decomposition to *data-correlated decomposition* to facilitate the analysis of the training dynamic for two-layer fully connected neural networks.

Definition 5.1 (Data-correlated Decomposition). Let $\mathbf{w}_{j,r}^{(t)}$, $j \in \{\pm 1\}$, $r \in [m]$ be the weights of first-layer neurons at the t -th iteration of gradient descent. There exist unique coefficients $\rho_{j,r,i}^{(t)}$ such that

$$\mathbf{w}_{j,r}^{(t)} = \mathbf{w}_{j,r}^{(0)} + \sum_{i=1}^n \rho_{j,r,i}^{(t)} \cdot \|\mathbf{x}_i\|_2^{-2} \cdot \mathbf{x}_i. \quad (5.1)$$

By defining $\bar{\rho}_{j,r,i}^{(t)} := \rho_{j,r,i}^{(t)} \mathbb{1}(\rho_{j,r,i}^{(t)} \geq 0)$, $\underline{\rho}_{j,r,i}^{(t)} := \rho_{j,r,i}^{(t)} \mathbb{1}(\rho_{j,r,i}^{(t)} \leq 0)$, (5.1) can be further written as

$$\mathbf{w}_{j,r}^{(t)} = \mathbf{w}_{j,r}^{(0)} + \sum_{i=1}^n \bar{\rho}_{j,r,i}^{(t)} \cdot \|\mathbf{x}_i\|_2^{-2} \cdot \mathbf{x}_i + \sum_{i=1}^n \underline{\rho}_{j,r,i}^{(t)} \cdot \|\mathbf{x}_i\|_2^{-2} \cdot \mathbf{x}_i. \quad (5.2)$$

As an extension of the signal-noise decomposition first proposed in Cao et al. (2022) for analyzing two-layer convolutional networks, *data-correlated decomposition* defined in Definition 5.1 can be used to analyze two-layer fully-connected network, where the normalization factors $\|\mathbf{x}_i\|_2^{-2}$ are introduced to ensure that $\rho_{j,r,i}^{(t)} \approx \langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle$. This is also inspired by previous works by Lyu and Li (2019); Frei et al. (2022b), which demonstrate that \mathbf{W} converges to a KKT point of the max-margin problem. This implies that $\mathbf{w}_{j,r}^{(\infty)} / \|\mathbf{w}_{j,r}^{(\infty)}\|_2$ can be expressed as a linear combination of the training data $\{\mathbf{x}_i\}_{i=1}^n$, with the coefficient λ_i corresponding to $\rho_{j,r,i}^{(t)}$ in our analysis. With the help of such decomposition techniques, one can reduce the study of neural network training process to a careful assessment of the coefficients $\bar{\rho}_{j,r,i}^{(t)}$, $\underline{\rho}_{j,r,i}^{(t)}$ throughout training. This technique does not rely on the strictly increasing and smoothness properties of the activation function and will serve as the foundation for our analysis. Let us first investigate the update rule of the coefficient $\bar{\rho}_{j,r,i}^{(t)}$, $\underline{\rho}_{j,r,i}^{(t)}$.

Lemma 5.2. The coefficients $\bar{\rho}_{j,r,i}^{(t)}$, $\underline{\rho}_{j,r,i}^{(t)}$ defined in Definition 5.1 satisfy the following iterative equations:

$$\bar{\rho}_{j,r,i}^{(0)}, \underline{\rho}_{j,r,i}^{(0)} = 0, \quad (5.3)$$

$$\bar{\rho}_{j,r,i}^{(t+1)} = \bar{\rho}_{j,r,i}^{(t)} - \frac{\eta}{nm} \cdot \ell_i'(t) \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle) \cdot \|\mathbf{x}_i\|_2^2 \cdot \mathbb{1}(y_i = j), \quad (5.4)$$

$$\underline{\rho}_{j,r,i}^{(t+1)} = \underline{\rho}_{j,r,i}^{(t)} + \frac{\eta}{nm} \cdot \ell_i'(t) \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle) \cdot \|\mathbf{x}_i\|_2^2 \cdot \mathbb{1}(y_i = -j), \quad (5.5)$$

for all $r \in [m]$, $j \in \{\pm 1\}$ and $i \in [n]$.

To study implicit bias, the first main challenge is to generalize the decomposition coefficient analysis to infinite time. The signal-noise decomposition used in Cao et al. (2022); Kou et al. (2023) require early stopping with threshold T^* to facilitate their analysis. They only provided upper bounds of $4 \log(T^*)$ for $\bar{\rho}_{j,r,i}^{(t)}$, $|\underline{\rho}_{j,r,i}^{(t)}|$ (See Proposition 5.3 in Cao et al. (2022), Proposition 5.2 in Kou et al. (2023)), and then carried out a two-stage analysis. To obtain upper bounds for $\bar{\rho}_{j,r,i}^{(t)}$, $|\underline{\rho}_{j,r,i}^{(t)}|$, they used an upper bound for $|\ell_i'(t)|$ and directly plugged it into (5.4) and (5.5) to demonstrate that $\bar{\rho}_{j,r,i}^{(t)}$

199 and $|\rho_{j,r,i}^{(t)}|$ would not exceed $4 \log(T^*)$, which is a fixed value related to the early stopping threshold.
 200 Therefore, dealing with infinite time requires new techniques. To overcome this difficulty, we propose
 201 a *refined analysis of decomposition coefficients* which generalizes Cao et al. (2022)’s technique. We
 202 first give the following key lemma.

203 **Lemma 5.3.** For non-negative real number sequence $\{x_t\}_{t=0}^\infty$ satisfying

$$C_1 \exp(-x_t) \leq x_{t+1} - x_t \leq C_2 \exp(-x_t), \quad (5.6)$$

204 it holds that

$$\log(\exp(-x_0) + C_1 \cdot t) \leq x_t \leq \log(\exp(-x_0) + C_2 \exp(C_2) \cdot t). \quad (5.7)$$

205 We can establish the relationship between (5.4), (5.5) and inequality (5.6) if we are able to express
 206 $|\ell_i^{(t)}|$ using coefficients $\bar{\rho}_{j,r,i}^{(t)}$ and $|\rho_{j,r,i}^{(t)}|$. To achieve this, we can first approximate $\ell_i^{(t)}$ using
 207 the margin $y_i f(\mathbf{W}^{(t)}, \mathbf{x}_i)$ and then approximate $F_j(\mathbf{W}_j^{(t)}, \mathbf{x}_i)$ using the coefficients $\rho_{j,r,i}^{(t)}$. The
 208 approximation is given as follows:

$$\ell_i^{(t)} = \Theta(\exp(-y_i f(\mathbf{W}^{(t)}, \mathbf{x}_i))) = \Theta(\exp(F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i) - F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i))), \quad (5.8)$$

$$\left| F_j(\mathbf{W}_j^{(t)}, \mathbf{x}_i) - \frac{1}{m} \sum_{r=1}^m \rho_{j,r,i}^{(t)} \right| \leq \sum_{i' \neq i} \left(\frac{1}{m} \sum_{r=1}^m |\rho_{j,r,i'}^{(t)}| R_{\min}^{-2} p \right), \quad (5.9)$$

209 From (5.9), one can see that we need to decouple $\ell_i^{(t)}$ from $|\rho_{j,r,i'}^{(t)}| (i' \neq i)$. In order to accomplish
 210 this, we also prove the following lemma, which demonstrates that the ratio between $\sum_{r=1}^m |\rho_{j,r,i}^{(t)}|$
 211 and $\sum_{r=1}^m |\rho_{j,r,i'}^{(t)}| (i' \neq i)$ will maintain a constant order throughout the training process. Here, we
 212 present the lemma for leaky ReLU networks.

213 **Lemma 5.4** (leaky ReLU automatic balance). For two-layer leaky ReLU network defined in (3.1),
 214 for any $t \geq 0$, we have

$$\sum_{r=1}^m |\rho_{j,r,i}^{(t)}| \geq c \gamma^2 \sum_{r=1}^m |\rho_{j,r,i'}^{(t)}|, \quad (5.10)$$

215 for any $j \in \{\pm 1\}$ and $r, r' \in [m]$, where c is a constant.

216 By Lemma 5.4, we can approximate the neural network output using equation (5.9). This approxima-
 217 tion expresses the output $F_j(\mathbf{W}_j^{(t)}, \mathbf{x}_i)$ as a sum of the coefficients $\rho_{j,r,i}^{(t)}$:

$$F_j(\mathbf{W}_j^{(t)}, \mathbf{x}_i) \approx \frac{1 \pm c \gamma^2 R_{\min}^{-2} p n}{m} \sum_{r=1}^m \rho_{j,r,i}^{(t)}. \quad (5.11)$$

218 By combining (5.4), (5.5), (5.8), and (5.11), we obtain the following relationship:

$$\frac{1}{m} \sum_{r=1}^m |\rho_{j,r,i}^{(t+1)}| - \frac{1}{m} \sum_{r=1}^m |\rho_{j,r,i}^{(t)}| = \Theta\left(\frac{\eta \|\mathbf{x}_i\|_2^2}{nm}\right) \cdot \exp\left(-\frac{1 \pm c \gamma^2 R_{\min}^{-2} p n}{m} \sum_{r=1}^m |\rho_{j,r,i}^{(t)}|\right).$$

219 This relationship aligns with the form of (5.6), if we set $x_t = \frac{1 \pm c \gamma^2 R_{\min}^{-2} p n}{m} \sum_{r=1}^m |\rho_{j,r,i}^{(t)}|$. Thus, we
 220 can directly apply Lemma 5.3 to gain insights into the logarithmic rate of increase for the average
 221 magnitudes of the coefficients $\frac{1}{m} \sum_{r=1}^m |\rho_{j,r,i}^{(t)}|$. This analysis provides a deeper understanding of the
 222 dynamics and convergence properties of the coefficient magnitudes during the training process. In
 223 the case of ReLU networks, we have the following lemma that provides automatic balance:

224 **Lemma 5.5** (ReLU automatic balance). For two-layer ReLU network defined in (3.1), there exists a
 225 constant c such that for any $t \geq 0$, we have

$$|\rho_{y_i, r, i}^{(t)}| \geq c |\rho_{j, r', i'}^{(t)}|,$$

226 for any $j \in \{\pm 1\}$, $r \in S_i^{(0)} := \{r \in [m] : \langle \mathbf{w}_{y_i, r}, \mathbf{x}_i \rangle \geq 0\}$, $r' \in [m]$ and $i, i' \in [n]$.

227 The automatic balance lemma guarantees that the magnitudes of coefficients related to the neurons of
 228 class y_i , which are activated by \mathbf{x}_i during initialization, dominate those of other classes. With the
 229 help of Lemma 5.5, we can get the following approximation for the margin $y_i f(\mathbf{W}^{(t)}, \mathbf{x}_i)$:

$$F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i) - F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i) \approx \frac{1 \pm cR_{\min}^{-2}pn}{m} \sum_{r \in S_i^{(0)}} \rho_{y_i, r, i}^{(t)}. \quad (5.12)$$

230 By combining (5.4), (5.5), (5.8) and (5.12), we obtain the following relationship:

$$\sum_{r \in S_i^{(0)}} |\rho_{y_i, r, i}^{(t+1)}| - \sum_{r \in S_i^{(0)}} |\rho_{y_i, r, i}^{(t)}| = \Theta\left(\frac{\eta \|\mathbf{x}_i\|_2^2 |S_i^{(0)}|}{nm}\right) \cdot \exp\left(-\frac{1 \pm cR_{\min}^{-2}pn}{m} \sum_{r \in S_i^{(0)}} \rho_{y_i, r, i}^{(t)}\right),$$

231 which precisely matches the form of (5.6) by setting $x_t = \frac{1 \pm cR_{\min}^{-2}pn}{m} \sum_{r \in S_i^{(0)}} \rho_{y_i, r, i}^{(t)}$. Therefore, we
 232 can directly apply Lemma 5.3 and obtain the logarithmic increasing rate of $|\rho_{y_i, r, i}^{(t)}|$ for $r \in S_i^{(0)}$,
 233 which provides further insights into the behavior of the coefficients.

234 5.2 Analysis of Activation Pattern

235 One notable previous work (Frei et al., 2022b) provided a constant upper bound for stable rank of
 236 two-layer smoothed leaky ReLU networks trained by gradient descent and an $O(t^{-1/2})$ convergence
 237 rate for training loss in their Theorem 4.2. To achieve better stable rank bound and better convergence
 238 rate, we characterize the activation pattern of leaky ReLU network neurons after certain threshold
 239 time T in the following lemma.

240 **Lemma 5.6** (leaky ReLU activation pattern). Let $T = C\eta^{-1}nmR_{\max}^{-2}$. For two-layer leaky ReLU
 241 network defined in (3.1), for any $t \geq T$, it holds that

$$\text{sign}(\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle) = jy_i,$$

242 for any $j \in \{\pm 1\}$ and $r \in [m]$.

243 Lemma 5.6 indicates that the activation pattern will not change after time T . Given Lemma 5.6, we
 244 can get $\sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle) = \gamma$ for $j \neq y_i$ and $\sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle) = 1$ for $j = y_i$. Plugging this into (5.4)
 245 and (5.5) can give the following useful lemma.

246 **Lemma 5.7.** Let T be defined in Lemma 5.6. For $t \geq T$, it holds that

$$\begin{aligned} \bar{\rho}_{y_i, r, i}^{(t)} - \bar{\rho}_{y_i, r, i}^{(T)} &= \bar{\rho}_{y_i, r', i}^{(t)} - \bar{\rho}_{y_i, r', i}^{(T)}, \underline{\rho}_{-y_i, r, i}^{(t)} - \underline{\rho}_{-y_i, r, i}^{(T)} = \underline{\rho}_{-y_i, r', i}^{(t)} - \underline{\rho}_{-y_i, r', i}^{(T)}, \\ \bar{\rho}_{y_i, r, i}^{(t)} - \bar{\rho}_{y_i, r, i}^{(T)} &= (\rho_{-y_i, r', i}^{(t)} - \rho_{-y_i, r', i}^{(T)})/\gamma, \end{aligned}$$

247 for any $i \in [n]$ and $r, r' \in [m]$.

248 This lemma reveals that beyond a certain time threshold T , the increase in $\rho_{j,r,i}^{(t)}$ is consistent across
 249 neurons within the same positive or negative class. However, for neurons belonging to the oppose
 250 class, this increment in $\rho_{j,r,i}^{(t)}$ is scaled by a factor equivalent to the slope of the leaky ReLU function
 251 γ . From this and (5.1), we can demonstrate that $\|\mathbf{w}_{j,r}^{(t)} - \mathbf{w}_{j,r'}^{(t)}\|_2 (r \neq r')$ can be upper bounded by a
 252 constant, leading to the following inequalities:

$$\|\mathbf{W}_j^{(t)}\|_F^2 \leq m\|\mathbf{w}_{j,1}^{(t)}\|_2^2 + mC_1\|\mathbf{w}_{j,1}^{(t)}\|_2 + mC_2, \quad \|\mathbf{W}_j^{(t)}\|_2^2 \geq m\|\mathbf{w}_{j,1}^{(t)}\|_2^2 - mC_3\|\mathbf{w}_{j,1}^{(t)}\|_2 - mC_4.$$

253 Considering that $\|\mathbf{w}_{j,r}^{(t)}\|_2 = \Theta(\log t)$, the stable rank of $\mathbf{W}_j^{(t)}$ naturally converges to a value of
 254 1. As for the convergence rate analysis, by applying the activation pattern to (5.4) and (5.5), and
 255 using Lemma 5.3 again, we can achieve a more fine-grained analysis to obtain a more precise
 256 approximation of $\rho_{j,r,i}^{(t)}/\log(t)$. This leads to $y_i f(\mathbf{W}^{(t)}, \mathbf{x}_i) \geq \frac{1 - c\gamma^{-2}R_{\min}^{-2}pn}{1 + c\gamma^{-2}R_{\min}^{-2}pn} \log(t)$ for sufficiently
 257 large t . Consequently, this results in an $O(t^{-(1 - c\gamma^{-2}R_{\min}^{-2}pn)/(1 + c\gamma^{-2}R_{\min}^{-2}pn)})$ convergence rate for
 258 the training loss, which is faster than $O(t^{-1/2})$ as demonstrated by Theorem 4.2 in (Frei et al., 2022b),

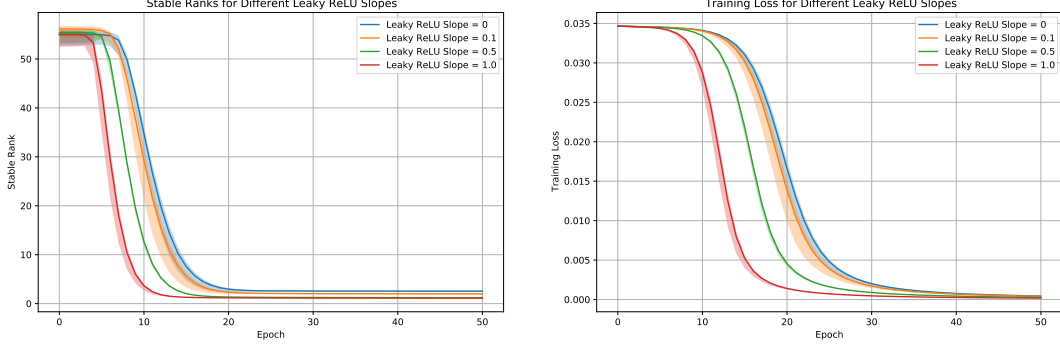


Figure 1: Stable ranks and training loss for different leaky ReLU slopes γ across multiple runs. A slope of 1 corresponds to linear activation, while a slope of 0 corresponds to ReLU activation. Each line represents the mean stable rank or training loss for a given leaky ReLU slope, while the shaded regions indicate the variability of the values (± 3 times the standard deviation) across the 5 runs.

if the training data is nearly orthogonal enough. For ReLU networks, we can partially characterize the activation pattern as illustrated in the following lemma.

Lemma 5.8. (ReLU activation pattern) For two-layer ReLU network defined in (3.1), for any $i \in [n]$ and $r \in S_i^{(0)} := \{r \in [m] : \langle \mathbf{w}_{y_i, r}^{(0)}, \mathbf{x}_i \rangle \geq 0\}$, we have

$$\text{sign}(\langle \mathbf{w}_{y_i, r}^{(t)}, \mathbf{x}_i \rangle) = \text{sign}(\langle \mathbf{w}_{y_i, r}^{(0)}, \mathbf{x}_i \rangle) = 1,$$

for any $t \geq 0$, or in other words, $S_i^{(0)} \subseteq S_i^{(t)} := \{r \in [m] : \langle \mathbf{w}_{y_i, r}^{(t)}, \mathbf{x}_i \rangle \geq 0\}$.

Lemma 5.8 suggests that once the neuron of class y_i is activated by \mathbf{x}_i during initialization, it will continue to remain activated throughout the training process. Leveraging such an activation pattern, $|S_i^{(0)}| = \Theta(m)$, and $\rho_{y_i, r, i}^{(t)} = \Theta(\log t)$, we can establish a lower bound for $\|\mathbf{W}_j^{(t)}\|_2$ as $\Omega(\sqrt{mn} \log(t))$. This matches the trivial upper bound of $\|\mathbf{W}_j^{(t)}\|_F$ of order $\Theta(\sqrt{mn} \log(t))$, thus providing us with a constant upper bound for the stable rank for ReLU networks. The convergence rate analysis of the training loss also utilizes this activation pattern and $|S_i^{(0)}| = \Theta(m)$ to ensure that the margin satisfies $y_i f(\mathbf{W}^{(t)}, \mathbf{x}_i) = \Omega(\rho_{y_i, r, i}^{(t)})$, $r \in S_i^{(0)}$. As a result, as time t goes to infinity, the ratio $y_i f(\mathbf{W}^{(t)}, \mathbf{x}_i) / \log(t)$ has no correlation with the width m .

6 Experiments

In this section, we present simulations of both synthetic and real data to back up our theoretical analysis in the previous section.

Synthetic-data experiments. Here we generate a synthetic mixture of Gaussian data as follows:

Let $\boldsymbol{\mu} \in \mathbb{R}^d$ be a fixed vector representing the signal contained in each data point. Each data point (\mathbf{x}, y) with predictor $\mathbf{x} \in \mathbb{R}^d$ and label $y \in \{-1, 1\}$ is generated from a distribution \mathcal{D} , which we specify as follows:

1. The label y is generated as a Rademacher random variable, i.e. $\mathbb{P}[y = 1] = \mathbb{P}[y = -1] = 1/2$.
2. A noise vector $\boldsymbol{\xi}$ is generated from the Gaussian distribution $\mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I}_d)$. And \mathbf{x} is assigned as $y \cdot \boldsymbol{\mu} + \boldsymbol{\xi}$ where $\boldsymbol{\mu}$ is a fixed feature vector.

Specifically, we set training data size $n = 10$, $d = 784$ and train the NN with stochastic gradient descent with batch size 64 and learning rate 0.1 for 50 epochs. We set $\boldsymbol{\mu}$ to be a feature randomly drawn from $\mathcal{N}(\mathbf{0}, 10^{-4} \mathbf{I}_d)$. We then generate the noise vector $\boldsymbol{\xi}$ from the Gaussian distribution $\mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I})$ with fixed standard deviation $\sigma_p = 1$. We train the FNN model defined in Section 3 with ReLU (or leaky-ReLU) activation function and width $m = 100$. As we can infer from Figure 1, the stable rank will decrease faster for larger leaky ReLU slopes and have a smaller value when epoch $t \rightarrow \infty$.

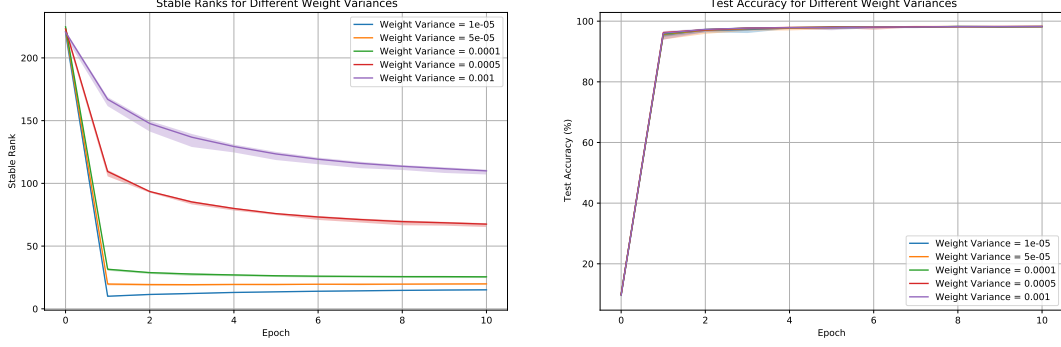


Figure 2: Stable ranks and test errors for different weight variances across multiple runs (ReLU Activation Function). Each line represents the mean stable rank or test accuracy for a given weight variance, while the shaded regions indicate the variability of the values (± 3 times the standard deviation) across the 5 runs.

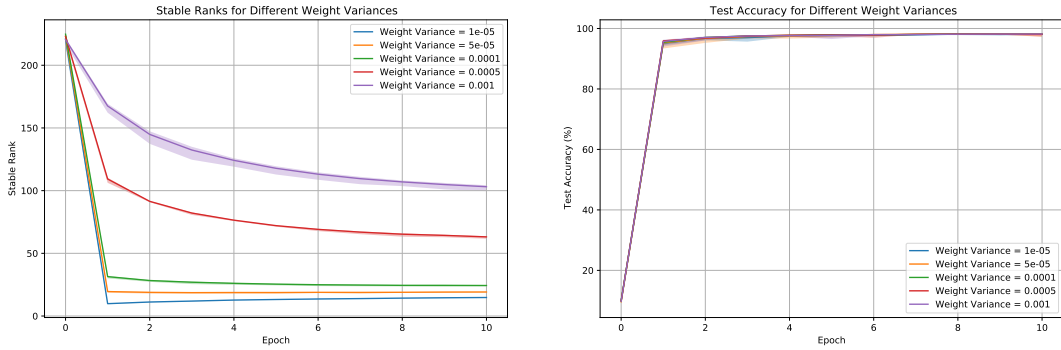


Figure 3: Stable ranks and test errors for different weight variances across multiple runs (leaky-ReLU Activation Function with slope 0.1). Each line represents the mean stable rank or test accuracy for a given weight variance, while the shaded regions indicate the variability of the values (± 3 times the standard deviation) across the 5 runs.

289 **Real-data experiments on MNIST dataset.** Here we train a two-layer feed-forward neural
 290 network defined in Section 3 with ReLU (or leaky-ReLU) functions. The number of widths is
 291 set as $m = 1000$. We use the Gaussian initialization and consider different weight variance
 292 $\sigma_0 \in \{0.00001, 0.00005, 0.0001, 0.0005, 0.001\}$. We train the NN with stochastic gradient de-
 293 scent with batch size 64 and learning rate 0.1 for 10 epochs. As we can infer from Figures 2 and 3,
 294 the stable rank of ReLU or leaky ReLU networks will largely depend on the initialization and the
 295 training time. When initialization is sufficiently small, the stable rank will quickly decrease to a small
 296 value compared to its initialization values.

297 7 Conclusion and Future Work

298 This paper employs a data-correlated decomposition technique to examine the implicit bias of two-
 299 layer ReLU and Leaky ReLU networks trained using gradient descent. By analyzing the training
 300 dynamics, we provide precise characterizations of the weight matrix stable rank limits for both ReLU
 301 and Leaky ReLU cases, demonstrating that both scenarios will yield a network with a low stable
 302 rank. Additionally, we present an empirical analysis of the convergence rate of the loss function. A
 303 crucial direction for future research is to explore the directional convergence of the weight matrix in
 304 neural networks trained via gradient descent. It would be valuable to investigate whether the weight
 305 matrix eventually converges to a Karush-Kuhn-Tucker (KKT) point of the max-margin problem.
 306 Furthermore, it is important to extend our analysis to fully understand the neuron activation patterns
 307 in the case of ReLU activations. Specifically, we can explore whether certain neurons continuously
 308 switch their activation states an infinite number of times throughout the training process or if the
 309 activation patterns stabilize after surpassing a certain threshold.

References

- ALLEN-ZHU, Z. and LI, Y. (2020). Towards understanding ensemble, knowledge distillation and self-distillation in deep learning. *arXiv preprint arXiv:2012.09816* .
- ARORA, S., COHEN, N., HU, W. and LUO, Y. (2019). Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems* **32**.
- BARTLETT, P. L., LONG, P. M., LUGOSI, G. and TSIGLER, A. (2020). Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences* .
- BELKIN, M., HSU, D., MA, S. and MANDAL, S. (2019). Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences* **116** 15849–15854.
- BELKIN, M., HSU, D. and XU, J. (2020). Two models of double descent for weak features. *SIAM Journal on Mathematics of Data Science* **2** 1167–1180.
- CAO, Y., CHEN, Z., BELKIN, M. and GU, Q. (2022). Benign overfitting in two-layer convolutional neural networks. *arXiv preprint arXiv:2202.06526* .
- CAO, Y., GU, Q. and BELKIN, M. (2021). Risk bounds for over-parameterized maximum margin classification on sub-gaussian mixtures. *Advances in Neural Information Processing Systems* **34**.
- CHATTERJI, N. S. and LONG, P. M. (2021). Finite-sample analysis of interpolating linear classifiers in the overparameterized regime. *Journal of Machine Learning Research* **22** 129–1.
- CHIZAT, L. and BACH, F. (2020). Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss. In *Conference on Learning Theory*. PMLR.
- FREI, S., CHATTERJI, N. S. and BARTLETT, P. (2022a). Benign overfitting without linearity: Neural network classifiers trained by gradient descent for noisy linear data. In *Conference on Learning Theory*. PMLR.
- FREI, S., VARDI, G., BARTLETT, P. L., SREBRO, N. and HU, W. (2022b). Implicit bias in leaky relu networks trained on high-dimensional data. *arXiv preprint arXiv:2210.07082* .
- GUNASEKAR, S., LEE, J. D., SOUDRY, D. and SREBRO, N. (2018). Implicit bias of gradient descent on linear convolutional networks. In *Advances in Neural Information Processing Systems*.
- HASTIE, T., MONTANARI, A., ROSSET, S. and TIBSHIRANI, R. J. (2022). Surprises in high-dimensional ridgeless least squares interpolation. *The Annals of Statistics* **50** 949–986.
- JELASSI, S. and LI, Y. (2022). Towards understanding how momentum improves generalization in deep learning. In *International Conference on Machine Learning*. PMLR.
- Ji, Z. and TELGARSKY, M. (2018). Gradient descent aligns the layers of deep linear networks. *arXiv preprint arXiv:1810.02032* .
- Ji, Z. and TELGARSKY, M. (2020). Directional convergence and alignment in deep learning. *Advances in Neural Information Processing Systems* **33** 17176–17186.
- KOU, Y., CHEN, Z., CHEN, Y. and GU, Q. (2023). Benign overfitting for two-layer relu networks. *arXiv preprint arXiv:2303.04145* .
- LI, Z., ZHOU, Z.-H. and GRETTON, A. (2021). Towards an understanding of benign overfitting in neural networks. *arXiv preprint arXiv:2106.03212* .
- LIAO, Z., COUILLET, R. and MAHONEY, M. (2020). A random matrix analysis of random fourier features: beyond the gaussian kernel, a precise phase transition, and the corresponding double descent. In *34th Conference on Neural Information Processing Systems (NeurIPS 2020)*.

352 LYU, K. and LI, J. (2019). Gradient descent maximizes the margin of homogeneous neural networks.
353 *arXiv preprint arXiv:1906.05890* .

354 LYU, K., LI, Z., WANG, R. and ARORA, S. (2021). Gradient descent on two-layer nets: Margin
355 maximization and simplicity bias. *Advances in Neural Information Processing Systems* **34** 12978–
356 12991.

357 MEI, S. and MONTANARI, A. (2019). The generalization error of random features regression:
358 Precise asymptotics and double descent curve. *arXiv preprint arXiv:1908.05355* .

359 PHUONG, M. and LAMPERT, C. H. (2021). The inductive bias of relu networks on orthogonally
360 separable data. In *International Conference on Learning Representations*.

361 RAZIN, N. and COHEN, N. (2020). Implicit regularization in deep learning may not be explainable
362 by norms. *Advances in neural information processing systems* **33** 21174–21187.

363 SAFRAN, I., VARDI, G. and LEE, J. D. (2022). On the effective number of linear regions in
364 shallow univariate relu networks: Convergence guarantees and implicit bias. *arXiv preprint*
365 *arXiv:2205.09072* .

366 SARUSSI, R., BRUTZKUS, A. and GLOBERSON, A. (2021). Towards understanding learning in
367 neural networks with linear teachers. In *International Conference on Machine Learning*. PMLR.

368 SHEN, R., BUBECK, S. and GUNASEKAR, S. (2022). Data augmentation as feature manipulation: a
369 story of desert cows and grass cows. *arXiv preprint arXiv:2203.01572* .

370 SOUDRY, D., HOFFER, E., NACSON, M. S., GUNASEKAR, S. and SREBRO, N. (2018). The
371 implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research* **19**
372 2822–2878.

373 TIMOR, N., VARDI, G. and SHAMIR, O. (2023). Implicit regularization towards rank minimization
374 in relu networks. In *International Conference on Algorithmic Learning Theory*. PMLR.

375 VARDI, G. (2022). On the implicit bias in deep-learning algorithms. *arXiv preprint arXiv:2208.12591*
376 .

377 VARDI, G., SHAMIR, O. and SREBRO, N. (2022a). On margin maximization in linear and relu
378 networks. *Advances in Neural Information Processing Systems* **35** 37024–37036.

379 VARDI, G., YEHUDAI, G. and SHAMIR, O. (2022b). Gradient methods provably converge to
380 non-robust networks. *arXiv preprint arXiv:2202.04347* .

381 VERSHYNIN, R. (2018). *High-Dimensional Probability: An Introduction with Applications in Data*
382 *Science*. Cambridge Series in Statistical and Probabilistic Mathematics, Cambridge University
383 Press.

384 WANG, K. and THRAMPOULIDIS, C. (2022). Binary classification of gaussian mixtures: Abundance
385 of support vectors, benign overfitting, and regularization. *SIAM Journal on Mathematics of Data*
386 *Science* **4** 260–284.

387 WU, D. and XU, J. (2020). On the optimal weighted ℓ_2 regularization in overparameterized linear
388 regression. *Advances in Neural Information Processing Systems* **33**.

A Preliminary Lemmas

In this section, we present some pivotal lemmas that illustrate some important properties of the data and neural network parameters at their random initialization and provide the update rule of coefficients from data-correlated decomposition.

Now turning to network initialization, the following lemma studies the inner product between a randomly initialized neural network neuron $\mathbf{w}_{j,r}^{(0)}$ ($j \in \{\pm 1\}$ and $r \in [m]$) and the training data. The calculations characterize how the neural network at initialization randomly captures the information in training data.

Lemma A.1. Suppose that $d = \Omega(\log(mn/\delta))$, $m = \Omega(\log(1/\delta))$. Then with probability at least $1 - \delta$,

$$\begin{aligned} \sigma_0^2 d/2 &\leq \|\mathbf{w}_{j,r}^{(0)}\|_2^2 \leq 3\sigma_0^2 d/2, \\ |\langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle| &\leq \sqrt{2 \log(8mn/\delta)} \cdot \sigma_0 R_{\max} \end{aligned}$$

for all $r \in [m]$, $j \in \{\pm 1\}$ and $i \in [n]$.

Proof of Lemma A.1. First of all, the initial weights $\mathbf{w}_{j,r}^{(0)} \sim \mathcal{N}(\mathbf{0}, \sigma_0^2 \mathbf{I})$. By Bernstein's inequality, with probability at least $1 - \delta/(4m)$ we have

$$\|\mathbf{w}_{j,r}^{(0)}\|_2^2 - \sigma_0^2 d = O(\sigma_0^2 \cdot \sqrt{d \log(8m/\delta)}).$$

Therefore, if we set appropriately $d = \Omega(\log(m/\delta))$, we have with probability at least $1 - \delta/2$, for all $j \in \{\pm 1\}$ and $r \in [m]$,

$$\sigma_0^2 d/2 \leq \|\mathbf{w}_{j,r}^{(0)}\|_2^2 \leq 3\sigma_0^2 d/2.$$

Under definition, we have $\|\mathbf{x}_i\|_2 \leq R_{\max}$ for all $i \in [n]$. It is clear that for each j, r , $\langle \mathbf{w}_{j,r}^{(0)}, \boldsymbol{\mu} \rangle$ is a Gaussian random variable with mean zero and variance $\sigma_0^2 \|\mathbf{x}_i\|_2^2$. Therefore, by Gaussian tail bound and union bound, with probability at least $1 - \delta/2$,

$$|\langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle| \leq \sqrt{2 \log(8mn/\delta)} \cdot \sigma_0 R_{\max}.$$

407

□

Next, we denote $S_i^{(0)}$ as $\{r \in [m] : \langle \mathbf{w}_{y_i,r}^{(0)}, \mathbf{x}_i \rangle > 0\}$. We give a lower bound of $|S_i^{(0)}|$ in the following two lemmas.

Lemma A.2. Suppose that $\delta > 0$ and $m \geq 50 \log(2n/\delta)$. Then with probability at least $1 - \delta$,

$$0.4m \leq |S_i^{(0)}| \leq 0.6m, \forall i \in [n].$$

Proof of Lemma A.2. Note that $|S_i^{(0)}| = \sum_{r=1}^m \mathbb{1}[\langle \mathbf{w}_{y_i,r}^{(0)}, \mathbf{x}_i \rangle > 0]$ and $P(\langle \mathbf{w}_{y_i,r}^{(0)}, \mathbf{x}_i \rangle > 0) = 1/2$, then by Hoeffding's inequality, with probability at least $1 - \delta/n$, we have

$$\left| \frac{|S_i^{(0)}|}{m} - \frac{1}{2} \right| \leq \sqrt{\frac{\log(2n/\delta)}{2m}}.$$

Therefore, as long as $m \geq 50 \log(2n/\delta)$, by applying union bound, with probability at least $1 - \delta$, we have

$$0.4m \leq |S_i^{(0)}| \leq 0.6m, \forall i \in [n].$$

415

□

Now we give the update rule of coefficients from data-correlated decomposition. We will begin by analyzing the coefficients in the data-correlated decomposition in Definition 5.1. The following lemma presents an iterative expression for the coefficients.

419 **Lemma A.3.** (Restatement of Lemma 5.2) The coefficients $\bar{\rho}_{j,r,i}^{(t)}, \underline{\rho}_{j,r,i}^{(t)}$ defined in Definition 5.1
 420 satisfy the following iterative equations:

$$\begin{aligned}\bar{\rho}_{j,r,i}^{(0)}, \underline{\rho}_{j,r,i}^{(0)} &= 0, \\ \bar{\rho}_{j,r,i}^{(t+1)} &= \bar{\rho}_{j,r,i}^{(t)} - \frac{\eta}{nm} \cdot \ell_i^{(t)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle) \cdot \|\mathbf{x}_i\|_2^2 \cdot \mathbb{1}(y_i = j), \\ \underline{\rho}_{j,r,i}^{(t+1)} &= \underline{\rho}_{j,r,i}^{(t)} + \frac{\eta}{nm} \cdot \ell_i^{(t)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle) \cdot \|\mathbf{x}_i\|_2^2 \cdot \mathbb{1}(y_i = -j),\end{aligned}$$

421 for all $r \in [m], j \in \{\pm 1\}$ and $i \in [n]$.

422 *Proof of Lemma A.3.* First, we iterate the gradient descent update rule (3.2) t times and get

$$\mathbf{w}_{j,r}^{(t+1)} = \mathbf{w}_{j,r}^{(0)} - \frac{\eta}{nm} \sum_{s=0}^t \sum_{i=1}^n \ell_i^{(s)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(s)}, \mathbf{x}_i \rangle) \cdot j y_i \mathbf{x}_i.$$

423 According to the definition of $\rho_{j,r,i}^{(t)}$, we have

$$\mathbf{w}_{j,r}^{(t)} = \mathbf{w}_{j,r}^{(0)} + \sum_{i=1}^n \rho_{j,r,i}^{(t)} \cdot \|\mathbf{x}_i\|_2^{-2} \cdot \mathbf{x}_i.$$

424 Therefore, we have the unique representation

$$\rho_{j,r,i}^{(t)} = -\frac{\eta}{nm} \sum_{s=0}^t \ell_i^{(s)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(s)}, \mathbf{x}_i \rangle) \cdot \|\mathbf{x}_i\|_2^2 \cdot j y_i.$$

425 Now with the notation $\bar{\rho}_{j,r,i}^{(t)} := \rho_{j,r,i}^{(t)} \mathbb{1}(\rho_{j,r,i}^{(t)} \geq 0)$, $\underline{\rho}_{j,r,i}^{(t)} := \rho_{j,r,i}^{(t)} \mathbb{1}(\rho_{j,r,i}^{(t)} \leq 0)$ and the fact
 426 $\ell_i^{(s)} < 0$, we get

$$\bar{\rho}_{j,r,i}^{(t)} = -\frac{\eta}{nm} \sum_{s=0}^t \ell_i^{(s)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(s)}, \mathbf{x}_i \rangle) \cdot \|\mathbf{x}_i\|_2^2 \cdot \mathbb{1}(y_i = j), \quad (\text{A.1})$$

$$\underline{\rho}_{j,r,i}^{(t)} = \frac{\eta}{nm} \sum_{s=0}^t \ell_i^{(s)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(s)}, \mathbf{x}_i \rangle) \cdot \|\mathbf{x}_i\|_2^2 \cdot \mathbb{1}(y_i = -j). \quad (\text{A.2})$$

427 Writing out the iterative versions of (A.1) and (A.2) completes the proof. \square

428 B Coefficient Analysis of Leaky ReLU

429 In this section, we establish a series of results on the data-correlated decomposition for two-layer
 430 leaky ReLU network defined as

$$\begin{aligned}f(\mathbf{W}^{(t)}, \mathbf{x}) &= F_{+1}(\mathbf{W}_{+1}^{(t)}, \mathbf{x}) - F_{-1}(\mathbf{W}_{-1}^{(t)}, \mathbf{x}) \\ &= \frac{1}{m} \sum_{r=1}^m \sigma(\langle \mathbf{w}_{+1,r}^{(t)}, \mathbf{x} \rangle) - \frac{1}{m} \sum_{r=1}^m \sigma(\langle \mathbf{w}_{-1,r}^{(t)}, \mathbf{x} \rangle), \\ \sigma(z) &= \max\{\gamma z, z\}, \gamma \in (0, 1).\end{aligned} \quad (\text{B.1})$$

431 The results in Section B, C and D are based on Lemma A.1, which hold with high probability. Denote
 432 by $\mathcal{E}_{\text{prelim}}$ the event that Lemma A.1 in Section A holds (for a given δ , we see $\mathbb{P}(\mathcal{E}_{\text{prelim}}) \geq 1 - \delta$).
 433 For simplicity and clarity, we state all the results in Section B, C and D conditional on $\mathcal{E}_{\text{prelim}}$.

434 Denote $\beta = \max_{i,j,r} \{|\langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle|\}$, $R_{\max} = \max_{i \in [n]} \|\mathbf{x}_i\|_2$, $R_{\min} = \min_{i \in [n]} \|\mathbf{x}_i\|_2$, $p =$
 435 $\max_{i \neq k} |\langle \mathbf{x}_i, \mathbf{x}_k \rangle|$ and suppose $R = R_{\max}/R_{\min}$ is at most an absolute constant. Here we list

the exact conditions for $\eta, \sigma_0, R_{\min}, R_{\max}, p$ required by the proofs in this section.

$$\sigma_0 \leq \gamma (CR_{\max} \sqrt{\log(mn/\delta)})^{-1}, \quad (\text{B.2})$$

$$\eta \leq (CR_{\max}^2/nm)^{-1}, \quad (\text{B.3})$$

$$R_{\min}^2 \geq Cr^{-4}R^2np, \quad (\text{B.4})$$

where C is a large enough constant. By Lemma A.1, we can upper bound β by $2\sqrt{\log(12mn/\delta)} \cdot \sigma_0 R_{\max}$. Then, by (B.2) and (B.4), it is straightforward to verify the following inequality:

$$\beta \leq c\gamma, \quad (\text{B.5})$$

$$\gamma^{-4}R_{\min}^{-2}np \leq c, \quad (\text{B.6})$$

$$\gamma^{-4}R_{\min}^{-2}R^2np \leq c, \quad (\text{B.7})$$

where c is a sufficiently small constant.

Suppose the conditions listed in (B.2) and (B.4) hold, we claim that for any $t \geq 0$ the following property holds.

Lemma B.1. Under the same conditions as Theorem 4.1, for any $t \geq 0$, we have that

$$\sum_{r=1}^m |\rho_{j,r,i}^{(t)}| \geq c_1 \gamma^2 \sum_{r=1}^m |\rho_{j',r,i'}^{(t)}|, \forall j, j' \in \{\pm 1\}, \forall i, i' \in [n], \quad (\text{B.8})$$

where c_1 is a constant.

To prove Lemma B.1, we divide it into two lemmas, each addressing a specific case: $0 \leq t \leq T_1$ (Lemma B.2) when the logit $|\ell_i^{(t)}| = \Theta(1)$, and $t \geq T_1$ (Lemma B.3) when the logit $|\ell_i^{(t)}|$ is smaller than constant order. Here, $T_1 = C'\eta^{-1}nmR_{\max}^{-2}$, and C' is a constant. For each case, we apply different techniques to establish the proof.

Lemma B.2 ($0 \leq t \leq T_1$). Under the same conditions as Theorem 4.1, for any $0 \leq t \leq T_1 = C'\eta^{-1}nmR_{\max}^{-2}$, where C' is a constant, we have that

$$|\rho_{j,r,i}^{(t)}| \geq c_2 \gamma |\rho_{j',r',i'}^{(t)}|, \forall j, j' \in \{\pm 1\}, \forall r, r' \in [m], \forall i, i' \in [n], \quad (\text{B.9})$$

where c_2 is a constant.

Proof of Lemma B.2. In this lemma, we first show that (B.8) hold for $t \leq T_1 = C'\eta^{-1}nmR_{\max}^{-2}$ where $C' = \Theta(1)$ is a constant. Recall from Lemma A.3 that

$$\begin{aligned} \bar{\rho}_{j,r,i}^{(t+1)} &= \bar{\rho}_{j,r,i}^{(t)} - \frac{\eta}{nm} \cdot \ell_i^{(t)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle) \cdot \|\mathbf{x}_i\|_2^2 \cdot \mathbb{1}(y_i = j), \\ \underline{\rho}_{j,r,i}^{(t+1)} &= \underline{\rho}_{j,r,i}^{(t)} + \frac{\eta}{nm} \cdot \ell_i^{(t)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle) \cdot \|\mathbf{x}_i\|_2^2 \cdot \mathbb{1}(y_i = -j), \end{aligned}$$

we can get

$$\bar{\rho}_{-y_i,r,i}^{(t)} - \underline{\rho}_{y_i,r,i}^{(t)} = 0, \quad (\text{B.10})$$

and

$$\bar{\rho}_{y_i,r,i}^{(t+1)} \leq \bar{\rho}_{y_i,r,i}^{(t)} + \frac{\eta}{nm} \cdot \|\mathbf{x}_i\|_2^2 \leq \bar{\rho}_{y_i,r,i}^{(t)} + \frac{\eta R_{\max}^2}{nm}, \quad (\text{B.11})$$

$$|\underline{\rho}_{-y_i,r,i}^{(t+1)}| \leq |\underline{\rho}_{-y_i,r,i}^{(t)}| + \frac{\eta}{nm} \cdot \|\mathbf{x}_i\|_2^2 \leq |\underline{\rho}_{-y_i,r,i}^{(t)}| + \frac{\eta R_{\max}^2}{nm}. \quad (\text{B.12})$$

Therefore, we have $\max_{j,r,i} \{|\bar{\rho}_{j,r,i}^{(t)}|, |\underline{\rho}_{j,r,i}^{(t)}|\} = O(1)$ for any $t \leq T_1$ and hence $\max_i \{F_{+1}(\mathbf{W}_{+1}^{(t)}, \mathbf{x}_i), F_{-1}(\mathbf{W}_{-1}^{(t)}, \mathbf{x}_i)\} = O(1)$ for any $t \leq T_1$. Thus there exists a positive constant \tilde{c} such that $|\ell_i^{(t)}| \geq \tilde{c}$ for any $t \leq T_1$. And it follows for any $j \in \{\pm 1\}, r \in [m], i \in [n]$ that

$$|\rho_{j,r,i}^{(t+1)}| \geq |\rho_{j,r,i}^{(t)}| + \frac{\gamma\eta}{nm} \cdot |\ell_i^{(t)}| \cdot \|\mathbf{x}_i\|_2^2 \geq |\rho_{j,r,i}^{(t)}| + \frac{\tilde{c}\gamma\eta}{nm} \cdot \|\mathbf{x}_i\|_2^2, \forall 0 \leq t \leq T_1,$$

$$|\rho_{j,r,i}^{(t)}| \geq \frac{\tilde{c}\gamma\eta t}{nm} \cdot \|\mathbf{x}_i\|_2^2 \geq \frac{\tilde{c}\gamma R_{\min}^2 t}{nm}, \forall 0 \leq t \leq T_1. \quad (\text{B.13})$$

On the other hand, by (B.10), (B.11) and (B.12), we have for any $j' \in \{\pm 1\}, r' \in [m], i' \in [n]$ that

$$|\rho_{j',r',i'}^{(t)}| \leq \frac{\eta R_{\max}^2 t}{nm}, \forall 0 \leq t \leq T_1. \quad (\text{B.14})$$

Dividing (B.14) by (B.13), we can get for any $j, j' \in \{\pm 1\}, r, r' \in [m], i, i' \in [n]$ that

$$|\rho_{j,r,i}^{(t)}| \geq \frac{\tilde{c}\gamma R_{\min}^2}{R_{\max}^2} |\rho_{j',r',i'}^{(t)}|,$$

which indicates that the first bullet holds for time $t \leq T_1$ as long as $c_2 \leq \tilde{c}R_{\min}^2 R_{\max}^{-2}$. \square

Lemma B.3 ($t \geq T_1$). Let T_1 be defined in Lemma B.2. Under the same conditions as Theorem 4.1, for any $t \geq T_1$, we have that

$$\sum_{r=1}^m |\rho_{j,r,i}^{(t)}| \geq c_3 \gamma^2 \sum_{r=1}^m |\rho_{j',r,i'}^{(t)}|, \forall j, j' \in \{\pm 1\}, \forall i, i' \in [n], \quad (\text{B.15})$$

where $c_3 = \Theta(1)$ is a constant. Moreover, we also have the following increasing rate estimation of

$$|\rho_{y_i,r,i}^{(t)}|, |\rho_{-y_i,r,i}^{(t)}|:$$

$$\bullet \frac{1}{m} \sum_{r=1}^m \rho_{y_i,r,i}^{(t)} \leq c_4^{-1} \log \left(1 + \frac{\eta \|\mathbf{x}_i\|_2^2 c_4 e^{2\beta}}{nm} \cdot t \right),$$

$$\bullet \frac{1}{m} \sum_{r=1}^m |\rho_{-y_i,r,i}^{(t)}| \leq c_5^{-1} \gamma^{-1} \log \left(1 + \frac{\gamma \eta \|\mathbf{x}_i\|_2^2 c_5 e^{2\beta}}{nm} \cdot t \right),$$

$$\bullet \frac{1}{m} \sum_{r=1}^m \rho_{y_i,r,i}^{(t)} \geq c_6^{-1} \log \left(1 + \frac{\gamma \eta \|\mathbf{x}_i\|_2^2 c_6 e^{-(\gamma+1)\beta}}{nm} \cdot t \right),$$

$$\bullet \frac{1}{m} \sum_{r=1}^m |\rho_{-y_i,r,i}^{(t)}| \geq c_6^{-1} \gamma \log \left(1 + \frac{\eta \|\mathbf{x}_i\|_2^2 c_6 e^{-(\gamma+1)\beta}}{nm} \cdot t \right),$$

where c_4, c_5, c_6 are constants.

Proof of Lemma B.3. We prove this lemma by induction. By Lemma B.2, we know that (B.15) holds for time $t = T_1$ as long as $c_3 \leq c_2$. Suppose that there exists $\tilde{t} \geq T_1$ such that (B.15) holds for all time $0 \leq t \leq \tilde{t} - 1$. We aim to prove that they also hold for $t = \tilde{t}$. For any $0 \leq t \leq \tilde{t} - 1$, we have

$$\begin{aligned} F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i) &= \frac{1}{m} \sum_{r=1}^m \sigma(\langle \mathbf{w}_{y_i,r}^{(t)}, \mathbf{x}_i \rangle) \\ &\geq \frac{1}{m} \sum_{r=1}^m \langle \mathbf{w}_{y_i,r}^{(t)}, \mathbf{x}_i \rangle \\ &= \frac{1}{m} \sum_{r=1}^m \left(\langle \mathbf{w}_{y_i,r}^{(0)}, \mathbf{x}_i \rangle + \sum_{i'=1}^n \rho_{y_i,r,i'}^{(t)} \|\mathbf{x}_{i'}\|_2^{-2} \cdot \langle \mathbf{x}_{i'}, \mathbf{x}_i \rangle \right) \\ &\geq \frac{1}{m} \sum_{r=1}^m \left(\rho_{y_i,r,i}^{(t)} - \sum_{i' \neq i} |\rho_{y_i,r,i'}^{(t)}| R_{\min}^{-2} p \right) - \beta \\ &= \frac{1}{m} \sum_{r=1}^m \rho_{y_i,r,i}^{(t)} - \sum_{i' \neq i} \left(\frac{1}{m} \sum_{r=1}^m \rho_{y_i,r,i'}^{(t)} \right) R_{\min}^{-2} p - \beta \\ &\geq \frac{1 - \gamma^{-2} c_3^{-1} R_{\min}^{-2} p n}{m} \sum_{r=1}^m \rho_{y_i,r,i}^{(t)} - \beta, \end{aligned} \quad (\text{B.16})$$

473 where the first inequality is by $\sigma(z) \geq z$; the second equality is by (5.1); the third inequality is
 474 by triangle inequality and the definition of β, p, R_{\min} ; the fourth inequality is by the induction
 475 hypothesis (B.15). Besides, for any $0 \leq t \leq t-1$, we also have the following upper bound of
 476 $F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i)$:

$$\begin{aligned}
 F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i) &= \frac{1}{m} \sum_{r=1}^m \sigma(\langle \mathbf{w}_{y_i,r}^{(t)}, \mathbf{x}_i \rangle) \\
 &= \frac{1}{m} \sum_{r=1}^m \sigma\left(\langle \mathbf{w}_{y_i,r}^{(0)}, \mathbf{x}_i \rangle + \sum_{i'=1}^n \rho_{y_i,r,i'}^{(t)} \|\mathbf{x}_{i'}\|_2^{-2} \cdot \langle \mathbf{x}_{i'}, \mathbf{x}_i \rangle\right) \\
 &\leq \frac{1}{m} \sum_{r=1}^m \sigma\left(\rho_{y_i,r,i}^{(t)} + \sum_{i' \neq i} |\rho_{y_i,r,i'}^{(t)}| R_{\min}^{-2} p + \beta\right) \\
 &= \frac{1}{m} \sum_{r=1}^m \left(\rho_{y_i,r,i}^{(t)} + \sum_{i' \neq i} |\rho_{y_i,r,i'}^{(t)}| R_{\min}^{-2} p + \beta\right) \\
 &= \frac{1}{m} \sum_{r=1}^m \rho_{y_i,r,i}^{(t)} + \sum_{i' \neq i} \left(\frac{1}{m} \sum_{r=1}^m \rho_{y_i,r,i'}^{(t)}\right) R_{\min}^{-2} p + \beta \\
 &\leq \frac{1 + \gamma^{-2} c_3^{-1} R_{\min}^{-2} p n}{m} \sum_{r=1}^m \rho_{y_i,r,i}^{(t)} + \beta,
 \end{aligned} \tag{B.17}$$

477 where the first inequality is by triangle inequality and the definition of β, p, R_{\min} ; the second
 478 inequality is by the induction hypothesis (B.15). On the other hand, for any $0 \leq t \leq t$, we can give
 479 following upper and lower bounds for $F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i)$ by applying similar arguments like (B.16)
 480 and (B.17):

$$\begin{aligned}
 F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i) &\geq \frac{\gamma}{m} \sum_{r=1}^m \langle \mathbf{w}_{-y_i,r}^{(t)}, \mathbf{x}_i \rangle \\
 &\geq \frac{\gamma}{m} \sum_{r=1}^m \left(\rho_{-y_i,r,i}^{(t)} - \sum_{i' \neq i} |\rho_{-y_i,r,i'}^{(t)}| R_{\min}^{-2} p - \beta\right), \\
 &\geq \frac{\gamma(1 + \gamma^{-2} c_3^{-1} R_{\min}^{-2} p n)}{m} \sum_{r=1}^m \rho_{-y_i,r,i}^{(t)} - \gamma\beta,
 \end{aligned} \tag{B.18}$$

481 and

$$\begin{aligned}
 F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i) &\leq \frac{1}{m} \sum_{r=1}^m \sigma\left(\rho_{-y_i,r,i}^{(t)} + \sum_{i' \neq i} |\rho_{-y_i,r,i'}^{(t)}| R_{\min}^{-2} p + \beta\right) \\
 &\leq \frac{1}{m} \sum_{r=1}^m \left[\sigma(\rho_{-y_i,r,i}^{(t)}) + \sigma\left(\sum_{i' \neq i} |\rho_{-y_i,r,i'}^{(t)}| R_{\min}^{-2} p\right) + \sigma(\beta)\right] \\
 &= \frac{\gamma}{m} \sum_{r=1}^m \rho_{-y_i,r,i}^{(t)} + \sum_{i' \neq i} \left(\frac{1}{m} \sum_{r=1}^m |\rho_{-y_i,r,i'}^{(t)}|\right) + \beta \\
 &= \frac{\gamma(1 - \gamma^{-3} c_3^{-1} R_{\min}^{-2} p n)}{m} \sum_{r=1}^m \rho_{-y_i,r,i}^{(t)} + \beta,
 \end{aligned} \tag{B.19}$$

482 where the second inequality is by a property of leaky ReLU function that $\sigma(a+b) \leq \sigma(a) +$
 483 $\sigma(b), \forall a, b \in \mathbb{R}$.

484 Next, we can bound $|\ell_i^{(t)}|$ for $0 \leq t \leq \tilde{t} - 1$:

$$\begin{aligned}
& |\ell_i^{(t)}| \\
&= \frac{1}{1 + \exp\{F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i) - F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i)\}} \\
&\leq \exp\{-F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i) + F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i)\} \\
&\leq \exp\left\{-\frac{1 - \gamma^{-2}c_3^{-1}R_{\min}^{-2}pn}{m} \sum_{r=1}^m \rho_{y_i,r,i}^{(t)} + \frac{\gamma(1 - \gamma^{-3}c_3^{-1}R_{\min}^{-2}pn)}{m} \sum_{r=1}^m \rho_{-y_i,r,i}^{(t)} + 2\beta\right\},
\end{aligned} \tag{B.20}$$

485 where the second inequality is by (B.16) and (B.19). And

$$\begin{aligned}
& |\ell_i^{(t)}| \\
&= \frac{1}{1 + \exp\{F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i) - F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i)\}} \\
&\geq \frac{1}{1 + \exp\left\{\frac{1 + \gamma^{-2}c_3^{-1}R_{\min}^{-2}pn}{m} \sum_{r=1}^m \rho_{y_i,r,i}^{(t)} - \frac{\gamma(1 + \gamma^{-2}c_3^{-1}R_{\min}^{-2}pn)}{m} \sum_{r=1}^m \rho_{-y_i,r,i}^{(t)} + (\gamma + 1)\beta\right\}} \\
&\geq \frac{1}{2} \exp\left\{-\frac{1 + \gamma^{-2}c_3^{-1}R_{\min}^{-2}pn}{m} \sum_{r=1}^m \rho_{y_i,r,i}^{(t)} + \frac{\gamma(1 + \gamma^{-2}c_3^{-1}R_{\min}^{-2}pn)}{m} \sum_{r=1}^m \rho_{-y_i,r,i}^{(t)} - (\gamma + 1)\beta\right\},
\end{aligned} \tag{B.21}$$

486 where the first inequality is by (B.17) and (B.18); the last inequality is by $1/(1 + \exp(z)) \geq$
487 $\exp(-z)/2$ if $z \geq 0$. By (B.20), we can get for $0 \leq t \leq \tilde{t} - 1$ that

$$|\ell_i^{(t)}| \leq \exp\left\{-\frac{1 - \gamma^{-2}c_3^{-1}R_{\min}^{-2}pn}{m} \sum_{r=1}^m \rho_{y_i,r,i}^{(t)} + 2\beta\right\}, \tag{B.22}$$

$$|\ell_i^{(t)}| \leq \exp\left\{\frac{\gamma(1 - \gamma^{-3}c_3^{-1}R_{\min}^{-2}pn)}{m} \sum_{r=1}^m \rho_{-y_i,r,i}^{(t)} + 2\beta\right\}. \tag{B.23}$$

488 By (B.21) and $\gamma\rho_{y_i,r,i}^{(t)} \leq |\rho_{-y_i,r,i}^{(t)}| \leq \gamma^{-4}\rho_{y_i,r,i}^{(t)}$, we can get for $0 \leq t \leq \tilde{t} - 1$ that

$$|\ell_i^{(t)}| \geq \frac{1}{2} \exp\left\{-\frac{2(1 + \gamma^{-2}c_3^{-1}R_{\min}^{-2}pn)}{m} \sum_{r=1}^m \rho_{y_i,r,i}^{(t)} - (\gamma + 1)\beta\right\}, \tag{B.24}$$

$$\begin{aligned}
|\ell_i^{(t)}| &\geq \frac{1}{2} \exp\left\{\frac{(\gamma^{-1} + \gamma)(1 + \gamma^{-2}c_3^{-1}R_{\min}^{-2}pn)}{m} \sum_{r=1}^m \rho_{-y_i,r,i}^{(t)} - (\gamma + 1)\beta\right\} \\
&\geq \frac{1}{2} \exp\left\{\frac{2(1 + \gamma^{-2}c_3^{-1}R_{\min}^{-2}pn)}{\gamma m} \sum_{r=1}^m \rho_{-y_i,r,i}^{(t)} - (\gamma + 1)\beta\right\}.
\end{aligned} \tag{B.25}$$

489 By (5.4), (5.5) and $\sigma' \in [\gamma, 1]$, we have for $0 \leq t \leq \tilde{t} - 1$ that

$$\begin{aligned}
\rho_{y_i,r,i}^{(t+1)} &\leq \rho_{y_i,r,i}^{(t)} + \frac{\eta}{nm} \cdot |\ell_i^{(t)}| \cdot \|\mathbf{x}_i\|_2^2, \\
\rho_{y_i,r,i}^{(t+1)} &\geq \rho_{y_i,r,i}^{(t)} + \frac{\gamma\eta}{nm} \cdot |\ell_i^{(t)}| \cdot \|\mathbf{x}_i\|_2^2, \\
|\rho_{-y_i,r,i}^{(t+1)}| &\leq |\rho_{-y_i,r,i}^{(t)}| + \frac{\eta}{nm} \cdot |\ell_i^{(t)}| \cdot \|\mathbf{x}_i\|_2^2, \\
|\rho_{-y_i,r,i}^{(t+1)}| &\geq |\rho_{-y_i,r,i}^{(t)}| + \frac{\gamma\eta}{nm} \cdot |\ell_i^{(t)}| \cdot \|\mathbf{x}_i\|_2^2.
\end{aligned} \tag{B.26}$$

490 By plugging (B.22), (B.24), (B.23) and (B.25) into (B.26), we have for $0 \leq t \leq \tilde{t} - 1$ that

$$\sum_{r=1}^m \rho_{y_i, r, i}^{(t+1)} \leq \sum_{r=1}^m \rho_{y_i, r, i}^{(t)} + \frac{\eta \|\mathbf{x}_i\|_2^2 e^{2\beta}}{n} \cdot \exp \left\{ - \frac{1 - \gamma^{-2} c_3^{-1} R_{\min}^{-2} p n}{m} \sum_{r=1}^m \rho_{y_i, r, i}^{(t)} \right\}, \quad (\text{B.27})$$

$$\sum_{r=1}^m |\rho_{-y_i, r, i}^{(t+1)}| \leq \sum_{r=1}^m |\rho_{-y_i, r, i}^{(t)}| + \frac{\eta \|\mathbf{x}_i\|_2^2 e^{2\beta}}{n} \cdot \exp \left\{ - \frac{\gamma(1 - \gamma^{-3} c_3^{-1} R_{\min}^{-2} p n)}{m} \sum_{r=1}^m |\rho_{-y_i, r, i}^{(t)}| \right\}, \quad (\text{B.28})$$

$$\sum_{r=1}^m \rho_{y_i, r, i}^{(t+1)} \geq \sum_{r=1}^m \rho_{y_i, r, i}^{(t)} + \frac{\gamma \eta \|\mathbf{x}_i\|_2^2 e^{-(\gamma+1)\beta}}{2n} \cdot \exp \left\{ - \frac{2(1 + \gamma^{-2} c_3^{-1} R_{\min}^{-2} p n)}{m} \sum_{r=1}^m \rho_{y_i, r, i}^{(t)} \right\}, \quad (\text{B.29})$$

$$\sum_{r=1}^m |\rho_{-y_i, r, i}^{(t+1)}| \geq \sum_{r=1}^m |\rho_{-y_i, r, i}^{(t)}| + \frac{\gamma \eta \|\mathbf{x}_i\|_2^2 e^{-(\gamma+1)\beta}}{2n} \cdot \exp \left\{ - \frac{2(1 + \gamma^{-2} c_3^{-1} R_{\min}^{-2} p n)}{\gamma m} \sum_{r=1}^m |\rho_{-y_i, r, i}^{(t)}| \right\}. \quad (\text{B.30})$$

491 By applying Lemma H.1 to (B.27) and taking

$$x_t = \frac{1 - \gamma^{-2} c_3^{-1} R_{\min}^{-2} p n}{m} \sum_{r=1}^m \rho_{y_i, r, i}^{(t)},$$

492 we can get for $0 \leq t \leq \tilde{t}$ that

$$\begin{aligned} \frac{1}{m} \sum_{r=1}^m \rho_{y_i, r, i}^{(t)} &\leq c_4^{-1} \log \left(1 + \frac{\eta \|\mathbf{x}_i\|_2^2 c_4 e^{2\beta}}{n m} \exp \left\{ \frac{\eta \|\mathbf{x}_i\|_2^2 c_4 e^{2\beta}}{n m} \right\} \cdot t \right) \\ &\leq c_4^{-1} \log \left(1 + \frac{\eta \|\mathbf{x}_i\|_2^2 c_4 e^{2\beta}}{n m} \cdot t \right), \end{aligned} \quad (\text{B.31})$$

493 where $c_4 := 1 - \gamma^{-2} c_3^{-1} R_{\min}^{-2} p n$ and the last inequality is by $\eta \leq (C R_{\max}^2 / n m)^{-1}$ and C is a
494 sufficiently large constant.

495 By applying Lemma H.1 to (B.28) and taking

$$x_t = \frac{\gamma(1 - \gamma^{-3} c_3^{-1} R_{\min}^{-2} p n)}{m} \sum_{r=1}^m |\rho_{-y_i, r, i}^{(t)}|,$$

496 we can get for $0 \leq t \leq \tilde{t}$ that

$$\begin{aligned} \frac{1}{m} \sum_{r=1}^m |\rho_{-y_i, r, i}^{(t)}| &\leq c_5^{-1} \gamma^{-1} \log \left(1 + \frac{\gamma \eta \|\mathbf{x}_i\|_2^2 c_5 e^{2\beta}}{n m} \exp \left\{ \frac{\gamma \eta \|\mathbf{x}_i\|_2^2 c_5 e^{2\beta}}{n m} \right\} \cdot t \right) \\ &\leq c_5^{-1} \gamma^{-1} \log \left(1 + \frac{\gamma \eta \|\mathbf{x}_i\|_2^2 c_5 e^{2\beta}}{n m} \cdot t \right), \end{aligned} \quad (\text{B.32})$$

497 where $c_5 := 1 - \gamma^{-3} c_3^{-1} R_{\min}^{-2} p n$ and the last inequality is by $\eta \leq (C R_{\max}^2 / n m)^{-1}$ and C is a
498 sufficiently large constant.

499 By applying Lemma H.2 to (B.29) and taking

$$x_t = \frac{2(1 + \gamma^{-2} c_3^{-1} R_{\min}^{-2} p n)}{m} \sum_{r=1}^m \rho_{y_i, r, i}^{(t)},$$

500 we can get

$$\frac{1}{m} \sum_{r=1}^m \rho_{y_i, r, i}^{(t)} \geq (2c_6)^{-1} \log \left(1 + \frac{\gamma \eta \|\mathbf{x}_i\|_2^2 c_6 e^{-(\gamma+1)\beta}}{n m} \cdot t \right), \quad (\text{B.33})$$

501 where $c_6 := 1 + \gamma^{-2} c_3^{-1} R_{\min}^{-2} p n$.

502 By applying Lemma H.2 to (B.30) and taking

$$x_t = \frac{2(1 + \gamma^{-4}c_3^{-1}R_{\min}^{-2}pn)}{\gamma m} \sum_{r=1}^m |\rho_{-y_i, r, i}^{(t)}|,$$

503 we can get

$$\frac{1}{m} \sum_{r=1}^m |\rho_{-y_i, r, i}^{(t)}| \geq (2c_6)^{-1} \gamma \log \left(1 + \frac{\eta \|\mathbf{x}_i\|_2^2 c_6 e^{-(\gamma+1)\beta}}{nm} \cdot t \right), \quad (\text{B.34})$$

504 where $c_6 := 1 + \gamma^{-2}c_3^{-1}R_{\min}^{-2}pn$.

505 In order to apply Lemma H.4 (requiring $b > a$), we loosen the bounds in (B.31), (B.32), (B.33) and
506 (B.34) as follows:

$$\frac{1}{m} \sum_{r=1}^m \rho_{y_i, r, i}^{(t)} \leq c_4^{-1} \gamma^{-1} \log \left(1 + \frac{\gamma \eta R_{\max}^2 c_5 e^{2\beta}}{nm} \cdot t \right), \forall 0 \leq t \leq \tilde{t}, \quad (\text{B.35})$$

$$\frac{1}{m} \sum_{r=1}^m |\rho_{-y_i, r, i}^{(t)}| \leq c_4^{-1} \gamma^{-1} \log \left(1 + \frac{\gamma \eta R_{\max}^2 c_5 e^{2\beta}}{nm} \cdot t \right), \forall 0 \leq t \leq \tilde{t}, \quad (\text{B.36})$$

$$\frac{1}{m} \sum_{r=1}^m \rho_{y_i, r, i}^{(t)} \geq (2c_6)^{-1} \gamma \log \left(1 + \frac{\eta R_{\min}^2 c_6 e^{-(\gamma+1)\beta}}{nm} \cdot t \right), \forall 0 \leq t \leq \tilde{t}, \quad (\text{B.37})$$

$$\frac{1}{m} \sum_{r=1}^m |\rho_{-y_i, r, i}^{(t)}| \geq (2c_6)^{-1} \gamma \log \left(1 + \frac{\eta R_{\min}^2 c_6 e^{-(\gamma+1)\beta}}{nm} \cdot t \right), \forall 0 \leq t \leq \tilde{t}, \quad (\text{B.38})$$

507 where (B.35) is by Bernoulli's inequality that $1 + \gamma^{-1}x \leq (1+x)^{\gamma^{-1}}$ for every real number $0 \leq x \leq 1$
508 and $x \geq -1$; (B.37) is by Bernoulli's inequality that $1 + \gamma x \geq (1+x)^\gamma$ for every real number
509 $0 \leq x \leq 1$ and $x \geq -1$. If $R_{\min}^2 c_6 e^{-(\gamma+1)\beta} \geq \gamma R_{\max}^2 c_5 e^{2\beta}$, we have

$$\frac{1}{m} \sum_{r=1}^m |\rho_{j, r, i}^{(t)}| \geq \frac{\gamma^2 (2c_6)^{-1} c_4}{m} \sum_{r=1}^m |\rho_{j', r, i'}^{(t)}|. \quad (\text{B.39})$$

510 If $R_{\min}^2 c_6 e^{-(\gamma+1)\beta} < \gamma R_{\max}^2 c_5 e^{2\beta}$, by Lemma H.4, we have

$$\begin{aligned} & \frac{\min\{\frac{1}{m} \sum_{r=1}^m \rho_{y_i, r, i}^{(t)}, \frac{1}{m} \sum_{r=1}^m |\rho_{-y_i, r, i}^{(t)}|\}}{\max\{\frac{1}{m} \sum_{r=1}^m \rho_{y_{i'}, r, i'}^{(t)}, \frac{1}{m} \sum_{r=1}^m |\rho_{-y_{i'}, r, i'}^{(t)}|\}} \\ & \geq \gamma^2 (2c_6)^{-1} c_4 \cdot \frac{\log \left(1 + \frac{\eta R_{\min}^2 c_6 e^{-(\gamma+1)\beta}}{nm} \cdot t \right)}{\log \left(1 + \frac{\gamma \eta R_{\max}^2 c_5 e^{2\beta}}{nm} \cdot t \right)} \\ & \geq \gamma^2 (2c_6)^{-1} c_4 \cdot \frac{\log \left(1 + \frac{\eta R_{\min}^2 c_6 e^{-(\gamma+1)\beta}}{nm} \cdot T_1 \right)}{\log \left(1 + \frac{\gamma \eta R_{\max}^2 c_5 e^{2\beta}}{nm} \cdot T_1 \right)} \\ & \geq \gamma^2 (2c_6)^{-1} c_4 \cdot \frac{\log(1 + R^{-2} c_6 e^{-(\gamma+1)\beta} C')}{\log(1 + \gamma c_5 e^{2\beta} C')}. \end{aligned}$$

511 Therefore, we can get for $0 \leq t \leq \tilde{t}$ that

$$\sum_{r=1}^m |\rho_{j, r, i}^{(t)}| \geq \gamma^2 c_3 \sum_{r=1}^m |\rho_{j', r, i'}^{(t)}|, \forall j, j' \in \{\pm 1\}, \forall i, i' \in [n], \quad (\text{B.40})$$

512 as long as

$$c_3 \leq (2c_6)^{-1} c_4 \cdot \min \left\{ 1, \frac{\log(1 + R^{-2} c_6 e^{-(\gamma+1)\beta} C')}{\log(1 + \gamma c_5 e^{2\beta} C')} \right\}.$$

513 This condition holds under the following conditions:

$$\gamma^{-3}c_3^{-1}R_{\min}^{-2}pn \leq \frac{1}{2} \implies c_4, c_5 \geq \frac{1}{2}, c_6 \leq \frac{3}{2},$$

$$c_3 = \frac{1}{6} \min \left\{ 1, \frac{\log(1 + R^{-2}e^{-(\gamma+1)\beta}C')}{\log(1 + \gamma e^{2\beta}C')} \right\}.$$

514 This implies that induction hypothesis (B.15) holds for $t = \tilde{t}$. \square

515 **Lemma B.4** (Implication of Lemma B.1). Under the same condition as Theorem 4.1, if (B.8) hold
516 for time t , then we have that

$$|\rho_{j,r,i}^{(t)}| \geq c_1 \gamma^4 |\rho_{j',r',i'}^{(t)}|,$$

517 where c_1 is the same constant as defined in Lemma B.1.

518 *Proof of Lemma B.4.* By $\sigma' \in [\gamma, 1]$, (5.4) and (5.5), we have

$$|\rho_{j,r,i}^{(t+1)}| \geq |\rho_{j,r,i}^{(t)}| + \frac{\gamma\eta}{nm} \cdot |\ell_i^{(t)}| \cdot \|\mathbf{x}_i\|_2^2, \forall j \in \{\pm 1\}, \forall r \in [m], \forall i \in [n],$$

$$|\rho_{j,r,i}^{(t+1)}| \leq |\rho_{j,r,i}^{(t)}| + \frac{\eta}{nm} \cdot |\ell_i^{(t)}| \cdot \|\mathbf{x}_i\|_2^2, \forall j \in \{\pm 1\}, \forall r \in [m], \forall i \in [n].$$

519 Thus, we have

$$|\rho_{j,r,i}^{(t)}| \geq \frac{\gamma\eta\|\mathbf{x}_i\|_2^2}{nm} \cdot \sum_{s=1}^{t-1} |\ell_i^{(s)}|, \forall j \in \{\pm 1\}, \forall r \in [m], \forall i \in [n],$$

$$|\rho_{j,r,i}^{(t)}| \leq \frac{\eta\|\mathbf{x}_i\|_2^2}{nm} \cdot \sum_{s=1}^{t-1} |\ell_i^{(s)}|, \forall j \in \{\pm 1\}, \forall r \in [m], \forall i \in [n].$$

520 Therefore, $|\rho_{j,r,i}^{(t)}| \geq \gamma |\rho_{j',r',i}^{(t)}|$ for any $j, j' \in \{\pm 1\}$, $r', r \in [m]$ and $i \in [n]$, and hence

$$m |\rho_{j,r,i}^{(t)}| \geq \gamma \sum_{r=1}^m |\rho_{j,r,i}^{(t)}|,$$

$$\sum_{r=1}^m |\rho_{j',r,i'}^{(t)}| \geq m \gamma |\rho_{j',r',i'}^{(t)}|. \quad (\text{B.41})$$

521 Plugging (B.41) back into (B.8) completes the proof. \square

522 **Lemma B.5.** Let T_1 be defined in Lemma B.2. Every neuron will never change its activation pattern
523 after time T_1 , that is,

$$\text{sign}(\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle) = \text{sign}(\langle \mathbf{w}_{j,r}^{(T_1)}, \mathbf{x}_i \rangle),$$

524 for any $t \geq T_1$, $j \in \{\pm 1\}$ and $r \in [m]$. Moreover, it holds that

$$\text{sign}(\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle) = jy_i, \quad (\text{B.42})$$

525 for any $t \geq T_1$, $j \in \{\pm 1\}$ and $r \in [m]$.

526 *Proof of Lemma B.5.* For $j = y_i$ and $t \geq 0$, we have $\rho_{j,r,i}^{(t)} = 0$, and so

$$\begin{aligned} \langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle &= \langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle + \sum_{i'=1}^n \rho_{j,r,i'}^{(t)} \|\mathbf{x}_{i'}\|_2^{-2} \cdot \langle \mathbf{x}_{i'}, \mathbf{x}_i \rangle \\ &= \langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle + \bar{\rho}_{j,r,i}^{(t)} + \sum_{i' \neq i} \rho_{j,r,i'}^{(t)} \|\mathbf{x}_{i'}\|_2^{-2} \cdot \langle \mathbf{x}_{i'}, \mathbf{x}_i \rangle \\ &\geq \bar{\rho}_{j,r,i}^{(t)} - \sum_{i' \neq i} |\rho_{j,r,i'}^{(t)}| R_{\min}^{-2} p - \beta \end{aligned}$$

$$\begin{aligned} &\geq \bar{\rho}_{j,r,i}^{(t)} - \gamma^{-4} c_1^{-1} \bar{\rho}_{j,r,i}^{(t)} R_{\min}^{-2} p n - \beta \\ &= (1 - \gamma^{-4} c_1^{-1} R_{\min}^{-2} p n) \cdot \bar{\rho}_{j,r,i}^{(t)} - \beta, \end{aligned}$$

527 where the first inequality is by triangle inequality; the second inequality is by $|\rho_{j,r,i'}^{(t)}| \leq \gamma^{-4} c_1^{-1} \bar{\rho}_{y_i,r,i}^{(t)}$
 528 from Lemma B.1 and Lemma B.4.

529 By (B.13), we have for $t \geq T_1$ that

$$\bar{\rho}_{y_i,r,i}^{(t)} \geq \frac{\tilde{c} \gamma \eta R_{\min}^2 T_1}{n m} = C' \tilde{c} \gamma R_{\min}^2 R_{\max}^{-2}. \quad (\text{B.43})$$

530 Therefore, by (B.5), (B.6) and (B.43), we know that

$$(1 - \gamma^{-1} c_1^{-4} R_{\min}^{-2} p n) \cdot \bar{\rho}_{y_i,r,i}^{(t)} > \beta, \forall r \in [m], i \in [n].$$

531 and thus $\text{sign}(\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle) = 1$ for any $r \in [m], i \in [n], j = y_i$.

532 For $j \neq y_i$ and any $t \geq 0$, we have $\bar{\rho}_{j,r,i}^{(t)} = 0$, and so

$$\begin{aligned} \langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle &= \langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle + \sum_{i'=1}^n \rho_{j,r,i'}^{(t)} \|\mathbf{x}_{i'}\|_2^{-2} \cdot \langle \mathbf{x}_{i'}, \mathbf{x}_i \rangle \\ &= \underline{\rho}_{j,r,i}^{(t)} + \sum_{i' \neq i} \rho_{j,r,i'}^{(t)} \|\mathbf{x}_{i'}\|_2^{-2} \cdot \langle \mathbf{x}_{i'}, \mathbf{x}_i \rangle \\ &\leq \underline{\rho}_{j,r,i}^{(t)} + \sum_{i' \neq i} |\rho_{j,r,i'}^{(t)}| R_{\min}^{-2} p + \beta \\ &\leq \underline{\rho}_{j,r,i}^{(t)} - \gamma^{-1} c_2^{-1} \underline{\rho}_{j,r,i}^{(t)} R_{\min}^{-2} p n - \beta \\ &= (1 - \gamma^{-1} c_2^{-1} R_{\min}^{-2} p n) \underline{\rho}_{j,r,i}^{(t)} - \beta, \end{aligned}$$

533 where the first inequality is by triangle inequality; the second inequality is by $|\rho_{j,r,i'}^{(t)}| \leq$
 534 $\gamma^{-4} c_1^{-1} |\underline{\rho}_{-y_i,r,i}^{(t)}|$ from Lemma B.1 and Lemma B.4.

535 By (B.13), we have

$$|\underline{\rho}_{-y_i,r,i}^{(t)}| \geq \frac{\tilde{c} \gamma \eta R_{\min}^2 T_1}{n m} = C' \tilde{c} \gamma R_{\min}^2 R_{\max}^{-2}. \quad (\text{B.44})$$

536 Therefore, by (B.5), (B.6) and (B.44), we know that

$$(1 - \gamma^{-4} c_1^{-1} R_{\min}^{-2} p n) \cdot |\underline{\rho}_{-y_i,r,i}^{(t)}| > \beta, \forall r \in [m], i \in [n],$$

537 and thus $\text{sign}(\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle) = -1$ for $j \neq y_i$, which completes the proof. \square

538 C Stable Rank of Leaky ReLU Network

539 In this section, we consider the properties of stable rank of the weight matrix $\mathbf{W}^{(t)}$ found by gradient
 540 descent at time t , defined as $\|\mathbf{W}^{(t)}\|_F^2 / \|\mathbf{W}^{(t)}\|_2^2$, the square of the ratio of the Frobenius norm to
 541 the spectral norm of $\mathbf{W}^{(t)}$. Given Lemma B.5, we have following coefficient update rule for $t \geq T_1$
 542 where T_1 is defined in Lemma B.2:

$$\bar{\rho}_{y_i,r,i}^{(t+1)} = \bar{\rho}_{y_i,r,i}^{(t)} + \frac{\eta}{n m} \cdot |\ell_i^{(t)}| \cdot \|\mathbf{x}_i\|_2^2, \quad (\text{C.1})$$

$$\underline{\rho}_{-y_i,r,i}^{(t+1)} = \underline{\rho}_{-y_i,r,i}^{(t)} - \frac{\gamma \eta}{n m} \cdot |\ell_i^{(t)}| \cdot \|\mathbf{x}_i\|_2^2, \quad (\text{C.2})$$

543 where

$$|\ell_i^{(t)}| = \frac{1}{1 + \exp\{F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i) - F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i)\}}.$$

544 Based on (C.1) and (C.2), we first introduce the following helpful lemmas.

545 **Lemma C.1.** Let T_1 be defined in Lemma B.3. For any $r, r' \in [m]$, $i \in [n]$ and $t \leq T_1$,

$$|\bar{\rho}_{y_i, r, i}^{(t)} - \bar{\rho}_{y_i, r', i}^{(t)}| \leq C', |\rho_{-y_i, r, i}^{(t)} - \rho_{-y_i, r', i}^{(t)}| \leq C'. \quad (\text{C.3})$$

546 *Proof of Lemma C.1.* By (B.14), we can get

$$|\rho_{j, r, i}^{(t)}| \leq \frac{\eta R_{\max}^2 T_1}{nm} = C',$$

547 for $t \leq T_1$. Notice that

$$\begin{aligned} |\bar{\rho}_{y_i, r, i}^{(t)} - \bar{\rho}_{y_i, r', i}^{(t)}| &\leq \max\{|\bar{\rho}_{y_i, r, i}^{(t)}|, |\bar{\rho}_{y_i, r', i}^{(t)}|\}, \\ |\rho_{-y_i, r, i}^{(t)} - \rho_{-y_i, r', i}^{(t)}| &\leq \max\{|\rho_{-y_i, r, i}^{(t)}|, |\rho_{-y_i, r', i}^{(t)}|\}, \end{aligned}$$

548 which completes the proof. \square

549 **Lemma C.2.** Let T_1 be defined in Lemma B.3. For any $r, r' \in [m]$, $i \in [n]$ and $t \geq T_1$,

$$|\bar{\rho}_{y_i, r, i}^{(t)} - \bar{\rho}_{y_i, r', i}^{(t)}| \leq C', |\rho_{-y_i, r, i}^{(t)} - \rho_{-y_i, r', i}^{(t)}| \leq C'.$$

550 *Proof of Lemma C.2.* By (C.1) and (C.2), we can get for any $r \in [m]$, $i \in [n]$ and $t \geq T_1$ that

$$\begin{aligned} \bar{\rho}_{y_i, r, i}^{(t)} &= \bar{\rho}_{y_i, r, i}^{(T_1)} + \frac{\eta}{nm} \sum_{s=T_1}^{t-1} |\ell_i^{(s)}| \cdot \|\mathbf{x}_i\|_2^2, \\ \rho_{-y_i, r, i}^{(t)} &= \rho_{-y_i, r, i}^{(T_1)} + \frac{\eta}{nm} \sum_{s=T_1}^{t-1} |\ell_i^{(s)}| \cdot \|\mathbf{x}_i\|_2^2. \end{aligned}$$

551 Since $\bar{\rho}_{y_i, r, i}^{(t)}$, $\bar{\rho}_{y_i, r', i}^{(t)}$ possess the same increment and $\rho_{-y_i, r, i}^{(t)}$, $\rho_{-y_i, r', i}^{(t)}$ possess the same increment,
552 we have

$$\begin{aligned} \bar{\rho}_{y_i, r, i}^{(t)} - \bar{\rho}_{y_i, r', i}^{(t)} &= \bar{\rho}_{y_i, r, i}^{(T_1)} - \bar{\rho}_{y_i, r', i}^{(T_1)}, \\ \rho_{-y_i, r, i}^{(t)} - \rho_{-y_i, r', i}^{(t)} &= \rho_{-y_i, r, i}^{(T_1)} - \rho_{-y_i, r', i}^{(T_1)}. \end{aligned}$$

553 Notice that

$$\max_{i, r, r'}\{|\bar{\rho}_{y_i, r, i}^{(T_1)} - \bar{\rho}_{y_i, r', i}^{(T_1)}|, |\rho_{-y_i, r, i}^{(T_1)} - \rho_{-y_i, r', i}^{(T_1)}|\} \leq \max_{i, r, r'}\{|\bar{\rho}_{y_i, r, i}^{(T_1)}|, |\rho_{-y_i, r, i}^{(T_1)}|\} \leq C',$$

554 which completes the proof. \square

555 Now we are ready to prove the second bullet of Theorem 4.1.

556 **Lemma C.3.** Throughout the gradient descent trajectory, the stable rank of the weights \mathbf{W}_j satisfies,

$$\lim_{t \rightarrow \infty} \frac{\|\mathbf{W}_j\|_F^2}{\|\mathbf{W}_j\|_2^2} = 1, \forall j \in \{\pm 1\},$$

557 with a decreasing rate of $O(1/\log(t))$.

558 *Proof of Lemma C.3.* By Definition 5.1, we have

$$\mathbf{w}_{j, r}^{(t)} = \mathbf{w}_{j, r}^{(0)} + \underbrace{\sum_{i=1}^n \rho_{j, r, i}^{(t)} \cdot \|\mathbf{x}_i\|_2^{-2} \cdot \mathbf{x}_i}_{:= \mathbf{v}_{j, r}^{(t)}}.$$

559 We first show that $\|\mathbf{v}_{j,r}^{(t)}\|_2 = \Theta(\log t)$.

$$\begin{aligned}
\|\mathbf{v}_{j,r}^{(t)}\|_2^2 &= \left(\sum_{i=1}^n \rho_{j,r,i}^{(t)} \cdot \|\mathbf{x}_i\|_2^{-2} \cdot \mathbf{x}_i \right)^2 \\
&= \sum_{i=1}^n (\rho_{j,r,i}^{(t)})^2 \cdot \|\mathbf{x}_i\|_2^{-2} + \sum_{i \neq i'} \rho_{j,r,i}^{(t)} \rho_{j,r,i'}^{(t)} \|\mathbf{x}_i\|_2^{-2} \|\mathbf{x}_{i'}\|_2^{-2} \langle \mathbf{x}_i, \mathbf{x}_{i'} \rangle \\
&\geq \sum_{i=1}^n (\rho_{j,r,i}^{(t)})^2 \cdot R_{\max}^{-2} - \sum_{i \neq i'} |\rho_{j,r,i}^{(t)}| |\rho_{j,r,i'}^{(t)}| \cdot R_{\min}^{-4} p \\
&\geq R_{\max}^{-2} (1 - R^2 R_{\min}^{-2} c_1^{-1} \gamma^{-4} n p) \sum_{i=1}^n (\rho_{j,r,i}^{(t)})^2 \\
&= \Theta(n R_{\max}^{-2} \log^2(t)),
\end{aligned}$$

560 where the second last inequality is by triangle inequality; the last inequality is by $|\rho_{j,r,i}^{(t)}| \leq$
561 $\gamma^{-4} c_1^{-1} |\rho_{j,r,i'}^{(t)}|$ from Lemma B.1 and Lemma B.4.

562 By the definition of $\mathbf{v}_{j,r}^{(t)}$, we have

$$\begin{aligned}
\|\mathbf{v}_{j,r}^{(t)} - \mathbf{v}_{j,r'}^{(t)}\|_2^2 &= \left\| \sum_{i=1}^n \rho_{j,r,i}^{(t)} \cdot \|\mathbf{x}_i\|_2^{-2} \cdot \mathbf{x}_i - \sum_{i=1}^n \rho_{j,r',i}^{(t)} \cdot \|\mathbf{x}_i\|_2^{-2} \cdot \mathbf{x}_i \right\|_2^2 \\
&= \left\| \sum_{i=1}^n (\rho_{j,r,i}^{(t)} - \rho_{j,r',i}^{(t)}) \cdot \|\mathbf{x}_i\|_2^{-2} \cdot \mathbf{x}_i \right\|_2^2 \\
&= \sum_{i=1}^n (\rho_{j,r,i}^{(t)} - \rho_{j,r',i}^{(t)})^2 \cdot \|\mathbf{x}_i\|_2^{-2} + \sum_{i \neq i'} (\rho_{j,r,i}^{(t)} - \rho_{j,r',i}^{(t)}) (\rho_{j,r,i'}^{(t)} - \rho_{j,r',i'}^{(t)}) \frac{\langle \mathbf{x}_i, \mathbf{x}_{i'} \rangle}{\|\mathbf{x}_i\|_2^2 \|\mathbf{x}_{i'}\|_2^2} \\
&\leq (C')^2 n R_{\min}^{-2} + (C')^2 n^2 R_{\min}^{-4} p \\
&\leq 2(C')^2 n R_{\min}^{-2},
\end{aligned}$$

563 where the first inequality is by Lemma C.1 and Lemma C.2.

564 Now, we are ready to estimate the stable rank of $\mathbf{W}^{(t)}$. On the one hand, for $\|\mathbf{W}_j^{(t)}\|_F^2$, we have

$$\begin{aligned}
\|\mathbf{W}_j^{(t)}\|_F^2 &= \sum_r \|\mathbf{w}_{j,r}^{(t)}\|_2^2 \\
&= \sum_r \|\mathbf{w}_{j,r}^{(0)} + \mathbf{v}_{j,r}^{(t)}\|_2^2 \\
&= \sum_r \|\mathbf{w}_{j,r}^{(0)}\|_2^2 + \|\mathbf{v}_{j,r}^{(t)}\|_2^2 + 2\langle \mathbf{w}_{j,r}^{(0)}, \mathbf{v}_{j,r}^{(t)} \rangle \\
&\leq \sum_r \|\mathbf{w}_{j,r}^{(0)}\|_2^2 + (\|\mathbf{v}_{j,1}^{(t)}\|_2 + \|\mathbf{v}_{j,r}^{(t)} - \mathbf{v}_{j,1}^{(t)}\|_2)^2 + 2\|\mathbf{w}_{j,r}^{(0)}\|_2 (\|\mathbf{v}_{j,1}^{(t)}\|_2 + \|\mathbf{v}_{j,r}^{(t)} - \mathbf{v}_{j,1}^{(t)}\|_2) \\
&= m\|\mathbf{v}_{j,1}^{(t)}\|_2^2 + 2\left(\sum_r \|\mathbf{v}_{j,r}^{(t)} - \mathbf{v}_{j,1}^{(t)}\|_2 + \|\mathbf{w}_{j,r}^{(0)}\|_2 \right) \|\mathbf{v}_{j,1}^{(t)}\|_2 \\
&\quad + \left(\sum_r \|\mathbf{w}_{j,r}^{(0)}\|_2^2 + \|\mathbf{v}_{j,r}^{(t)} - \mathbf{v}_{j,1}^{(t)}\|_2^2 + 2\|\mathbf{w}_{j,r}^{(0)}\|_2 \|\mathbf{v}_{j,r}^{(t)} - \mathbf{v}_{j,1}^{(t)}\|_2 \right) \\
&\leq m\|\mathbf{v}_{j,1}^{(t)}\|_2^2 + mC_1\|\mathbf{v}_{j,1}^{(t)}\|_2 + mC_2.
\end{aligned}$$

565 where the first inequality is by triangle inequality and Cauchy inequality; the last inequality is by
566 Lemma C.1, Lemma C.2 and taking

$$C_1 = 3(\sigma_0 \sqrt{d} + C' \sqrt{n} R_{\min}^{-1}),$$

$$C_2 = 2(\sigma_0\sqrt{d} + C'\sqrt{n}R_{\min}^{-1})^2.$$

567 On the other hand, for $\|\mathbf{W}_j^{(t)}\|_2^2$, we have

$$\begin{aligned} \|\mathbf{W}_j^{(t)}\|_2^2 &= \max_{\mathbf{x} \in S^{d-1}} \|\mathbf{W}_j^{(0)}\mathbf{x} + \mathbf{V}_j^{(t)}\mathbf{x}\|_2^2 \\ &= \max_{\mathbf{x} \in S^{d-1}} \|\mathbf{W}_j^{(0)}\mathbf{x}\|_2^2 + \|\mathbf{V}_j^{(t)}\mathbf{x}\|_2^2 + 2\langle \mathbf{W}_j^{(0)}\mathbf{x}, \mathbf{V}_j^{(t)}\mathbf{x} \rangle \\ &= \max_{\mathbf{x} \in S^{d-1}} \sum_r \langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x} \rangle^2 + \sum_r \langle \mathbf{v}_{j,r}^{(t)}, \mathbf{x} \rangle^2 + \sum_r \langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x} \rangle \cdot \langle \mathbf{v}_{j,r}^{(t)}, \mathbf{x} \rangle \\ &\geq \sum_r \left\langle \mathbf{w}_{j,r}^{(0)}, \frac{\mathbf{v}_{j,1}^{(t)}}{\|\mathbf{v}_{j,1}^{(t)}\|_2} \right\rangle^2 + \sum_r \left\langle \mathbf{v}_{j,r}^{(t)}, \frac{\mathbf{v}_{j,1}^{(t)}}{\|\mathbf{v}_{j,1}^{(t)}\|_2} \right\rangle^2 + \sum_r \left\langle \mathbf{w}_{j,r}^{(0)}, \frac{\mathbf{v}_{j,1}^{(t)}}{\|\mathbf{v}_{j,1}^{(t)}\|_2} \right\rangle \cdot \left\langle \mathbf{v}_{j,r}^{(t)}, \frac{\mathbf{v}_{j,1}^{(t)}}{\|\mathbf{v}_{j,1}^{(t)}\|_2} \right\rangle \\ &\geq \sum_r \left\langle \mathbf{v}_{j,r}^{(t)}, \frac{\mathbf{v}_{j,1}^{(t)}}{\|\mathbf{v}_{j,1}^{(t)}\|_2} \right\rangle^2 - \sum_r \|\mathbf{w}_{j,r}^{(0)}\|_2^2 - \sum_r \|\mathbf{w}_{j,r}^{(0)}\|_2 \|\mathbf{v}_{j,r}^{(t)}\|_2 \\ &\geq m\|\mathbf{v}_{j,1}^{(t)}\|_2^2 + 2 \sum_r \|\mathbf{v}_{j,1}^{(t)}\|_2 \cdot \left\langle \mathbf{v}_{j,r}^{(t)} - \mathbf{v}_{j,1}^{(t)}, \frac{\mathbf{v}_{j,1}^{(t)}}{\|\mathbf{v}_{j,1}^{(t)}\|_2} \right\rangle + \sum_r \left\langle \mathbf{v}_{j,r}^{(t)} - \mathbf{v}_{j,1}^{(t)}, \frac{\mathbf{v}_{j,1}^{(t)}}{\|\mathbf{v}_{j,1}^{(t)}\|_2} \right\rangle^2 \\ &\quad - \sum_r \|\mathbf{w}_{j,r}^{(0)}\|_2^2 - \sum_r \|\mathbf{w}_{j,r}^{(0)}\|_2 (\|\mathbf{v}_{j,1}^{(t)}\|_2 + \|\mathbf{v}_{j,r}^{(t)} - \mathbf{v}_{j,1}^{(t)}\|_2) \\ &\geq m\|\mathbf{v}_{j,1}^{(t)}\|_2^2 - \left(\sum_r 2\|\mathbf{v}_{j,r}^{(t)} - \mathbf{v}_{j,1}^{(t)}\|_2 + \|\mathbf{w}_{j,r}^{(0)}\|_2 \right) \cdot \|\mathbf{v}_{j,1}^{(t)}\|_2 \\ &\quad - \left(\sum_r \|\mathbf{w}_{j,r}^{(0)}\|_2^2 + \|\mathbf{w}_{j,r}^{(0)}\|_2 \|\mathbf{v}_{j,r}^{(t)} - \mathbf{v}_{j,1}^{(t)}\|_2 \right) \\ &\geq m\|\mathbf{v}_{j,1}^{(t)}\|_2^2 - mC_3\|\mathbf{v}_{j,1}^{(t)}\|_2 - mC_4 \end{aligned}$$

568 where the first inequality is by taking $\mathbf{x} = \mathbf{v}_{j,1}^{(t)}/\|\mathbf{v}_{j,1}^{(t)}\|_2$; the second inequality is by Cauchy
569 inequality; the third inequality by breaking $\mathbf{v}_{j,r}^{(t)}$ down into $\mathbf{v}_{j,1}^{(t)} + \mathbf{v}_{j,r}^{(t)} - \mathbf{v}_{j,1}^{(t)}$ and then expanding the
570 first term as well as applying triangle inequality to the last term; the fourth inequality is by Cauchy
571 inequality; the last inequality is by Lemma C.1, Lemma C.2 and taking

$$\begin{aligned} C_3 &= 1.5\sigma_0\sqrt{d} + 3C'\sqrt{n}R_{\min}^{-1}, \\ C_4 &= 1.5\sigma_0^2d + 3C'\sigma_0\sqrt{d}\sqrt{n}R_{\min}^{-1}. \end{aligned}$$

572 By leverage the upper bound of $\|\mathbf{W}_j^{(t)}\|_F^2$ as well as the lower bound of $\|\mathbf{W}_j^{(t)}\|_2^2$, we can get

$$\frac{\|\mathbf{W}_j^{(t)}\|_F^2}{\|\mathbf{W}_j^{(t)}\|_2^2} \leq \frac{\|\mathbf{v}_{j,1}^{(t)}\|_2^2 + C_1\|\mathbf{v}_{j,1}^{(t)}\|_2 + C_2}{\|\mathbf{v}_{j,1}^{(t)}\|_2^2 - C_3\|\mathbf{v}_{j,1}^{(t)}\|_2 - C_4}.$$

573 Since $\|\mathbf{W}_j^{(t)}\|_F^2/\|\mathbf{W}_j^{(t)}\|_2^2 \geq 1$, $\|\mathbf{v}_{j,1}^{(t)}\|_2 = \Theta(\log t)$ and C_1, C_2, C_3, C_4 are constants, it follow that

$$\lim_{t \rightarrow \infty} \frac{\|\mathbf{W}_j^{(t)}\|_F^2}{\|\mathbf{W}_j^{(t)}\|_2^2} = 1,$$

574 and

$$\begin{aligned} \frac{\|\mathbf{W}_j^{(t)}\|_F^2}{\|\mathbf{W}_j^{(t)}\|_2^2} - 1 &\leq \frac{(C_1 + C_3)\|\mathbf{v}_{j,1}^{(t)}\|_2 + (C_2 + C_4)}{\|\mathbf{v}_{j,1}^{(t)}\|_2^2 - C_3\|\mathbf{v}_{j,1}^{(t)}\|_2 - C_4} \\ &\preceq \frac{C_1 + C_3}{\|\mathbf{v}_{j,1}^{(t)}\|_2} \end{aligned}$$

$$\begin{aligned}
&\leq \frac{6(C'\sqrt{n}R_{\min}^{-1} + \sigma_0\sqrt{d})}{\|\mathbf{v}_{j,1}^{(t)}\|_2} \\
&= \Theta\left(\frac{\sqrt{n}R_{\min}^{-1} + \sigma_0\sqrt{d}}{\sqrt{n}R_{\max}^{-1}\log(t)}\right) = \Theta\left(\frac{1 + \sigma_0\sqrt{d/n}R_{\max}}{\log(t)}\right),
\end{aligned}$$

575 which completes the proof. \square

576 D Loss Convergence of Leaky ReLU Network

577 Given Lemma B.5 about the activation pattern after time T_1 , we can now establish a new lemma that
578 will aid us in proving the convergence rate later.

579 **Lemma D.1.** Let T_1 be defined in Lemma B.3. For $t \geq T_1$, it holds that

$$\begin{aligned}
\bar{\rho}_{y_i,r,i}^{(t)} - \bar{\rho}_{y_i,r,i}^{(T_1)} &= \bar{\rho}_{y_i,r',i}^{(t)} - \bar{\rho}_{y_i,r',i}^{(T_1)}, \\
\rho_{-y_i,r,i}^{(t)} - \rho_{-y_i,r,i}^{(T_1)} &= \rho_{-y_i,r',i}^{(t)} - \rho_{-y_i,r',i}^{(T_1)}, \\
\bar{\rho}_{y_i,r,i}^{(t)} - \bar{\rho}_{y_i,r,i}^{(T_1)} &= (|\rho_{-y_i,r',i}^{(t)}| - |\rho_{-y_i,r',i}^{(T_1)}|)/\gamma,
\end{aligned}$$

580 for any $i \in [n]$ and $r, r' \in [m]$.

581 *Proof of Lemma D.1.* By Lemma B.3 about the activation pattern after time T_1 , we can get

$$\bar{\rho}_{y_i,r,i}^{(t+1)} = \bar{\rho}_{y_i,r,i}^{(t)} + \frac{\eta}{nm} \cdot |\ell_i^{(t)}| \cdot \|\mathbf{x}_i\|_2^2, \quad (\text{D.1})$$

$$\rho_{-y_i,r,i}^{(t+1)} = \rho_{-y_i,r,i}^{(t)} - \frac{\gamma\eta}{nm} \cdot |\ell_i^{(t)}| \cdot \|\mathbf{x}_i\|_2^2, \quad (\text{D.2})$$

582 for $t \geq T_1$. Recursively using (D.1) and (D.2) $t - T_1$ times, we can get

$$\begin{aligned}
\bar{\rho}_{y_i,r,i}^{(t)} - \bar{\rho}_{y_i,r,i}^{(T_1)} &= \frac{\eta\|\mathbf{x}_i\|_2^2}{nm} \sum_{s=T_1}^{t-1} |\ell_i^{(s)}|, \\
|\rho_{-y_i,r,i}^{(t)}| - |\rho_{-y_i,r,i}^{(T_1)}| &= \frac{\gamma\eta\|\mathbf{x}_i\|_2^2}{nm} \sum_{s=T_1}^{t-1} |\ell_i^{(s)}|.
\end{aligned}$$

583 This indicates that for different $r, r' \in [m]$, $\bar{\rho}_{y_i,r,i}^{(t)} - \bar{\rho}_{y_i,r,i}^{(T_1)}$ and $\bar{\rho}_{y_i,r',i}^{(t)} - \bar{\rho}_{y_i,r',i}^{(T_1)}$ are the same, whereas
584 $\gamma(|\rho_{-y_i,r,i}^{(t)}| - |\rho_{-y_i,r,i}^{(T_1)}|)$ and $\bar{\rho}_{y_i,r',i}^{(t)} - \bar{\rho}_{y_i,r',i}^{(T_1)}$ are the same, which completes the proof. \square

585 By leveraging Lemma B.5 and Lemma D.1, we can now give an upper bound for the convergence
586 rate.

587 **Lemma D.2.** Let T be defined in Lemma B.3. For leaky ReLU neural network defined in (3.1), for
588 any $t \geq T$, we have

$$L_S(\mathbf{W}^{(t)}) \leq \frac{c}{n} \sum_{i=1}^n \left(1 + \frac{\eta\|\mathbf{x}_i\|_2^2 e^{C_i} (1 + \gamma^2) (1 + c_1^{-1} \gamma^{-2} R_{\min}^{-2} pn)}{2nm} \cdot (t - T) \right)^{-\frac{1 - c_1^{-1} \gamma^{-2} R_{\min}^{-2} pn}{1 + c_1^{-1} \gamma^{-2} R_{\min}^{-2} pn}},$$

589 where c_1 is the same constant as Lemma B.1 and c, C_i are constants. This indicates that the training
590 loss will converge with rate $O(t^{-\alpha})$ where $\alpha = (1 - c_1^{-1} \gamma^{-2} R_{\min}^{-2} pn) / (1 + c_1^{-1} \gamma^{-2} R_{\min}^{-2} pn)$ is a
591 positive constant.

592 *Proof of Lemma D.2.* To establish an upper bound for $L_S(\mathbf{W}^{(t)})$, we need to first determine a
593 lower bound for the margin $y_i f(\mathbf{W}^{(t)}, \mathbf{x}_i)$. To facilitate improved convergence rate analysis, we
594 will initially present a more precise analysis for the increasing rate of $\rho_{y_i,r,i}^{(t)}$, which takes into
595 account the activation pattern after T_1 given in Lemma B.5. Given the activation pattern outlined in

Lemma B.5, we can provide the following refined upper bound for $F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i)$ and lower bound for $F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i)$:

$$\begin{aligned}
F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i) &= \frac{1}{m} \sum_{r=1}^m \sigma(\langle \mathbf{w}_{y_i,r}^{(t)}, \mathbf{x}_i \rangle) \\
&= \frac{1}{m} \sum_{r=1}^m \langle \mathbf{w}_{y_i,r}^{(t)}, \mathbf{x}_i \rangle \\
&= \frac{1}{m} \sum_{r=1}^m \left[\langle \mathbf{w}_{y_i,r}^{(0)}, \mathbf{x}_i \rangle + \sum_{i'=1}^n \rho_{y_i,r,i'}^{(t)} \|\mathbf{x}_{i'}\|_2^{-2} \cdot \langle \mathbf{x}_{i'}, \mathbf{x}_i \rangle \right] \\
&\leq \frac{1}{m} \sum_{r=1}^m \rho_{y_i,r,i}^{(t)} + \frac{1}{m} \sum_{r=1}^m \sum_{i' \neq i} |\rho_{y_i,r,i'}^{(t)}| R_{\min}^{-2} p + \beta \\
&\leq \frac{1 + c_1^{-1} \gamma^{-2} R_{\min}^{-2} p n}{m} \sum_{r=1}^m \rho_{y_i,r,i}^{(t)} + \beta,
\end{aligned} \tag{D.3}$$

where the second equality is by Lemma B.5; the third equality is by (5.1); the second last inequality is by triangle inequality; the last inequality is by Lemma B.1. And

$$\begin{aligned}
F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i) &= \frac{1}{m} \sum_{r=1}^m \sigma(\langle \mathbf{w}_{-y_i,r}^{(t)}, \mathbf{x}_i \rangle) \\
&= \frac{\gamma}{m} \sum_{r=1}^m \langle \mathbf{w}_{-y_i,r}^{(t)}, \mathbf{x}_i \rangle \\
&= \frac{\gamma}{m} \sum_{r=1}^m \left[\langle \mathbf{w}_{-y_i,r}^{(0)}, \mathbf{x}_i \rangle + \sum_{i'=1}^n \rho_{-y_i,r,i'}^{(t)} \|\mathbf{x}_{i'}\|_2^{-2} \cdot \langle \mathbf{x}_{i'}, \mathbf{x}_i \rangle \right] \\
&\geq \frac{\gamma}{m} \sum_{r=1}^m \rho_{-y_i,r,i}^{(t)} - \frac{\gamma}{m} \sum_{r=1}^m \sum_{i' \neq i} |\rho_{-y_i,r,i'}^{(t)}| R_{\min}^{-2} p - \gamma \beta \\
&\geq \frac{\gamma(1 + c_1^{-1} \gamma^{-2} R_{\min}^{-2} p n)}{m} \sum_{r=1}^m \rho_{-y_i,r,i}^{(t)} - \gamma \beta,
\end{aligned} \tag{D.4}$$

where the second equality is by Lemma B.3; the third equality is by (5.1); the second last inequality is by triangle inequality; the last inequality is by Lemma B.1. By (D.3) and (D.4), we can obtain the following margin upper bound:

$$\begin{aligned}
&y_i f(\mathbf{W}^{(t)}, \mathbf{x}_i) \\
&= F_{y_i}(\mathbf{W}_{y_i}, \mathbf{x}_i) - F_{-y_i}(\mathbf{W}_{-y_i}, \mathbf{x}_i) \\
&\leq \frac{1 + c_1^{-1} \gamma^{-2} R_{\min}^{-2} p n}{m} \sum_{r=1}^m \rho_{y_i,r,i}^{(t)} + \frac{\gamma(1 + c_1^{-1} \gamma^{-2} R_{\min}^{-2} p n)}{m} \sum_{r=1}^m |\rho_{-y_i,r,i}^{(t)}| + (\gamma + 1) \beta \\
&= \frac{1 + c_1^{-1} \gamma^{-2} R_{\min}^{-2} p n}{m} \sum_{r=1}^m (\rho_{y_i,r,i}^{(t)} - \rho_{y_i,r,i}^{(T_1)}) + \frac{\gamma(1 + c_1^{-1} \gamma^{-2} R_{\min}^{-2} p n)}{m} \sum_{r=1}^m (|\rho_{-y_i,r,i}^{(t)}| - |\rho_{-y_i,r,i}^{(T_1)}|) \\
&\quad + \underbrace{\frac{1 + c_1^{-1} \gamma^{-2} R_{\min}^{-2} p n}{m} \sum_{r=1}^m \rho_{y_i,r,i}^{(T_1)} + \frac{\gamma(1 + c_1^{-1} \gamma^{-2} R_{\min}^{-2} p n)}{m} \sum_{r=1}^m |\rho_{-y_i,r,i}^{(T_1)}| + (\gamma + 1) \beta}_{:= C_i} \\
&= \frac{(1 + \gamma^2)(1 + c_1^{-1} \gamma^{-2} R_{\min}^{-2} p n)}{m} \sum_{r=1}^m (\rho_{y_i,r,i}^{(t)} - \rho_{y_i,r,i}^{(T_1)}) + C_i,
\end{aligned}$$

603 where the last equality is by Lemma D.1. Here

$$C_i \leq (\gamma + 1)C'(1 + c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn) + (\gamma + 1)\beta$$

604 by (C.3) and thus can be seen as a constant. Then, it follows that

$$\begin{aligned} |\ell_i'^{(t)}| &= \frac{1}{1 + \exp\{y_i f(\mathbf{W}^{(t)}, \mathbf{x}_i)\}} \\ &\geq \frac{1}{1 + \exp\left\{\frac{(1+\gamma^2)(1+c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn)}{m} \sum_{r=1}^m (\rho_{y_i,r,i}^{(t)} - \rho_{y_i,r,i}^{(T_1)}) + C_i\right\}} \\ &\geq \frac{1}{2} \exp\left\{-\frac{(1+\gamma^2)(1+c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn)}{m} \sum_{r=1}^m (\rho_{y_i,r,i}^{(t)} - \rho_{y_i,r,i}^{(T_1)}) - C_i\right\}, \end{aligned}$$

605 where the last inequality is by $1/(1 + \exp(z)) \geq \exp(-z)/2$ if $z \geq 0$. By the (5.4) and the above

606 lower bound for $|\ell_i'^{(t)}|$, we can get

$$\begin{aligned} &\sum_{r=1}^m (\rho_{y_i,r,i}^{(t+1)} - \rho_{y_i,r,i}^{(T_1)}) \\ &\geq \sum_{r=1}^m (\rho_{y_i,r,i}^{(t)} - \rho_{y_i,r,i}^{(T_1)}) + \frac{\eta \|\mathbf{x}_i\|_2^2 e^{-C_i}}{2n} \exp\left\{-\frac{(1+\gamma^2)(1+c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn)}{m} \sum_{r=1}^m (\rho_{y_i,r,i}^{(t)} - \rho_{y_i,r,i}^{(T_1)})\right\}. \end{aligned} \quad (\text{D.5})$$

607 By applying Lemma H.2 and taking $x_t = \frac{(1+\gamma^2)(1+c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn)}{m} \sum_{r=1}^m (\rho_{y_i,r,i}^{(t+T_1)} - \rho_{y_i,r,i}^{(T_1)})$, we can
608 get the following increasing rate of $\rho_{y_i,r,i}^{(t)}$ for $t \geq T_1$:

$$\begin{aligned} &\sum_{r=1}^m (\rho_{y_i,r,i}^{(t)} - \rho_{y_i,r,i}^{(T_1)}) \\ &\geq \frac{m}{(1+\gamma^2)(1+c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn)} \log\left(1 + \frac{\eta \|\mathbf{x}_i\|_2^2 e^{-C_i} (1+\gamma^2)(1+c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn)}{2nm} \cdot (t - T_1)\right). \end{aligned} \quad (\text{D.6})$$

609 Now we are ready to provide a refined lower bound for the margin $y_i f(\mathbf{W}^{(t)}, \mathbf{x}_i)$. By taking into
610 account the activation pattern after T_1 and conducting similar analysis as (D.3) and (D.4), we can
611 obtain for $t \geq T_1$:

$$\begin{aligned} F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i) &\geq \frac{1 - c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn}{m} \sum_{r=1}^m \rho_{y_i,r,i}^{(t)} - \beta, \\ F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i) &\leq \frac{\gamma(1 - c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn)}{m} \sum_{r=1}^m \rho_{-y_i,r,i}^{(t)} + \gamma\beta. \end{aligned}$$

612 Therefore, we have the following margin lower bound

$$\begin{aligned} &y_i f(\mathbf{W}^{(t)}, \mathbf{x}_i) \\ &= F_{y_i}(\mathbf{W}_{y_i}, \mathbf{x}_i) - F_{-y_i}(\mathbf{W}_{-y_i}, \mathbf{x}_i) \\ &\geq \frac{1 - c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn}{m} \sum_{r=1}^m \rho_{y_i,r,i}^{(t)} + \frac{\gamma(1 - c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn)}{m} \sum_{r=1}^m |\rho_{-y_i,r,i}^{(t)}| - (\gamma + 1)\beta \\ &= \frac{(1 - c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn)}{m} \sum_{r=1}^m (\rho_{y_i,r,i}^{(t)} - \rho_{y_i,r,i}^{(T_1)}) + \frac{\gamma(1 - c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn)}{m} \sum_{r=1}^m (|\rho_{-y_i,r,i}^{(t)}| - |\rho_{-y_i,r,i}^{(T_1)}|) \\ &\quad + \underbrace{\frac{(1 - c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn)}{m} \sum_{r=1}^m \rho_{y_i,r,i}^{(T_1)} + \frac{\gamma(1 - c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn)}{m} \sum_{r=1}^m |\rho_{-y_i,r,i}^{(T_1)}| - (\gamma + 1)\beta}_{:=C'_i} \end{aligned}$$

$$\begin{aligned}
&= \frac{(1 + \gamma^2)(1 - c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn)}{m} \sum_{r=1}^m (\rho_{y_i, r, i}^{(t)} - \rho_{y_i, r, i}^{(T_1)}) + C'_i \\
&\geq \frac{1 - c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn}{1 + c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn} \cdot \log \left(1 + \frac{\eta \|\mathbf{x}_i\|_2^2 e^{-C_i} (1 + \gamma^2) (1 + c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn)}{2nm} \cdot (t - T_1) \right) + C'_i,
\end{aligned}$$

where the last inequality is by (D.6). Here

$$|C'_i| \leq C'(1 + \gamma)(1 - c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn) + (\gamma + 1)\beta$$

by (C.3) and thus can be seen as a constant. Correspondingly, we can get the following for the increasing rate of training loss:

$$\begin{aligned}
L_S(\mathbf{W}^{(t)}) &= \frac{1}{n} \sum_{i=1}^n \ell(y_i f(\mathbf{W}^{(t)}, \mathbf{x}_i)) \\
&= \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i f(\mathbf{W}^{(t)}, \mathbf{x}_i))) \\
&\leq \frac{1}{n} \sum_{i=1}^n \exp(-y_i f(\mathbf{W}^{(t)}, \mathbf{x}_i)) \\
&\leq \frac{e^{-C'_i}}{n} \sum_{i=1}^n \exp \left(-\frac{1 - c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn}{1 + c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn} \cdot \log \left(1 + \frac{\eta \|\mathbf{x}_i\|_2^2 e^{-C_i} (1 + \gamma^2) (1 + c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn)}{2nm} \cdot (t - T) \right) \right) \\
&\leq \frac{e^{-C'_i}}{n} \sum_{i=1}^n \left(1 + \frac{\eta \|\mathbf{x}_i\|_2^2 e^{-C_i} (1 + \gamma^2) (1 + c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn)}{2nm} \cdot (t - T) \right)^{-\frac{1 - c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn}{1 + c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn}} \\
&= O(t^{-(1 - c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn)/(1 + c_1^{-1}\gamma^{-2}R_{\min}^{-2}pn)}),
\end{aligned}$$

which completes the proof. \square

E Coefficient Analysis of ReLU

In this section, we discuss the stable rank of two-layer ReLU neural network, which is defined as

$$\begin{aligned}
f(\mathbf{W}, \mathbf{x}) &= F_{+1}(\mathbf{W}_{+1}, \mathbf{x}) - F_{+1}(\mathbf{W}_{-1}, \mathbf{x}), \\
F_j(\mathbf{W}_j, \mathbf{x}) &= \frac{1}{m} \sum_{r=1}^m \sigma(\langle \mathbf{w}_{j,r}, \mathbf{x} \rangle),
\end{aligned} \tag{E.1}$$

where $\sigma(z) = \max\{0, z\}$ is ReLU activation function.

These results are based on the conclusions in Section A, which hold with high probability. Denote by $\mathcal{E}'_{\text{prelim}}$ the event that all the results in Section A hold (for a given δ , we see $\mathbb{P}(\mathcal{E}'_{\text{prelim}}) \geq 1 - 2\delta$ by a union bound). For simplicity and clarity, we state all the results in this and the following sections conditional on $\mathcal{E}'_{\text{prelim}}$.

Denote $\beta = \max_{i,j,r} \{|\langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle|\}$, $R_{\max} = \max_{i \in [n]} \|\mathbf{x}_i\|_2$, $R_{\min} = \min_{i \in [n]} \|\mathbf{x}_i\|_2$, $p = \max_{i \neq k} |\langle \mathbf{x}_i, \mathbf{x}_k \rangle|$ and suppose $R = R_{\max}/R_{\min}$ is at most an absolute constant. Here we list the exact conditions for $\eta, \sigma_0, R_{\min}, R_{\max}, p$ required by the proofs in this section:

$$\sigma_0 \leq (CR_{\max} \sqrt{\log(mn/\delta)})^{-1}, \tag{E.2}$$

$$\eta \leq (CR_{\max}^2/nm)^{-1}, \tag{E.3}$$

$$R_{\min}^2 \geq CR^2 np, \tag{E.4}$$

where C is a large enough constant. By Lemma A.1, we can upper bound β by $2\sqrt{\log(12mn/\delta)} \cdot \sigma_0 R_{\max}$. Then, by (B.2) and (B.4), it is straightforward to verify the following inequality:

$$\beta \leq c, \quad (\text{E.5})$$

$$R_{\min}^{-2} np \leq c, \quad (\text{E.6})$$

$$R_{\min}^{-2} R^2 np \leq c, \quad (\text{E.7})$$

where c is a sufficiently small constant.

We first introduce the following lemma which characterizes the increasing rate of coefficients $\rho_{j,r,i}^{(t)}$.

Lemma E.1. For two-layer ReLU neural network defined in (E.1), under the same condition as Theorem 4.3, the decomposition coefficients $\rho_{j,r,i}^{(t)}$ satisfy following properties:

$$\bullet \bar{\rho}_{y_i,r,i}^{(t)} \geq c_1 |\rho_{j,r',i'}^{(t)}| \text{ for any } r \in S_i^{(0)}, r' \in [m], j \in \{\pm 1\} \text{ and } i, i' \in [n],$$

$$\bullet \bar{\rho}_{y_i,r,i}^{(t)} \geq c_2 \log \left(1 + \frac{\eta |S_i^{(0)}| \|\mathbf{x}_i\|_2^2 e^{-\beta}}{2nm^2} \cdot t \right) \text{ for any } r \in S_i^{(0)} \text{ and } i \in [n],$$

$$\bullet \bar{\rho}_{y_i,r,i}^{(t)} \leq c_3 \log \left(1 + \frac{2\eta |S_i^{(0)}| \|\mathbf{x}_i\|_2^2 e^{2\beta}}{nm^2} \cdot t \right) \text{ for any } r \in S_i^{(0)} \text{ and } i \in [n],$$

where c_1, c_2, c_3 are constants. And the following activation pattern is also observed: $S_i^{(0)} \subseteq S_i^{(t)}$ where $S_i^{(t)} := \{r \in [m] : \langle \mathbf{w}_{y_i,r}, \mathbf{x}_i \rangle \geq 0\}$, that is, the on-diagonal neuron activated at initialization will remain activated throughout the training.

Proof of Lemma E.1. We first show that the first bullet and $S_i^{(0)} \subseteq S_i^{(t)}$ hold for $t \leq T_1 = C\eta^{-1}nmR_{\max}^{-2}$ where $C = \Theta(1)$ is a constant. Now we prove this by induction. When $t = 0$, the two hypotheses hold naturally. Suppose that there exists time $\tilde{t} \leq T_1$ such that the two hypotheses hold for all time $t \leq \tilde{t} - 1$. We aim to prove they also hold for $t = \tilde{t}$. Recall from Lemma A.3 that

$$\begin{aligned} \bar{\rho}_{j,r,i}^{(t+1)} &= \bar{\rho}_{j,r,i}^{(t)} - \frac{\eta}{nm} \cdot \ell_i^{(t)} \cdot \sigma'(\langle \mathbf{w}_{j,r}, \mathbf{x}_i \rangle) \cdot \|\mathbf{x}_i\|_2^2 \cdot \mathbb{1}(y_i = j), \\ \underline{\rho}_{j,r,i}^{(t+1)} &= \underline{\rho}_{j,r,i}^{(t)} + \frac{\eta}{nm} \cdot \ell_i^{(t)} \cdot \sigma'(\langle \mathbf{w}_{j,r}, \mathbf{x}_i \rangle) \cdot \|\mathbf{x}_i\|_2^2 \cdot \mathbb{1}(y_i = -j), \end{aligned}$$

we can get

$$\bar{\rho}_{j,r,i}^{(t+1)} \leq \bar{\rho}_{j,r,i}^{(t)} + \frac{\eta}{nm} \cdot \|\mathbf{x}_i\|_2^2 \leq \bar{\rho}_{j,r,i}^{(t)} + \frac{\eta R_{\max}^2}{nm}, \quad (\text{E.8})$$

$$|\underline{\rho}_{j,r,i}^{(t+1)}| \leq |\underline{\rho}_{j,r,i}^{(t)}| + \frac{\eta}{nm} \cdot \|\mathbf{x}_i\|_2^2 \leq |\underline{\rho}_{j,r,i}^{(t)}| + \frac{\eta R_{\max}^2}{nm}. \quad (\text{E.9})$$

Therefore, $\max_{j,r,i} \{\bar{\rho}_{j,r,i}^{(t)}, |\underline{\rho}_{j,r,i}^{(t)}|\} = O(1)$ for any $t \leq T_1$ and hence $\max_i \{F_{+1}(\mathbf{W}_{+1}^{(t)}, \mathbf{x}_i), F_{-1}(\mathbf{W}_{-1}^{(t)}, \mathbf{x}_i)\} = O(1)$ for any $t \leq T_1$. Thus there exists a positive constant c such that $|\ell_i^{(t)}| \geq c$ for any $t \leq T_1$. By induction hypothesis, we have $S_i^{(0)} \subseteq S_i^{(t)}$ for all $0 \leq t \leq \tilde{t} - 1$ and hence $\sigma'(\langle \mathbf{w}_{y_i,r}, \mathbf{x}_i \rangle) = 1$ for all $0 \leq t \leq \tilde{t} - 1$. And it follows that for $r \in S_i^{(0)}$

$$\begin{aligned} \bar{\rho}_{y_i,r,i}^{(t+1)} &= \bar{\rho}_{y_i,r,i}^{(t)} + \frac{\eta}{nm} \cdot |\ell_i^{(t)}| \cdot \|\mathbf{x}_i\|_2^2 \geq \bar{\rho}_{y_i,r,i}^{(t)} + \frac{c\eta}{nm} \cdot \|\mathbf{x}_i\|_2^2, \forall 0 \leq t \leq \tilde{t} - 1, \\ \bar{\rho}_{y_i,r,i}^{(\tilde{t})} &\geq \frac{c\eta\tilde{t}}{nm} \cdot \|\mathbf{x}_i\|_2^2 \geq \frac{c\eta R_{\min}^2 \tilde{t}}{nm}. \end{aligned} \quad (\text{E.10})$$

On the other hand, by (E.8) and (E.9), we have

$$\bar{\rho}_{j,r',i'}^{(\tilde{t})} \leq \frac{\eta R_{\max}^2 \tilde{t}}{nm}, |\underline{\rho}_{j,r',i'}^{(\tilde{t})}| \leq \frac{\eta R_{\max}^2 \tilde{t}}{nm} \implies |\rho_{j,r',i'}^{(\tilde{t})}| \leq \frac{\eta R_{\max}^2 \tilde{t}}{nm}. \quad (\text{E.11})$$

650 Dividing (E.10) by (E.11), we can get

$$\frac{\bar{\rho}_{y_i,r,i}^{(\tilde{t})}}{|\rho_{j,r',i'}^{(\tilde{t})}|} \geq \frac{cR_{\min}^2}{R_{\max}^2}, \forall r \in S_i^{(0)}, j \in \{\pm 1\}, i, i' \in [n], \quad (\text{E.12})$$

651 which indicates that the first bullet holds for time $t = \tilde{t}$ as long as $c_1 \leq (cR_{\min}^2)/R_{\max}^2$. For $r \in S_i^{(0)}$,
 652 we have

$$\begin{aligned} \langle \mathbf{w}_{y_i,r}^{(\tilde{t})}, \mathbf{x}_i \rangle &= \langle \mathbf{w}_{y_i,r}^{(0)}, \mathbf{x}_i \rangle + \sum_{i'=1}^n \rho_{y_i,r,i'}^{(\tilde{t})} \|\mathbf{x}_{i'}\|_2^{-2} \cdot \langle \mathbf{x}_{i'}, \mathbf{x}_i \rangle \\ &= \langle \mathbf{w}_{y_i,r}^{(0)}, \mathbf{x}_i \rangle + \bar{\rho}_{y_i,r,i}^{(\tilde{t})} + \sum_{i' \neq i} \rho_{y_i,r,i'}^{(\tilde{t})} \|\mathbf{x}_{i'}\|_2^{-2} \cdot \langle \mathbf{x}_{i'}, \mathbf{x}_i \rangle \\ &\geq \bar{\rho}_{y_i,r,i}^{(\tilde{t})} - \sum_{i' \neq i} |\rho_{y_i,r,i'}^{(\tilde{t})}| R_{\min}^{-2} p \\ &\geq \bar{\rho}_{y_i,r,i}^{(\tilde{t})} - \sum_{i' \neq i} \frac{R_{\max}^2}{cR_{\min}^2} \bar{\rho}_{y_i,r,i}^{(\tilde{t})} \cdot R_{\min}^{-2} p \\ &\geq \left(1 - \frac{R_{\max}^2}{cR_{\min}^4} pn\right) \cdot \bar{\rho}_{y_i,r,i}^{(\tilde{t})} \geq 0, \end{aligned}$$

653 where the second inequality is by (E.12). This implies that $S_i^{(0)} \subseteq S_i^{(\tilde{t})}$ holds for time $t = \tilde{t}$, which
 654 completes the induction. By then, we have already proved that the first bullet and $S_i^{(0)} \subseteq S_i^{(t)}$ hold
 655 for $t \leq T_1 = C\eta^{-1}nmR_{\max}^{-2}$.

656 Next, we will prove by induction that the three bullets as well as $S_i^{(0)} \subseteq S_i^{(t)}$ hold for any time $t \geq 0$.
 657 The second and third bullets are obvious at $t = 0$ as all the coefficients are zero. Suppose there exists
 658 \tilde{t} such that the three bullets as well as $S_i^{(0)} \subseteq S_i^{(\tilde{t})}$ hold for all time $0 \leq t \leq \tilde{t} - 1$. We aim to prove
 659 that they also hold for $t = \tilde{t}$. We first prove that the second and third bullets hold for $t = \tilde{t}$. To prove
 660 this, we first provide more precise upper and lower bounds for $|\ell_i^{(\tilde{t})}|$. For lower bound, we have

$$\begin{aligned} |\ell_i^{(\tilde{t})}| &= \frac{1}{1 + \exp \{F_{y_i}(\mathbf{W}_{y_i}^{(\tilde{t})}, \mathbf{x}_i) - F_{-y_i}(\mathbf{W}_{-y_i}^{(\tilde{t})}, \mathbf{x}_i)\}} \\ &\geq \frac{1}{1 + \exp \{F_{y_i}(\mathbf{W}_{y_i}^{(\tilde{t})}, \mathbf{x}_i)\}} \\ &= \frac{1}{1 + \exp \{\frac{1}{m} \sum_{r \in S_i^{(\tilde{t})}} \langle \mathbf{w}_{y_i,r}^{(\tilde{t})}, \mathbf{x}_i \rangle\}} \end{aligned} \quad (\text{E.13})$$

661 and

$$\begin{aligned} \sum_{r \in S_i^{(\tilde{t})}} \langle \mathbf{w}_{y_i,r}^{(\tilde{t})}, \mathbf{x}_i \rangle &= \sum_{r \in S_i^{(\tilde{t})}} \left(\langle \mathbf{w}_{y_i,r}^{(0)}, \mathbf{x}_i \rangle + \bar{\rho}_{y_i,r,i}^{(\tilde{t})} + \sum_{i' \neq i} \rho_{y_i,r,i'}^{(\tilde{t})} \|\mathbf{x}_{i'}\|_2^{-2} \cdot \langle \mathbf{x}_{i'}, \mathbf{x}_i \rangle \right) \\ &\leq \sum_{r \in S_i^{(\tilde{t})}} \bar{\rho}_{y_i,r,i}^{(\tilde{t})} + \sum_{r \in S_i^{(\tilde{t})}} \sum_{i' \neq i} |\rho_{y_i,r,i'}^{(\tilde{t})}| R_{\min}^{-2} p + |S_i^{(\tilde{t})}| \cdot \beta \\ &\leq \sum_{r \in S_i^{(\tilde{t})}} \bar{\rho}_{y_i,r,i}^{(\tilde{t})} + \frac{|S_i^{(\tilde{t})}|}{c_1 |S_i^{(0)}|} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i,r,i}^{(\tilde{t})} R_{\min}^{-2} pn + |S_i^{(\tilde{t})}| \cdot \beta \\ &\leq \frac{|S_i^{(\tilde{t})}|}{|S_i^{(0)}|} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i,r,i}^{(\tilde{t})} + \frac{|S_i^{(\tilde{t})}|}{c_1 |S_i^{(0)}|} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i,r,i}^{(\tilde{t})} R_{\min}^{-2} pn + |S_i^{(\tilde{t})}| \cdot \beta \end{aligned}$$

$$\leq c'(1 + R_{\min}^{-2}pn/c_1) \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} + |S_i^{(t)}| \cdot \beta, \quad (\text{E.14})$$

where the first inequality is by triangle inequality; the second inequality is by the first induction hypothesis that $\bar{\rho}_{y_i, r, i}^{(t)} \geq c_1 |\rho_{y_i, r', i'}^{(t)}|$ for $r \in S_i^{(0)}$ and $0 \leq t \leq \tilde{t} - 1$ and hence $|\rho_{y_i, r', i'}^{(t)}| \leq \frac{1}{c_1 |S_i^{(0)}|} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)}$; the third inequality is by

$$\begin{aligned} \bar{\rho}_{y_i, r', i}^{(t)} &= \frac{\eta}{nm} \sum_{s=0}^{t-1} |\ell_i^{(s)}| \cdot \sigma'(\langle \mathbf{w}_{y_i, r'}^{(s)}, \mathbf{x}_i \rangle) \cdot \|\mathbf{x}_i\|_2^2 \leq \frac{\eta}{nm} \sum_{s=0}^{t-1} |\ell_i^{(s)}| \cdot \|\mathbf{x}_i\|_2^2, \\ \bar{\rho}_{y_i, r, i}^{(t)} &= \frac{\eta}{nm} \sum_{s=0}^{t-1} |\ell_i^{(s)}| \cdot \|\mathbf{x}_i\|_2^2, \end{aligned}$$

and hence $\bar{\rho}_{y_i, r', i}^{(t)} \leq \bar{\rho}_{y_i, r, i}^{(t)}, \forall r' \in S_i^{(t)} \setminus S_i^{(0)}, r \in S_i^{(0)}$ for $0 \leq t \leq \tilde{t} - 1$; the last inequality is by $|S_i^{(t)}| \leq m \leq c'|S_i^{(0)}|$ and c' can be taken as 2.5 by Lemma A.2. By plugging (E.14) back into (E.13), we can get

$$\begin{aligned} |\ell_i^{(t)}| &\geq \frac{1}{1 + \exp \left\{ \frac{c'(1 + R_{\min}^{-2}pn/c_1)}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} + \frac{|S_i^{(t)}|}{m} \cdot \beta \right\}} \\ &\geq \frac{1}{1 + \exp \left\{ \frac{c'(1 + R_{\min}^{-2}pn/c_1)}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} + \beta \right\}} \\ &\geq \frac{1}{2} \exp \left\{ - \frac{c'(1 + R_{\min}^{-2}pn/c_1)}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} - \beta \right\}, \forall 0 \leq t \leq \tilde{t} - 1. \end{aligned} \quad (\text{E.15})$$

For upper bound of $|\ell_i^{(t)}|$, we first bound $F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i) - F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i)$ as follows:

$$\begin{aligned} &F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i) - F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i) \\ &\geq \frac{1}{m} \left(\sum_{r \in S_i^{(t)}} \bar{\rho}_{y_i, r, i}^{(t)} - \sum_{r \in S_i^{(t)}} \sum_{i' \neq i} |\rho_{y_i, r, i'}^{(t)}| R_{\min}^{-2}p - |S_i^{(t)}| \cdot \beta \right) - \frac{1}{m} \sum_{r=1}^m \left(\beta + \sum_{i' \neq i} |\rho_{-y_i, r, i'}^{(t)}| R_{\min}^{-2}p \right) \\ &\geq \frac{1}{m} \sum_{r \in S_i^{(t)}} \bar{\rho}_{y_i, r, i}^{(t)} - \frac{|S_i^{(t)}|}{c_1 m |S_i^{(0)}|} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} R_{\min}^{-2}pn - \frac{1}{c_1 |S_i^{(0)}|} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} R_{\min}^{-2}pn - 2\beta \\ &\geq \frac{1}{m} \sum_{r \in S_i^{(t)}} \bar{\rho}_{y_i, r, i}^{(t)} - \frac{2}{c_1 |S_i^{(0)}|} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} R_{\min}^{-2}pn - 2\beta \\ &\geq \frac{1}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} - \frac{2c' R_{\min}^{-2}pn}{c_1 m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} - 2\beta \\ &= \frac{1 - 2c' R_{\min}^{-2}pn/c_1}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} - 2\beta, \end{aligned}$$

where the first inequality is by

$$\begin{aligned} F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i) &= \frac{1}{m} \sum_{r \in S_i^{(t)}} \langle \mathbf{w}_{y_i, r}^{(t)}, \mathbf{x}_i \rangle \\ &= \frac{1}{m} \sum_{r \in S_i^{(t)}} \left(\langle \mathbf{w}_{y_i, r}^{(0)}, \mathbf{x}_i \rangle + \bar{\rho}_{y_i, r, i}^{(t)} + \sum_{i' \neq i} \rho_{y_i, r, i'}^{(t)} \|\mathbf{x}_{i'}\|_2^{-2} \cdot \langle \mathbf{x}_{i'}, \mathbf{x}_i \rangle \right) \end{aligned}$$

$$\geq \frac{1}{m} \left(\sum_{r \in S_i^{(t)}} \bar{\rho}_{y_i, r, i}^{(t)} - \sum_{r \in S_i^{(t)}} \sum_{i' \neq i} |\rho_{y_i, r, i'}^{(t)}| R_{\min}^{-2} p - |S_i^{(t)}| \cdot \beta \right),$$

670 and

$$\begin{aligned} \langle \mathbf{w}_{-y_i, r}^{(t)}, \mathbf{x}_i \rangle &= \langle \mathbf{w}_{-y_i, r}^{(0)}, \mathbf{x}_i \rangle + \sum_{i'=1}^n \rho_{-y_i, r, i'}^{(t)} \|\mathbf{x}_{i'}\|_2^{-2} \cdot \langle \mathbf{x}_{i'}, \mathbf{x}_i \rangle \\ &= \langle \mathbf{w}_{-y_i, r}^{(0)}, \mathbf{x}_i \rangle + \rho_{-y_i, r, i}^{(t)} + \sum_{i' \neq i} \rho_{-y_i, r, i'}^{(t)} \|\mathbf{x}_{i'}\|_2^{-2} \cdot \langle \mathbf{x}_{i'}, \mathbf{x}_i \rangle \\ &\leq \beta + \sum_{i' \neq i} |\rho_{-y_i, r, i'}^{(t)}| R_{\min}^{-2} p, \end{aligned}$$

671 and hence

$$F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i) = \frac{1}{m} \sum_{r=1}^m \sigma(\langle \mathbf{w}_{-y_i, r}^{(t)}, \mathbf{x}_i \rangle) \leq \frac{1}{m} \sum_{r=1}^m \left(\beta + \sum_{i' \neq i} |\rho_{-y_i, r, i'}^{(t)}| R_{\min}^{-2} p \right); \quad (\text{E.16})$$

672 the second inequality is by the first induction hypothesis that $\bar{\rho}_{y_i, r, i}^{(t)} \geq c_1 |\rho_{y_i, r', i'}^{(t)}|$, $\bar{\rho}_{y_i, r, i}^{(t)} \geq$
673 $c_1 |\rho_{-y_i, r', i'}^{(t)}|$ and hence $|\rho_{y_i, r', i'}^{(t)}| \leq \frac{1}{c_1 |S_i^{(0)}|} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)}$, $|\rho_{-y_i, r', i'}^{(t)}| \leq \frac{1}{c_1 |S_i^{(0)}|} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)}$
674 for $r \in S_i^{(0)}$ and $0 \leq t \leq \tilde{t} - 1$; the third inequality is by $|S_i^{(t)}| \leq m$; the fourth inequality is by
675 $m \leq c' |S_i^{(0)}|$. Therefore,

$$\begin{aligned} |\ell_i^{(t)}| &= \frac{1}{1 + \exp \{ F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i) - F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i) \}} \\ &\leq \exp \{ -F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i) + F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i) \} \\ &\leq \exp \left\{ -\frac{1 - 2c' R_{\min}^{-2} p n / c_1}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} + 2\beta \right\}, \forall 0 \leq t \leq \tilde{t} - 1. \end{aligned}$$

676 By the induction hypothesis, we know that $S_i^{(0)} \subseteq S_i^{(t)}$ for all $0 \leq t \leq \tilde{t} - 1$ and hence
677 $\sigma'(\langle \mathbf{w}_{y_i, r}^{(t)}, \mathbf{x}_i \rangle) = 1$ for all $r \in S_i^{(0)}$ and $0 \leq t \leq \tilde{t} - 1$. By (E.15) and (E.16), it follows that
678 for all $0 \leq t \leq \tilde{t} - 1$

$$\sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t+1)} \geq \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} + \frac{\eta |S_i^{(0)}| \|\mathbf{x}_i\|_2^2 e^{-\beta}}{2nm} \cdot \exp \left\{ -\frac{c'(1 + R_{\min}^{-2} p n / c_1)}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} \right\}, \quad (\text{E.17})$$

$$\sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t+1)} \leq \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} + \frac{\eta |S_i^{(0)}| \|\mathbf{x}_i\|_2^2 e^{2\beta}}{nm} \cdot \exp \left\{ -\frac{1 - 2c' R_{\min}^{-2} p n / c_1}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} \right\}. \quad (\text{E.18})$$

679 By applying Lemma H.2 to (E.17) and taking $x_t = \frac{c'(1 + R_{\min}^{-2} p n / c_1)}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)}$, we can get

$$\sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} \geq \frac{m}{c'(1 + R_{\min}^{-2} p n / c_1)} \log \left(1 + \frac{c'(1 + R_{\min}^{-2} p n / c_1) \eta |S_i^{(0)}| \|\mathbf{x}_i\|_2^2 e^{-\beta}}{2nm^2} \cdot t \right), \forall 0 \leq t \leq \tilde{t}. \quad (\text{E.19})$$

By applying Lemma H.1 to (E.18) and taking $x_t = \frac{1-2c'R_{\min}^{-2}pn/c_1}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i,r,i}^{(t)}$, we can get for any $0 \leq t \leq \tilde{t}$ that

$$\begin{aligned} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i,r,i}^{(t)} &\leq \frac{m}{1-2c'R_{\min}^{-2}pn/c_1} \log \left(1 + \frac{(1-2c'R_{\min}^{-2}pn/c_1)\eta|S_i^{(0)}|\|\mathbf{x}_i\|_2^2 e^{2\beta}}{nm^2} \right. \\ &\quad \left. \exp \left(\frac{(1-2c'R_{\min}^{-2}pn/c_1)\eta|S_i^{(0)}|\|\mathbf{x}_i\|_2^2 e^{2\beta}}{nm^2} \right) \cdot t \right) \\ &\leq \frac{m}{1-2c'R_{\min}^{-2}pn/c_1} \log \left(1 + \frac{2(1-2c'R_{\min}^{-2}pn/c_1)\eta|S_i^{(0)}|\|\mathbf{x}_i\|_2^2 e^{2\beta}}{nm^2} \cdot t \right), \end{aligned} \quad (\text{E.20})$$

where the last inequality is by $\eta \leq (CR_{\max}^2/nm)^{-1}$, C is a large enough constant and hence $(1-2c'R_{\min}^{-2}pn/c_1)\eta|S_i^{(0)}|\|\mathbf{x}_i\|_2^2 e^{2\beta}/nm^2 \leq \log 2$. Since $S_i^{(0)} \subseteq S_i^{(t)}$ for all $0 \leq t \leq \tilde{t}-1$ and hence $\sigma'(\langle \mathbf{w}_{y_i,r}^{(t)}, \mathbf{x}_i \rangle) = 1$ for all $r \in S_i^{(0)}$ and $0 \leq t \leq \tilde{t}-1$, we have

$$\bar{\rho}_{y_i,r,i}^{(t)} = \frac{\eta}{nm} \sum_{s=0}^{t-1} |\ell_i^{(s)}| \cdot \|\mathbf{x}_i\|_2^2, \forall 0 \leq t \leq \tilde{t}.$$

Accordingly, it holds that

$$\bar{\rho}_{y_i,r,i}^{(t)} = \bar{\rho}_{y_i,r',i}^{(t)}, \forall r, r' \in S_i^{(0)}, \forall 0 \leq t \leq \tilde{t}.$$

Applying this to (E.19) and (E.20), we can get

$$\begin{aligned} \bar{\rho}_{y_i,r,i}^{(t)} &\geq \frac{m}{c'(1+R_{\min}^{-2}pn/c_1)|S_i^{(0)}|} \log \left(1 + \frac{c'(1+R_{\min}^{-2}pn/c_1)\eta|S_i^{(0)}|\|\mathbf{x}_i\|_2^2 e^{-\beta}}{2nm^2} \cdot t \right) \\ &\geq \frac{1}{c'(1+R_{\min}^{-2}pn/c_1)} \log \left(1 + \frac{c'(1+R_{\min}^{-2}pn/c_1)\eta|S_i^{(0)}|\|\mathbf{x}_i\|_2^2 e^{-\beta}}{2nm^2} \cdot t \right), \forall 0 \leq t \leq \tilde{t}, \\ \bar{\rho}_{y_i,r,i}^{(t)} &\leq \frac{m}{(1-2c'R_{\min}^{-2}pn/c_1)|S_i^{(0)}|} \log \left(1 + \frac{2(1-2c'R_{\min}^{-2}pn/c_1)\eta|S_i^{(0)}|\|\mathbf{x}_i\|_2^2 e^{2\beta}}{nm^2} \cdot t \right), \\ &\leq \frac{c'}{(1-2c'R_{\min}^{-2}pn/c_1)} \log \left(1 + \frac{2(1-2c'R_{\min}^{-2}pn/c_1)\eta|S_i^{(0)}|\|\mathbf{x}_i\|_2^2 e^{2\beta}}{nm^2} \cdot t \right), \forall 0 \leq t \leq \tilde{t}, \end{aligned} \quad (\text{E.21})$$

By taking

$$c_2 = \frac{1}{c'(1+R_{\min}^{-2}pn/c_1)}, c_3 = \frac{c'}{(1-2c'R_{\min}^{-2}pn/c_1)},$$

the above inequalities indicates that the second and third bullets hold at time $t = \tilde{t}$. For the first bullet, it is only necessary to consider the situation where $\tilde{t} \geq T_1 = C\eta^{-1}nmR_{\max}^{-2}$. In order to apply Lemma H.4 (requiring $b > a$), we loosen the bounds in (E.21) as follows:

$$\bar{\rho}_{y_i,r,i}^{(t)} \geq c_2 \log \left(1 + \frac{\eta R_{\min}^2 e^{-\beta}}{5nm} \cdot t \right), \forall 0 \leq t \leq \tilde{t}, \quad (\text{E.22})$$

$$\bar{\rho}_{y_i,r,i}^{(t)} \leq c_3 \log \left(1 + \frac{6\eta R_{\max}^2 e^{2\beta}}{5nm} \cdot t \right), \forall 0 \leq t \leq \tilde{t}, \quad (\text{E.23})$$

where we use $0.4m \leq |S_i^{(0)}| \leq 0.6m$.

By applying Lemma H.4 to (E.22) and (E.23), we can get for any $i, i' \in [n]$, $r \in S_i^{(0)}$, $r' \in S_{i'}^{(0)}$ and $T_1 \leq t \leq \tilde{t}$ that

$$\begin{aligned} \frac{\bar{\rho}_{y_i, r, i}^{(t)}}{\bar{\rho}_{y_{i'}, r', i'}^{(t)}} &\geq \frac{c_2}{c_3} \cdot \frac{\log \left(1 + \frac{\eta R_{\min}^2 e^{-\beta}}{5nm} \cdot t \right)}{\log \left(1 + \frac{6\eta R_{\max}^2 e^{2\beta}}{5nm} \cdot t \right)} \\ &\geq \frac{c_2}{c_3} \cdot \frac{\log \left(1 + \frac{\eta R_{\max}^2 e^{2\beta}}{5nm} \cdot T_1 \right)}{\log \left(1 + \frac{6\eta R_{\max}^2 e^{2\beta}}{5nm} \cdot T_1 \right)} \\ &= \frac{c_2}{c_3} \cdot \frac{\log \left(1 + 0.2C e^{-\beta} R_{\min}^2 \right)}{\log \left(1 + 1.2C e^{2\beta} R_{\max}^2 \right)}. \end{aligned}$$

Notice that $S_{i'}^{(0)} \subseteq S_{i'}^{(t)}$ for all $0 \leq t \leq \tilde{t} - 1$ and hence $\sigma'(\langle \mathbf{w}_{y_{i'}, r}^{(t)}, \mathbf{x}_{i'} \rangle) = 1$ for all $r \in S_{i'}^{(0)}$ and $0 \leq t \leq \tilde{t} - 1$, we have

$$\begin{aligned} |\rho_{j, r'', i'}^{(t)}| &= \frac{\eta}{nm} \sum_{s=0}^{t-1} |\ell_{i'}^{(s)}| \cdot \sigma'(\langle \mathbf{w}_{j, r''}^{(s)}, \mathbf{x}_{i'} \rangle) \cdot \|\mathbf{x}_i\|_2^2 \\ &\leq \frac{\eta}{nm} \sum_{s=0}^{t-1} |\ell_{i'}^{(s)}| \cdot \|\mathbf{x}_{i'}\|_2^2, \quad \forall j \in \{\pm 1\}, r'' \in [m], i' \in [n], \\ \bar{\rho}_{y_{i'}, r', i'}^{(t)} &= \frac{\eta}{nm} \sum_{s=0}^{t-1} |\ell_{i'}^{(s)}| \cdot \|\mathbf{x}_{i'}\|_2^2, \quad \forall r' \in S_{i'}^{(0)}, i' \in [n], \end{aligned}$$

and hence $|\rho_{j, r', i'}^{(t)}| \leq \bar{\rho}_{y_{i'}, r', i'}^{(t)}$ for $0 \leq t \leq \tilde{t}$. Therefore, as long as

$$c_1 \leq \frac{c_2}{c_3} \cdot \frac{\log \left(1 + 0.2C e^{-\beta} R_{\min}^2 \right)}{\log \left(1 + 1.2C e^{2\beta} R_{\max}^2 \right)},$$

the first bullet hold for time $t = \tilde{t}$. This condition holds as long as

$$\begin{aligned} c' &= 2.5, \\ 2c' c_1^{-1} R_{\min}^{-2} pn &\leq 0.5 \implies c_2 \geq 0.37, c_3 \leq 5, \\ c_1 &= \frac{\log \left(1 + 0.2C e^{-\beta} R_{\min}^2 \right)}{14 \log \left(1 + 1.2C e^{2\beta} R_{\max}^2 \right)}. \end{aligned}$$

Finally, we verify that $S_i^{(0)} \subseteq S_i^{(\tilde{t})}$. For $r \in S_i^{(0)}$, we have

$$\begin{aligned} \langle \mathbf{w}_{y_i, r}^{(\tilde{t})}, \mathbf{x}_i \rangle &= \langle \mathbf{w}_{y_i, r}^{(0)}, \mathbf{x}_i \rangle + \sum_{i'=1}^n \rho_{y_i, r, i'}^{(\tilde{t})} \|\mathbf{x}_{i'}\|_2^{-2} \cdot \langle \mathbf{x}_{i'}, \mathbf{x}_i \rangle \\ &= \langle \mathbf{w}_{y_i, r}^{(0)}, \mathbf{x}_i \rangle + \bar{\rho}_{y_i, r, i}^{(\tilde{t})} + \sum_{i' \neq i} \rho_{y_i, r, i'}^{(\tilde{t})} \|\mathbf{x}_{i'}\|_2^{-2} \cdot \langle \mathbf{x}_{i'}, \mathbf{x}_i \rangle \\ &\geq \bar{\rho}_{y_i, r, i}^{(\tilde{t})} - \sum_{i' \neq i} |\rho_{y_i, r, i'}^{(\tilde{t})}| R_{\min}^{-2} p \\ &\geq \bar{\rho}_{y_i, r, i}^{(\tilde{t})} - (R_{\min}^{-2} pn / c_1) \bar{\rho}_{y_i, r, i'}^{(\tilde{t})} \\ &= (1 - R_{\min}^{-2} pn / c_1) \cdot \bar{\rho}_{y_i, r, i}^{(\tilde{t})} \geq 0, \end{aligned}$$

where the second inequality is by $|\rho_{y_i, r, i'}^{(\tilde{t})}| \leq \bar{\rho}_{y_i, r, i}^{(\tilde{t})}/c_1$. This implies that $S_i^{(0)} \subseteq S_i^{(\tilde{t})}$ holds for time $t = \tilde{t}$, which completes the induction. \square

Next, we show that $|\rho_{-y_i, r, i'}^{(t)}|$ will be much smaller than $|\bar{\rho}_{y_i, r', i}^{(t)}|$ as the training goes on.

Lemma E.2. There exists time T_2 and constant c such that for any time $t \geq T_2$

$$|\rho_{y_i, r, i'}^{(t)}| \leq cR_{\min}^{-2}pn|\bar{\rho}_{y_i, r', i}^{(t)}|,$$

where $r \in [m]$, $r' \in S_i^{(0)}$ and $i, i' \in [n]$ satisfying $-y_i = y_{i'}$.

Proof of Lemma E.2. First, we will prove by induction that for $r \in [m]$, $r' \in S_i^{(0)}$ and $i, i' \in [n]$ satisfying $-y_i = y_{i'}$ it holds that

$$|\rho_{y_i, r, i'}^{(t)}| \leq \beta + 1 + R_{\min}^{-2}pn|\bar{\rho}_{y_i, r', i}^{(t)}|/c_1. \quad (\text{E.24})$$

This result holds naturally when $t = 0$ since all the coefficients are zero. Suppose that there exists time \tilde{t} such that the induction hypothesis (E.24) holds for all time $t \leq \tilde{t} - 1$. We aim to prove that (E.24) also holds for $t = \tilde{t}$. In the following analysis, two cases will be considered: $|\rho_{y_i, r, i'}^{(\tilde{t}-1)}| > \beta + \sum_{k \neq i'} |\rho_{y_i, r, k}^{(\tilde{t}-1)}| \|\mathbf{x}_k\|_2^{-2}p$ and $|\rho_{y_i, r, i'}^{(\tilde{t}-1)}| \leq \beta + \sum_{k \neq i'} |\rho_{y_i, r, k}^{(\tilde{t}-1)}| \|\mathbf{x}_k\|_2^{-2}p$.

For if $|\rho_{y_i, r, i'}^{(\tilde{t}-1)}| > \beta + \sum_{k \neq i'} |\rho_{y_i, r, k}^{(\tilde{t}-1)}| \|\mathbf{x}_k\|_2^{-2}p$, then by the decomposition (5.1) we have

$$\begin{aligned} \langle \mathbf{w}_{y_i, r}^{(\tilde{t}-1)}, \mathbf{x}_{i'} \rangle &= \langle \mathbf{w}_{y_i, r}^{(0)}, \mathbf{x}_{i'} \rangle + \rho_{y_i, r, i'}^{(\tilde{t}-1)} + \sum_{k \neq i'} \rho_{y_i, r, k}^{(\tilde{t}-1)} \|\mathbf{x}_k\|_2^{-2} \langle \mathbf{x}_k, \mathbf{x}_{i'} \rangle \\ &\leq \rho_{y_i, r, i'}^{(\tilde{t}-1)} + \beta + \sum_{k \neq i'} |\rho_{y_i, r, k}^{(\tilde{t}-1)}| \|\mathbf{x}_k\|_2^{-2}p < 0. \end{aligned}$$

and hence

$$\rho_{y_i, r, i'}^{(\tilde{t})} = \rho_{y_i, r, i'}^{(\tilde{t}-1)} + \frac{\eta}{nm} \cdot \ell_i'(\tilde{t}-1) \cdot \sigma'(\langle \mathbf{w}_{y_i, r}^{(\tilde{t}-1)}, \mathbf{x}_{i'} \rangle) \cdot \|\mathbf{x}_{i'}\|_2^2 = \rho_{y_i, r, i'}^{(\tilde{t}-1)}.$$

Therefore, by induction hypothesis (E.24) at time $t = \tilde{t} - 1$, we have

$$\rho_{y_i, r, i'}^{(\tilde{t})} = \rho_{y_i, r, i'}^{(\tilde{t}-1)} \leq \beta + 1 + R_{\min}^{-2}pn|\bar{\rho}_{y_i, r', i}^{(\tilde{t}-1)}|/c_1 \leq \beta + 1 + R_{\min}^{-2}pn|\bar{\rho}_{y_i, r', i}^{(\tilde{t})}|/c_1.$$

For if $|\rho_{y_i, r, i'}^{(\tilde{t}-1)}| \leq \beta + \sum_{k \neq i'} |\rho_{y_i, r, k}^{(\tilde{t}-1)}| \|\mathbf{x}_k\|_2^{-2}p$, by the first bullet in Lemma E.1, we have

$$|\rho_{y_i, r, i'}^{(\tilde{t}-1)}| \leq \beta + \sum_{k \neq i'} |\bar{\rho}_{y_i, r', i}^{(\tilde{t}-1)}| \|\mathbf{x}_k\|_2^{-2}p/c_1 \leq \beta + |\bar{\rho}_{y_i, r', i}^{(\tilde{t}-1)}| R_{\min}^{-2}pn/c_1, \quad (\text{E.25})$$

and thus

$$\begin{aligned} -\rho_{y_i, r, i'}^{(\tilde{t})} &= -\rho_{y_i, r, i'}^{(\tilde{t}-1)} + \frac{\eta}{nm} \cdot |\ell_i'(\tilde{t}-1)| \cdot \sigma'(\langle \mathbf{w}_{y_i, r}^{(\tilde{t}-1)}, \mathbf{x}_{i'} \rangle) \cdot \|\mathbf{x}_{i'}\|_2^2 \\ &\leq -\rho_{y_i, r, i'}^{(\tilde{t}-1)} + \frac{\eta R_{\max}^2}{nm} \\ &\leq \beta + 1 + |\bar{\rho}_{y_i, r', i}^{(\tilde{t})}| R_{\min}^{-2}pn/c_1, \end{aligned}$$

where the last inequality is by (E.25) and $\eta \leq (CR_{\max}^2/nm)^{-1}$ with a sufficiently large constant C .

This demonstrates that inequality (E.24) holds for $t = \tilde{t}$, thereby completing the induction process.

By Lemma E.1, we know that there exists time T' such that

$$|\bar{\rho}_{y_i, r', i}^{(t)}| \geq c_1(\beta + 1)R_{\min}^2/pn,$$

for any time $t \geq T'$. Taking $T_2 = T'$ and $c = 2/c_1$, we have

$$|\rho_{y_i, r, i'}^{(t)}| \leq \beta + 1 + R_{\min}^{-2} p n |\bar{\rho}_{y_i, r', i}^{(t)}| / c_1 \leq 2 R_{\min}^{-2} p n |\bar{\rho}_{y_i, r', i}^{(t)}| / c_1 = c R_{\min}^{-2} p n |\bar{\rho}_{y_i, r', i}^{(t)}|,$$

which completes the proof. \square

Given Lemma E.2, the following corollary can be directly obtained.

Corollary E.3. There exists time T_3 and constant c such that for any time $t \geq T_3$

$$|\rho_{y_i, r, i'}^{(t)}| \leq c R_{\min}^{-2} p n |\bar{\rho}_{y_i, r', i}^{(t)}|, \forall r \in [m], r' \in S_i^{(0)}, i, i' \in [n] \text{ with } -y_i = y_{i'}.$$

F Stable Rank of ReLU Network

In this section, we consider the properties of stable rank of the weight matrix $\mathbf{W}^{(t)}$ found by gradient descent at time t , defined as $\|\mathbf{W}^{(t)}\|_F^2 / \|\mathbf{W}^{(t)}\|_2^2$. Given Lemma E.1, we have following coefficient update rule for any $t \geq 0$, $i \in [n]$ and $r \in S_i^{(0)}$:

$$\bar{\rho}_{y_i, r, i}^{(t+1)} = \bar{\rho}_{y_i, r, i}^{(t)} + \frac{\eta}{nm} \cdot |\ell_i^{(t)}| \cdot \|\mathbf{x}_i\|_2^2, \quad (\text{F.1})$$

where

$$|\ell_i^{(t)}| = \frac{1}{1 + \exp\{F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i) - F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i)\}}.$$

Now we are ready to prove the first bullet of Theorem 4.3.

Lemma F.1. For two-layer ReLU neural network defined in (E.1), under the same condition as Theorem 4.3, the stable rank of $\mathbf{W}_j^{(t)}$ satisfies the following property:

$$\limsup_{t \rightarrow \infty} \frac{\|\mathbf{W}_j^{(t)}\|_F^2}{\|\mathbf{W}_j^{(t)}\|_2^2} \leq C,$$

where $C = \Theta(1)$ is a constant.

Proof of Lemma F.1. By decomposition (5.1), we have

$$\mathbf{w}_{j, r}^{(t)} - \mathbf{w}_{j, r}^{(0)} = \sum_{i=1}^n \rho_{j, r, i}^{(t)} \cdot \|\mathbf{x}_i\|_2^{-2} \cdot \mathbf{x}_i = \left[\rho_{j, r, 1}^{(t)} \|\mathbf{x}_1\|_2^{-2}, \dots, \rho_{j, r, n}^{(t)} \|\mathbf{x}_n\|_2^{-2} \right] \cdot \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix},$$

and

$$\mathbf{W}_j^{(t)} - \mathbf{W}_j^{(0)} = \underbrace{\begin{bmatrix} \rho_{j, 1, 1}^{(t)} \|\mathbf{x}_1\|_2^{-2} & \rho_{j, 1, 2}^{(t)} \|\mathbf{x}_1\|_2^{-2} & \cdots & \rho_{j, 1, n}^{(t)} \|\mathbf{x}_n\|_2^{-2} \\ \rho_{j, 2, 1}^{(t)} \|\mathbf{x}_1\|_2^{-2} & \rho_{j, 2, 2}^{(t)} \|\mathbf{x}_1\|_2^{-2} & \cdots & \rho_{j, 2, n}^{(t)} \|\mathbf{x}_n\|_2^{-2} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{j, m, 1}^{(t)} \|\mathbf{x}_1\|_2^{-2} & \rho_{j, m, 2}^{(t)} \|\mathbf{x}_1\|_2^{-2} & \cdots & \rho_{j, m, n}^{(t)} \|\mathbf{x}_n\|_2^{-2} \end{bmatrix}}_{\mathbf{A}_t} \cdot \underbrace{\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix}}_{:= \mathbf{X}}.$$

Let $\mathbf{a}_i(t)^\top$ be the i -th column of \mathbf{A}_t . It follows that

$$\begin{aligned} \|\mathbf{W}_j^{(t)} - \mathbf{W}_j^{(0)}\|_F^2 &= \text{Tr}(\mathbf{A}_t \mathbf{X} \mathbf{X}^\top \mathbf{A}_t^\top) \\ &= \text{Tr} \left(\left(\sum_{i=1}^n \mathbf{a}_i(t)^\top \mathbf{x}_i \right) \left(\sum_{i=1}^n \mathbf{x}_i^\top \mathbf{a}_i(t) \right) \right) \\ &= \text{Tr} \left(\sum_{i=1}^n \sum_{i'=1}^n \mathbf{a}_i(t)^\top \mathbf{x}_i \mathbf{x}_{i'}^\top \mathbf{a}_{i'}(t) \right) \end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^n \sum_{i'=1}^n \langle \mathbf{x}_i, \mathbf{x}_{i'} \rangle \cdot \text{Tr}(\mathbf{a}_i(t)^\top \mathbf{a}_{i'}(t)) \\
&= \sum_{i=1}^n \|\mathbf{x}_i\|_2^2 \cdot \text{Tr}(\mathbf{a}_i(t)^\top \mathbf{a}_i(t)) + \sum_{i \neq i'} \langle \mathbf{x}_i, \mathbf{x}_{i'} \rangle \cdot \text{Tr}(\mathbf{a}_i(t)^\top \mathbf{a}_{i'}(t)) \\
&\leq R_{\max}^2 \sum_{i=1}^n \text{Tr}(\mathbf{a}_i(t)^\top \mathbf{a}_i(t)) + p \sum_{i \neq i'} |\text{Tr}(\mathbf{a}_i(t)^\top \mathbf{a}_{i'}(t))|,
\end{aligned}$$

734 and

$$\begin{aligned}
\text{Tr}(\mathbf{a}_i(t)^\top \mathbf{a}_i(t)) &= \sum_{r=1}^m ([\mathbf{a}_i(t)]_r)^2 = \sum_{r=1}^m (\rho_{j,r,i}^{(t)} \|\mathbf{x}_i\|_2^{-2})^2 \leq \sum_{r=1}^m (\rho_{j,r,i}^{(t)})^2 R_{\min}^{-4} \leq C m R_{\min}^{-4} (\log(t))^2, \\
|\text{Tr}(\mathbf{a}_i(t)^\top \mathbf{a}_{i'}(t))| &\leq \sum_{r=1}^m |[\mathbf{a}_i(t)]_r [\mathbf{a}_{i'}(t)]_r| = \sum_{r=1}^m |\rho_{j,r,i}^{(t)} \rho_{j,r,i'}^{(t)}| \|\mathbf{x}_i\|_2^{-2} \|\mathbf{x}_{i'}\|_2^{-2} \leq C m R_{\min}^{-4} (\log(t))^2.
\end{aligned}$$

735 Accordingly, we have

$$\begin{aligned}
\|\mathbf{W}_j^{(t)} - \mathbf{W}_j^{(0)}\|_F &\leq C m n R_{\max}^2 R_{\min}^{-4} (\log(t))^2 + C m n^2 p R_{\min}^{-4} (\log(t))^2 \\
&= C m n R_{\max}^2 R_{\min}^{-4} (1 + R_{\max}^{-2} n p) (\log(t))^2 \\
&\leq C' m n R_{\max}^2 R_{\min}^{-4} (\log(t))^2.
\end{aligned}$$

736 On the other hand, we will give an lower bound for $\|\mathbf{W}_j^{(t)} - \mathbf{W}_j^{(0)}\|_2$.

$$\begin{aligned}
\|\mathbf{W}_j^{(t)} - \mathbf{W}_j^{(0)}\|_2 &= \max_{\mathbf{y} \in S^{d-1}} \|(\mathbf{W}_j^{(t)} - \mathbf{W}_j^{(0)})\mathbf{y}\|_2 \\
&\geq \frac{\|(\mathbf{W}_j^{(t)} - \mathbf{W}_j^{(0)})\mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{1}_n\|_2}{\|\mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{1}_n\|_2} \\
&= \frac{\|\mathbf{A}_t \mathbf{1}_n\|_2}{\|\mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{1}_n\|_2}.
\end{aligned}$$

737 We first provide a lower bound for $\|\mathbf{A}_t \mathbf{1}_n\|_2$. Note that

$$\mathbf{A}_t \mathbf{1}_n = \begin{bmatrix} \sum_{i=1}^n \rho_{j,1,i}^{(t)} \|\mathbf{x}_i\|_2^{-2} \\ \vdots \\ \sum_{i=1}^n \rho_{j,m,i}^{(t)} \|\mathbf{x}_i\|_2^{-2} \end{bmatrix},$$

738 we need to bound $\sum_{i=1}^n \rho_{j,r,i}^{(t)} \|\mathbf{x}_i\|_2^{-2}$, $r \in [m]$. By Corollary E.3, there exists time T such that for
739 any $t \geq T$, $|\rho_{y_i, r, i}^{(t)}| \leq c R_{\min}^{-2} p n \bar{\rho}_{y_i, r, i}^{(t)}$, $\forall i \in S_r^{(0)}$ and $\forall i' \in [n]$ with $y_{i'} = -y_i$. Therefore, we have
740 for $t \geq T$ that

$$\begin{aligned}
\sum_{r=1}^m \sum_{i=1}^n \rho_{j,r,i}^{(t)} \|\mathbf{x}_i\|_2^{-2} &= \sum_{r=1}^m \left(\sum_{i \in S_j} \bar{\rho}_{j,r,i}^{(t)} \|\mathbf{x}_i\|_2^{-2} + \sum_{i \in S_{-j}} \rho_{j,r,i}^{(t)} \|\mathbf{x}_i\|_2^{-2} \right) \\
&= \sum_{i \in S_j} \sum_{r=1}^m \bar{\rho}_{y_i, r, i}^{(t)} \|\mathbf{x}_i\|_2^{-2} + \sum_{i \in S_{-j}} \sum_{r=1}^m \rho_{-y_i, r, i}^{(t)} \|\mathbf{x}_i\|_2^{-2} \\
&\geq \sum_{i \in S_j} \sum_{r=1}^m \bar{\rho}_{y_i, r, i}^{(t)} R_{\max}^{-2} + \sum_{i \in S_{-j}} \sum_{r=1}^m \rho_{-y_i, r, i}^{(t)} R_{\min}^{-2} \\
&\geq \sum_{i \in S_j} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} R_{\max}^{-2} - \frac{m |S_{-j}|}{|S_j|} \sum_{i \in S_j} \frac{c R_{\min}^{-2} p n}{|S_i^{(0)}|} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} R_{\min}^{-2}
\end{aligned}$$

$$\begin{aligned}
&\geq \sum_{i \in S_j} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} R_{\max}^{-2} - \frac{m|S_{-j}| \cdot cR_{\min}^{-4} R_{\max}^2 pn}{|S_j| \cdot \min_{i \in S_j} |S_i^{(0)}|} \sum_{i \in S_j} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} R_{\max}^{-2} \\
&\geq \frac{R_{\max}^{-2}}{2} \sum_{i \in S_j} \sum_{r \in S_i^{(0)}} \bar{\rho}_{j, r, i}^{(t)}, \tag{F.2}
\end{aligned}$$

741 where the second inequality is by Corollary E.3 and hence

$$\begin{aligned}
|\rho_{-y_i, r, i}^{(t)}| &\leq cR_{\min}^{-2} pn \bar{\rho}_{y_i, r', i'}^{(t)}, \forall r, r' \in [m], \forall i, i' \in [n], \\
|\rho_{-y_i, r, i}^{(t)}| &\leq \frac{cR_{\min}^{-2} pn}{|S_{i'}^{(0)}|} \sum_{r \in S_{i'}^{(0)}} \bar{\rho}_{y_i, r, i'}^{(t)}, \forall r \in [m], \forall i, i' \in [n], \\
|\rho_{-y_i, r, i}^{(t)}| &\leq \frac{1}{|S_j|} \sum_{i \in S_j} \frac{cR_{\min}^{-2} pn}{|S_i^{(0)}|} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)}, \forall r \in [m], \forall i \in [n],
\end{aligned}$$

742 the last inequality is by

$$\frac{m|S_{-j}| \cdot cR_{\min}^{-4} R_{\max}^2 pn}{|S_j| \cdot \min_{i \in S_j} |S_i^{(0)}|} \leq c'R_{\min}^{-4} R_{\max}^2 pn \cdot \frac{|S_{-j}|}{|S_j|} \leq \frac{1}{2}.$$

743 Then, we have for $t \geq T$ that

$$\begin{aligned}
\|\mathbf{A}_t \mathbf{1}_n\|_2 &= \sqrt{\sum_{r=1}^m \left(\sum_{i=1}^n \rho_{j, r, i}^{(t)} \|\mathbf{x}_i\|_2^{-2} \right)^2} \\
&\geq \left| \sum_{r=1}^m \sum_{i=1}^n \rho_{j, r, i}^{(t)} \|\mathbf{x}_i\|_2^{-2} / \sqrt{m} \right| \\
&\geq \frac{R_{\max}^{-2}}{2\sqrt{m}} \sum_{i \in S_j} \sum_{r \in S_i^{(0)}} \bar{\rho}_{j, r, i}^{(t)} \\
&\geq \frac{R_{\max}^{-2} |S_j| |S_i^{(0)}|}{2\sqrt{m}} \cdot \log \left(1 + \frac{\eta R_{\min}^2 e^{-\beta}}{2c'nm} \cdot t \right) \\
&\geq CR_{\max}^{-2} \sqrt{mn} \log(t)
\end{aligned}$$

744 where the second inequality is by (F.2); the third inequality is by the second bullet of Lemma E.1.

745 For $\|\mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{1}_n\|_2$, we have

$$\begin{aligned}
\|\mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{1}_n\|_2 &= \sqrt{\mathbf{1}_n^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{X}\mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{1}_n} \\
&= \sqrt{\mathbf{1}_n^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{1}_n} \\
&\leq \frac{\|\mathbf{1}_n\|_2}{\sqrt{\lambda_{\min}(\mathbf{X}\mathbf{X}^\top)}}
\end{aligned}$$

746 By the Gershgorin circle theorem, we know that $\lambda_{\min}(\mathbf{X}\mathbf{X}^\top)$ lies within at least one of the Gershgorin
747 discs $D((\mathbf{X}\mathbf{X}^\top)_{ii}, R_i)$, $i \in [n]$ where $D((\mathbf{X}\mathbf{X}^\top)_{ii}, R_i)$ is a closed disc centered at $(\mathbf{X}\mathbf{X}^\top)_{ii} =$
748 $\|\mathbf{x}_i\|_2^2$ with radius $R_i = \sum_{i' \neq i} |(\mathbf{X}\mathbf{X}^\top)_{ii'}| = \sum_{i' \neq i} |\langle \mathbf{x}_i, \mathbf{x}_{i'} \rangle|$. Assume $\lambda_{\min}(\mathbf{X}\mathbf{X}^\top)$ lies within
749 $D((\mathbf{X}\mathbf{X}^\top)_{ii}, R_i)$, then we can get following lower bound for $\lambda_{\min}(\mathbf{X}\mathbf{X}^\top)$:

$$\lambda_{\min}(\mathbf{X}\mathbf{X}^\top) \geq \|\mathbf{x}_i\|_2^2 - \sum_{i' \neq i} |\langle \mathbf{x}_i, \mathbf{x}_{i'} \rangle| \geq R_{\min}^2 - np = (1 - R_{\min}^{-2} np) R_{\min}^2 \geq R_{\min}^2 / 2.$$

Therefore, we have

$$\|\mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{1}_n\|_2 \leq \frac{\|\mathbf{1}_n\|_2}{\sqrt{\lambda_{\min}(\mathbf{X}\mathbf{X}^\top)}} \leq \frac{\sqrt{2n}}{R_{\min}}.$$

Accordingly,

$$\|\mathbf{W}_j^{(t)} - \mathbf{W}_j^{(0)}\|_2 \geq \frac{\|\mathbf{A}_t \mathbf{1}_n\|_2}{\|\mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{1}_n\|_2} \geq \frac{C R_{\max}^{-2} \sqrt{mn} \log(t)}{\sqrt{2n}/R_{\min}} = C'' R_{\max}^{-2} R_{\min} \sqrt{mn} \log(t).$$

Therefore, we have for $t \geq T$ that

$$\frac{\|\mathbf{W}_j^{(t)} - \mathbf{W}_j^{(0)}\|_F^2}{\|\mathbf{W}_j^{(t)} - \mathbf{W}_j^{(0)}\|_2^2} \leq \frac{C' mn R_{\max}^2 R_{\min}^{-4} (\log(t))^2}{C''^2 mn R_{\max}^2 R_{\min}^{-4} (\log(t))^2} \leq \frac{C' R_{\max}^6}{C''^2 R_{\min}^6},$$

which completes the proof. \square

Next, we will provide an example of training data satisfying the condition in Theorem 4.3 and prove that the stable rank of $\mathbf{W}_j^{(t)}$ trained by gradient descent using such data will converge to $2 \pm o(1)$.

We first provide the following lemma about the increasing rate of coefficients $\rho_{j,r,i}^{(t)}$.

Lemma F.2. If training data $\mathbf{x}_1, \dots, \mathbf{x}_n$ are mutually orthogonal, the activation pattern after time $T = C\eta^{-1} R_{\min}^{-2}/nm$ is determined by the activation pattern at initialization, that is,

$$\begin{aligned} \langle \mathbf{w}_{y_i,r}^{(t)}, \mathbf{x}_i \rangle &\geq 0, & \text{if } \langle \mathbf{w}_{y_i,r}^{(0)}, \mathbf{x}_i \rangle &\geq 0, \\ \langle \mathbf{w}_{y_i,r}^{(t)}, \mathbf{x}_i \rangle &< 0, & \text{if } \langle \mathbf{w}_{y_i,r}^{(0)}, \mathbf{x}_i \rangle &< 0, \\ \langle \mathbf{w}_{-y_i,r}^{(t)}, \mathbf{x}_i \rangle &< 0, & \text{if } \langle \mathbf{w}_{-y_i,r}^{(0)}, \mathbf{x}_i \rangle &\geq 0, \\ \langle \mathbf{w}_{-y_i,r}^{(t)}, \mathbf{x}_i \rangle &< 0, & \text{if } \langle \mathbf{w}_{-y_i,r}^{(0)}, \mathbf{x}_i \rangle &< 0, \end{aligned}$$

for any time $t \geq T$. Besides, $\rho_{j,r,i}^{(t)}$ satisfy the following properties:

$$\begin{aligned} \bar{\rho}_{y_i,r,i}^{(t)} &= 0, \forall t \geq 0, & \text{if } \langle \mathbf{w}_{y_i,r}^{(0)}, \mathbf{x}_i \rangle &< 0, \\ \bar{\rho}_{y_i,r,i}^{(t)} &= \bar{\rho}_{y_i,r,i}^{(T)}, \forall t \geq T, & \text{if } \langle \mathbf{w}_{-y_i,r}^{(0)}, \mathbf{x}_i \rangle &\geq 0, \\ \lim_{t \rightarrow \infty} \bar{\rho}_{y_i,r,i}^{(t)} / \log t &= m/|S_i^{(0)}|, & \text{if } \langle \mathbf{w}_{y_i,r}^{(0)}, \mathbf{x}_i \rangle &\geq 0. \end{aligned}$$

Proof of Lemma F.2. Part 1. For if $\langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle < 0$, we first prove by induction that

$$\rho_{j,r,i}^{(t)} = 0, \langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle = \langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle < 0, \forall t \geq 0. \quad (\text{F.3})$$

The result is obvious at $t = 0$ as all the coefficients are zero. Suppose that there exists \tilde{t} such that (F.3) holds for all time $0 \leq t \leq \tilde{t} - 1$. We aim to prove that (F.3) also holds for $t = \tilde{t}$. Recall that by (5.4), (5.5) and with (F.3) at time $\tilde{t} - 1$, we have

$$\begin{aligned} \bar{\rho}_{j,r,i}^{(\tilde{t})} &= \bar{\rho}_{j,r,i}^{(\tilde{t}-1)} - \frac{\eta}{nm} \cdot \ell_i'(\tilde{t}-1) \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(\tilde{t}-1)}, \mathbf{x}_i \rangle) \cdot \|\mathbf{x}_i\|_2^2 \cdot \mathbb{1}(y_i = j) = \bar{\rho}_{j,r,i}^{(\tilde{t}-1)} = 0, \\ \underline{\rho}_{j,r,i}^{(\tilde{t})} &= \underline{\rho}_{j,r,i}^{(\tilde{t}-1)} + \frac{\eta}{nm} \cdot \ell_i'(\tilde{t}-1) \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(\tilde{t}-1)}, \mathbf{x}_i \rangle) \cdot \|\mathbf{x}_i\|_2^2 \cdot \mathbb{1}(y_i = -j) = \underline{\rho}_{j,r,i}^{(\tilde{t}-1)} = 0. \end{aligned}$$

By (5.1) and the orthogonality of training data, we can get

$$\langle \mathbf{w}_{j,r}^{(\tilde{t})}, \mathbf{x}_i \rangle = \langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle + \sum_{i'=1}^n \rho_{j,r,i'}^{(\tilde{t})} \|\mathbf{x}_{i'}\|_2^{-2} \cdot \langle \mathbf{x}_{i'}, \mathbf{x}_i \rangle = \langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle + \rho_{j,r,i}^{(\tilde{t})} = \langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle < 0.$$

Therefore, (F.3) holds at time \tilde{t} , which completes the induction.

766 For if $\langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle \geq 0$ and $j = y_i$, we will next prove by induction that

$$\bar{\rho}_{j,r,i}^{(t)} \geq 0, \langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle \geq \langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle \geq 0, \forall t \geq 0. \quad (\text{F.4})$$

767 The result is natural at $t = 0$. Suppose that there exists \tilde{t} such that (F.4) hold for all time $0 \leq t \leq \tilde{t} - 1$.
 768 By (5.4) and (F.4) at time $\tilde{t} - 1$, we have

$$\bar{\rho}_{j,r,i}^{(\tilde{t})} = \bar{\rho}_{j,r,i}^{(\tilde{t}-1)} + \frac{\eta}{nm} \cdot |\ell_i^{(\tilde{t}-1)}| \cdot \|\mathbf{x}_i\|_2^2 \geq \bar{\rho}_{j,r,i}^{(\tilde{t}-1)} \geq 0$$

769 and hence the orthogonality of training data, we can get

$$\langle \mathbf{w}_{j,r}^{(\tilde{t})}, \mathbf{x}_i \rangle = \langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle + \rho_{j,r,i}^{(\tilde{t})} = \langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle \geq 0.$$

770 Therefore, (F.4) hold at time \tilde{t} , which completes the induction.

771 For if $\langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle \geq 0$ and $j = -y_i$, we first show that under the same condition as Theorem 4.3 it
 772 holds that

$$\rho_{j,r,i}^{(T)} < -\langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle.$$

773 Since $T = C\eta^{-1}R_{\min}^{-2}/nm$, we have $\max_{j,r,i}\{|\rho_{j,r,i}^{(t)}|\} = O(1)$ for $t \leq T$. Therefore, we know that
 774 $F_{+1}(\mathbf{W}_{+1}^{(t)}, \mathbf{x}_i), F_{-1}(\mathbf{W}_{-1}^{(t)}, \mathbf{x}_i) = O(1)$. Thus there exists a positive constant c such that $-\ell_i^{(t)} \geq c$
 775 for all $i \in [n]$. Here we use the method of proof by contradiction. Assume $\rho_{j,r,i}^{(T)} \geq -\langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle$.
 776 Since $\rho_{j,r,i}^{(T)} \leq \rho_{j,r,i}^{(t)}$ for $0 \leq t \leq T$ which can be seen from (5.5), we have $\rho_{j,r,i}^{(t)} \geq -\langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle$ for
 777 all $t \leq T$. Then, we can get

$$\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle = \langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle + \rho_{j,r,i}^{(t)} \geq 0, \forall t \leq T.$$

778 Therefore, by the non-negativeness of $\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle$ and (5.5), we can get

$$|\rho_{j,r,i}^{(t+1)}| = |\rho_{j,r,i}^{(t)}| + \frac{\eta}{nm} \cdot |\ell_i^{(t)}| \cdot \|\mathbf{x}_i\|_2^2 \geq |\rho_{j,r,i}^{(t)}| + \frac{c\eta\|\mathbf{x}_i\|_2^2}{nm}$$

779 and hence

$$|\rho_{j,r,i}^{(T)}| \geq \frac{c\eta\|\mathbf{x}_i\|_2^2 T}{nm} = cC\|\mathbf{x}_i\|_2^2 R_{\min}^{-2} \geq cC \geq \beta,$$

780 which is a contradiction. Therefore, $\rho_{j,r,i}^{(T)} < -\langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle$. By (5.5), we have $\rho_{j,r,i}^{(t)} \leq \rho_{j,r,i}^{(T)} <$
 781 $-\langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle$ for $t \geq T$. Therefore,

$$\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle = \langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle + \rho_{j,r,i}^{(t)} < 0, \forall t \geq T.$$

782 Plugging this into (5.5) gives us

$$\rho_{j,r,i}^{(t+1)} = \rho_{j,r,i}^{(t)} + \frac{\eta}{nm} \cdot \ell_i^{(t)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \mathbf{x}_i \rangle) \cdot \|\mathbf{x}_i\|_2^2 = \rho_{j,r,i}^{(t)}, \forall t \geq T.$$

783 This completes the proof of the first half of the lemma about the activation pattern as well as the first
 784 two properties of $\rho_{j,r,i}^{(t)}$.

785 **Part 2.** Now we will show that

$$\lim_{t \rightarrow \infty} \bar{\rho}_{y_i,r,i}^{(t)} / \log t = m/|S_i^{(0)}|, \quad (\text{F.5})$$

786 if $\langle \mathbf{w}_{y_i, r}^{(t)}, \mathbf{x}_i \rangle \geq 0$. By the activation pattern, we can get

$$\begin{aligned}
\bar{\rho}_{y_i, r, i}^{(t+1)} &= \bar{\rho}_{y_i, r, i}^{(t)} + \frac{\eta}{nm} \cdot |\ell_i'^{(t)}| \cdot \|\mathbf{x}_i\|_2^2, & \forall t \geq 0, & \text{ for } \langle \mathbf{w}_{y_i, i}^{(t)}, \mathbf{x}_i \rangle \geq 0, \\
\bar{\rho}_{y_i, r, i}^{(t)} &= 0, & \forall t \geq 0, & \text{ for } \langle \mathbf{w}_{y_i, i}^{(t)}, \mathbf{x}_i \rangle < 0, \\
\rho_{-y_i, r, i}^{(t)} &= \rho_{-y_i, r, i}^{(T)}, & \forall t \geq 0, & \text{ for } \langle \mathbf{w}_{-y_i, i}^{(t)}, \mathbf{x}_i \rangle \geq 0, \\
\rho_{-y_i, r, i}^{(t)} &= 0, & \forall t \geq 0, & \text{ for } \langle \mathbf{w}_{-y_i, i}^{(t)}, \mathbf{x}_i \rangle < 0.
\end{aligned} \tag{F.6}$$

787 Given this activation pattern, we can get for $t \geq 0$ that

$$\begin{aligned}
F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i) - F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i) &= \frac{1}{m} \sum_{r=1}^m \sigma(\langle \mathbf{w}_{y_i, r}^{(t)}, \mathbf{x}_i \rangle) - \frac{1}{m} \sum_{r=1}^m \sigma(\langle \mathbf{w}_{-y_i, r}^{(t)}, \mathbf{x}_i \rangle) \\
&\leq \frac{1}{m} \sum_{r \in S_i^{(0)}} \langle \mathbf{w}_{y_i, r}^{(t)}, \mathbf{x}_i \rangle \\
&= \frac{1}{m} \sum_{r \in S_i^{(0)}} [\langle \mathbf{w}_{y_i, r}^{(0)}, \mathbf{x}_i \rangle + \bar{\rho}_{y_i, r, i}^{(t)}] \\
&\leq \frac{1}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} + \beta,
\end{aligned}$$

788 and

$$\begin{aligned}
F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i) - F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i) &= \frac{1}{m} \sum_{r=1}^m \sigma(\langle \mathbf{w}_{y_i, r}^{(t)}, \mathbf{x}_i \rangle) - \frac{1}{m} \sum_{r=1}^m \sigma(\langle \mathbf{w}_{-y_i, r}^{(t)}, \mathbf{x}_i \rangle) \\
&\geq \frac{1}{m} \sum_{r \in S_i^{(0)}} \langle \mathbf{w}_{y_i, r}^{(t)}, \mathbf{x}_i \rangle - \beta \\
&= \frac{1}{m} \sum_{r \in S_i^{(0)}} [\langle \mathbf{w}_{y_i, r}^{(0)}, \mathbf{x}_i \rangle + \bar{\rho}_{y_i, r, i}^{(t)}] - \beta \\
&\geq \frac{1}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} - 2\beta.
\end{aligned}$$

789 Therefore, we can get following upper and lower bounds for $|\ell_i'^{(t)}|$:

$$\begin{aligned}
|\ell_i'^{(t)}| &\leq \exp \left(-\frac{1}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} + 2\beta \right), \forall t \geq 0, \\
|\ell_i'^{(t)}| &\geq \frac{1}{2} \exp \left(-\frac{1}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} - \beta \right), \forall t \geq 0.
\end{aligned}$$

790 And it follows that

$$\begin{aligned}
\frac{1}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t+1)} &\leq \frac{1}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} + \frac{\eta \|\mathbf{x}_i\|_2^2 e^{2\beta}}{nm} \cdot \exp \left(-\frac{1}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} \right), \forall t \geq 0, \\
\frac{1}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t+1)} &\geq \frac{1}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} + \frac{\eta \|\mathbf{x}_i\|_2^2 e^{-\beta}}{2nm} \cdot \exp \left(-\frac{1}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} \right), \forall t \geq 0.
\end{aligned}$$

By leveraging Lemma H.1 as well as Lemma H.2 and taking

$$x_t = \frac{1}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)},$$

we can get

$$\begin{aligned} \frac{1}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} &\leq \log \left(1 + \frac{\eta \|\mathbf{x}_i\|_2^2 e^{2\beta}}{nm} \exp \left(\frac{\eta \|\mathbf{x}_i\|_2^2 e^{2\beta}}{nm} \right) \cdot t \right) \leq \log \left(1 + \frac{2\eta \|\mathbf{x}_i\|_2^2 e^{-\beta}}{nm} \cdot t \right), \\ \frac{1}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} &\geq \log \left(1 + \frac{\eta \|\mathbf{x}_i\|_2^2 e^{-\beta}}{2nm} \cdot t \right). \end{aligned}$$

Therefore, we have

$$\lim_{t \rightarrow \infty} \frac{1}{m} \sum_{r \in S_i^{(0)}} \bar{\rho}_{y_i, r, i}^{(t)} / \log t = 1.$$

Since $\bar{\rho}_{y_i, r, i}^{(t)} = \bar{\rho}_{y_i, r', i}^{(t)}$ for any $r \neq r' \in S_i^{(0)}$, we have

$$\lim_{t \rightarrow \infty} \bar{\rho}_{y_i, r, i}^{(t)} / \log t = m / |S_i^{(0)}|,$$

which completes the proof. \square

Lemma F.3. For two-layer ReLU neural network defined in (E.1), there exists mutually orthogonal data $\mathbf{x}_1, \dots, \mathbf{x}_n$ such that stable rank of $\mathbf{W}_j^{(t)}$ will converge to $2 \pm o(1)$.

Proof of Lemma F.3. By (5.1), we have

$$\mathbf{w}_{j, r}^{(t)} = \mathbf{w}_{j, r}^{(0)} + \underbrace{\sum_{i=1}^n \rho_{j, r, i}^{(t)} \cdot \|\mathbf{x}_i\|_2^{-2} \cdot \mathbf{x}_i}_{:= \mathbf{v}_{j, r}^{(t)}}.$$

Given the definition of $\mathbf{v}_{j, r}^{(t)}$, we have the following representation of $\mathbf{v}_{j, r}^{(t)}$ and $\mathbf{V}_j^{(t)}$.

$$\mathbf{v}_{j, r}^{(t)} = [\rho_{j, r, 1}^{(t)} \cdot \|\mathbf{x}_1\|_2^{-2} \cdots \rho_{j, r, n}^{(t)} \cdot \|\mathbf{x}_n\|_2^{-2}] \cdot \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix},$$

$$\mathbf{V}_j^{(t)} = \begin{bmatrix} \rho_{j, r, 1}^{(t)} \cdot \|\mathbf{x}_1\|_2^{-2} & \cdots & \rho_{j, r, n}^{(t)} \cdot \|\mathbf{x}_n\|_2^{-2} \\ \vdots & \ddots & \vdots \\ \rho_{j, m, 1}^{(t)} \cdot \|\mathbf{x}_1\|_2^{-2} & \cdots & \rho_{j, m, n}^{(t)} \cdot \|\mathbf{x}_n\|_2^{-2} \end{bmatrix} \cdot \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix}.$$

Assume n is an even number and $\mathbf{x}_1, \dots, \mathbf{x}_{n/2}$ are with label $+1$ while $\mathbf{x}_{(n/2)+1}, \dots, \mathbf{x}_n$ are with label -1 . And we take $\mathbf{x}_1, \dots, \mathbf{x}_n$ as $\mathbf{e}_1, \dots, \mathbf{e}_n$. Given Lemma F.2, $\mathbf{W}_j^{(t)} = \mathbf{W}_j^{(0)} + \mathbf{V}_j^{(t)}$ and such selection of training data, we have

$$\begin{aligned} \lim_{t \rightarrow \infty} \frac{\mathbf{W}_{+1}^{(t)}}{\log t} &= [\mathbf{A}_{m \times (n/2)}, \mathbf{0}_{m \times (n/2)}] \cdot [\mathbf{I}_n, \mathbf{0}_{n \times (d-n)}] = [\mathbf{A}_{m \times (n/2)}, \mathbf{0}_{m \times (d-(n/2))}], \\ \lim_{t \rightarrow \infty} \frac{\mathbf{W}_{-1}^{(t)}}{\log t} &= [\mathbf{0}_{m \times (n/2)}, \mathbf{B}_{m \times (n/2)}] \cdot [\mathbf{I}_n, \mathbf{0}_{n \times (d-n)}] = [\mathbf{0}_{m \times (n/2)}, \mathbf{B}_{m \times (n/2)}, \mathbf{0}_{n \times (d-n)}], \\ \lim_{t \rightarrow \infty} \frac{\mathbf{W}^{(t)}}{\log t} &= \begin{bmatrix} \mathbf{A}_{m \times (n/2)} & \mathbf{0}_{m \times (n/2)} & \mathbf{0}_{n \times (d-n)} \\ \mathbf{0}_{m \times (n/2)} & \mathbf{B}_{m \times (n/2)} & \mathbf{0}_{n \times (d-n)} \end{bmatrix}. \end{aligned}$$

804 where

$$\begin{aligned}\mathbf{A}_{m \times (n/2)} &= \underbrace{\begin{bmatrix} \mathbb{1}[\langle \mathbf{w}_{+1,1}^{(0)}, \mathbf{x}_1 \rangle \geq 0] & \cdots & \mathbb{1}[\langle \mathbf{w}_{+1,1}^{(0)}, \mathbf{x}_{n/2} \rangle \geq 0] \\ \vdots & \ddots & \vdots \\ \mathbb{1}[\langle \mathbf{w}_{+1,m}^{(0)}, \mathbf{x}_1 \rangle \geq 0] & \cdots & \mathbb{1}[\langle \mathbf{w}_{+1,m}^{(0)}, \mathbf{x}_{n/2} \rangle \geq 0] \end{bmatrix}}_{:=\mathbf{C}_{m \times (n/2)}} \cdot \text{diag} \begin{bmatrix} m/|S_1^{(0)}| \\ \vdots \\ m/|S_{n/2}^{(0)}| \end{bmatrix}, \\ \mathbf{B}_{m \times (n/2)} &= \underbrace{\begin{bmatrix} \mathbb{1}[\langle \mathbf{w}_{-1,1}^{(0)}, \mathbf{x}_{(n/2)+1} \rangle \geq 0] & \cdots & \mathbb{1}[\langle \mathbf{w}_{-1,1}^{(0)}, \mathbf{x}_n \rangle \geq 0] \\ \vdots & \ddots & \vdots \\ \mathbb{1}[\langle \mathbf{w}_{-1,m}^{(0)}, \mathbf{x}_{(n/2)+1} \rangle \geq 0] & \cdots & \mathbb{1}[\langle \mathbf{w}_{-1,m}^{(0)}, \mathbf{x}_n \rangle \geq 0] \end{bmatrix}}_{:=\mathbf{D}_{m \times (n/2)}} \cdot \text{diag} \begin{bmatrix} m/|S_{(n/2)+1}^{(0)}| \\ \vdots \\ m/|S_n^{(0)}| \end{bmatrix}.\end{aligned}$$

805 Then, we can get the stable rank limits as follows:

$$\begin{aligned}\lim_{t \rightarrow \infty} \frac{\|\mathbf{W}_{+1}^{(t)}\|_F^2}{\|\mathbf{W}_{+1}^{(t)}\|_2^2} &= \frac{\|\mathbf{A}\|_F^2}{\|\mathbf{A}\|_2^2}, \\ \lim_{t \rightarrow \infty} \frac{\|\mathbf{W}_{-1}^{(t)}\|_F^2}{\|\mathbf{W}_{-1}^{(t)}\|_2^2} &= \frac{\|\mathbf{B}\|_F^2}{\|\mathbf{B}\|_2^2}, \\ \lim_{t \rightarrow \infty} \frac{\|\mathbf{W}^{(t)}\|_F^2}{\|\mathbf{W}^{(t)}\|_2^2} &= \frac{\|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2}{(\max\{\|\mathbf{A}\|_2, \|\mathbf{B}\|_2\})^2}.\end{aligned}$$

806 Since $\mathbf{x}_1 = \mathbf{e}_1, \dots, \mathbf{x}_n = \mathbf{e}_n$, we can get

$$\mathbb{1}[\langle \mathbf{w}_{j,r}^{(0)}, \mathbf{x}_i \rangle \geq 0] = \mathbb{1}[[\mathbf{w}_{j,r}^{(0)}]_i \geq 0].$$

807 Therefore, the entries of matrix \mathbf{C} and matrix \mathbf{D} can be regarded as i.i.d. random variables taking 0
808 or 1 with equal probability. For $\|\mathbf{A}\|_F$ and $\|\mathbf{B}\|_F$, we have

$$\begin{aligned}\|\mathbf{A}\|_F^2 &= \sum_{r=1}^m \sum_{i=1}^{n/2} \mathbb{1}[[\mathbf{w}_{+1,r}^{(0)}]_i \geq 0] \cdot (m/|S_i^{(0)}|)^2, \\ \|\mathbf{B}\|_F^2 &= \sum_{r=1}^m \sum_{i=(n/2)+1}^n \mathbb{1}[[\mathbf{w}_{-1,r}^{(0)}]_i \geq 0] \cdot (m/|S_i^{(0)}|)^2.\end{aligned}$$

809 By Lemma A.2, we have with probability at least $1 - \delta$ that $0.4m \leq |S_i^{(0)}| \leq 0.6m$. By Hoeffding's
810 inequality, we have with probability at least $1 - 2\delta$ that

$$\begin{aligned}\left| \|\mathbf{A}\|_F^2 - \frac{m}{2} \sum_{i=1}^{n/2} (m/|S_i^{(0)}|)^2 \right| &\leq \sqrt{\frac{m \log(2/\delta)}{2} \sum_{i=1}^{n/2} (m/|S_i^{(0)}|)^4} \leq \sqrt{\frac{625mn \log(2/\delta)}{32}}, \\ \left| \|\mathbf{B}\|_F^2 - \frac{m}{2} \sum_{i=(n/2)+1}^n (m/|S_i^{(0)}|)^2 \right| &\leq \sqrt{\frac{m \log(2/\delta)}{2} \sum_{i=(n/2)+1}^n (m/|S_i^{(0)}|)^4} \leq \sqrt{\frac{625mn \log(2/\delta)}{32}}.\end{aligned}$$

812 Next, we estimate $\|\mathbf{A}\|_2$ and $\|\mathbf{B}\|_2$. Let $\mathbf{A} = \tilde{\mathbf{A}} + \mathbb{E}[\mathbf{A}]$ and $\mathbf{B} = \tilde{\mathbf{B}} + \mathbb{E}[\mathbf{B}]$. Assume \mathbf{G} be the
813 $m \times (n/2)$ matrix with all entries equal to $1/2$. Then,

$$\mathbb{E}[\mathbf{A}] = \mathbf{G} \cdot \underbrace{\text{diag} \begin{bmatrix} m/|S_1^{(0)}| \\ \vdots \\ m/|S_{n/2}^{(0)}| \end{bmatrix}}_{:=\mathbf{a}}, \mathbb{E}[\mathbf{B}] = \mathbf{G} \cdot \underbrace{\text{diag} \begin{bmatrix} m/|S_{(n/2)+1}^{(0)}| \\ \vdots \\ m/|S_n^{(0)}| \end{bmatrix}}_{:=\mathbf{b}}.$$

814 And the entries of matrix $\tilde{\mathbf{A}}$ and matrix $\tilde{\mathbf{B}}$ are independent, mean zero, sub-gaussian random variables.
 815 By Lemma H.3, we have with probability at least $1 - \delta$ that

$$\|\tilde{\mathbf{A}}\|_2 \leq \frac{C}{2}(\sqrt{m} + \sqrt{n} + \sqrt{\log(2/\delta)}),$$

816

$$\|\tilde{\mathbf{B}}\|_2 \leq \frac{C}{2}(\sqrt{m} + \sqrt{n} + \sqrt{\log(2/\delta)}),$$

817 where C is a constant. Let $\mathbf{1}_k$ denote the row vector with k entries equal to 1. Then for $\mathbb{E}[\mathbf{A}]$ and
 818 $\mathbb{E}[\mathbf{B}]$, we have

$$\begin{aligned} \|\mathbb{E}[\mathbf{A}]\|_2 &= \max_{\mathbf{x} \in S^{\frac{n}{2}-1}} \|\mathbf{G} \text{diag}(\mathbf{a}) \mathbf{x}\|_2 \\ &= \max_{\mathbf{x} \in S^{\frac{n}{2}-1}} \frac{1}{2} \|\mathbf{1}_m^\top \mathbf{1}_{\frac{n}{2}} \text{diag}(\mathbf{a}) \mathbf{x}\|_2 \\ &= \max_{\mathbf{x} \in S^{\frac{n}{2}-1}} \frac{\sqrt{m}}{2} |\mathbf{1}_{\frac{n}{2}} \text{diag}(\mathbf{a}) \mathbf{x}| \\ &= \max_{\mathbf{x} \in S^{\frac{n}{2}-1}} \frac{\sqrt{m}}{2} |\mathbf{a}^\top \mathbf{x}| \\ &= \frac{\sqrt{m} \|\mathbf{a}\|_2}{2}, \end{aligned}$$

819 and

$$\|\mathbb{E}[\mathbf{B}]\|_2 = \frac{\sqrt{m} \|\mathbf{b}\|_2}{2}.$$

820 By triangle inequality, we have

$$\begin{aligned} \|\mathbf{A}\|_2 &\geq \|\mathbf{C}\|_2 - \|\tilde{\mathbf{A}}\|_2 \geq (\sqrt{m} \|\mathbf{a}\|_2)/2 - \frac{C}{2}(\sqrt{m} + \sqrt{n} + \sqrt{\log(2/\delta)}), \\ \|\mathbf{A}\|_2 &\leq \|\mathbf{C}\|_2 + \|\tilde{\mathbf{A}}\|_2 \leq (\sqrt{m} \|\mathbf{a}\|_2)/2 + \frac{C}{2}(\sqrt{m} + \sqrt{n} + \sqrt{\log(2/\delta)}), \\ \|\mathbf{B}\|_2 &\geq \|\mathbf{C}\|_2 - \|\tilde{\mathbf{B}}\|_2 \geq (\sqrt{m} \|\mathbf{b}\|_2)/2 - \frac{C}{2}(\sqrt{m} + \sqrt{n} + \sqrt{\log(2/\delta)}), \\ \|\mathbf{B}\|_2 &\leq \|\mathbf{C}\|_2 + \|\tilde{\mathbf{B}}\|_2 \leq (\sqrt{m} \|\mathbf{b}\|_2)/2 + \frac{C}{2}(\sqrt{m} + \sqrt{n} + \sqrt{\log(2/\delta)}). \end{aligned}$$

821 Notice that $\|\mathbf{a}\|_2^2 = \Theta(n)$ and $\|\mathbf{b}\|_2^2 = \Theta(n)$, then with probability at least $1 - 2\delta$, we have

$$\begin{aligned} \frac{\|\mathbf{A}\|_F^2}{\|\mathbf{A}\|_2^2} &\leq \frac{m \|\mathbf{a}\|_2^2/2 + \sqrt{625mn \log(2/\delta)/32}}{(\sqrt{m} \|\mathbf{a}\|_2/2 - \frac{C}{2}(\sqrt{m} + \sqrt{n} + \sqrt{\log(2/\delta)}))^2} = 2 + o(1), \\ \frac{\|\mathbf{A}\|_F^2}{\|\mathbf{A}\|_2^2} &\geq \frac{m \|\mathbf{a}\|_2^2/2 - \sqrt{625mn \log(2/\delta)/32}}{(\sqrt{m} \|\mathbf{a}\|_2/2 + \frac{C}{2}(\sqrt{m} + \sqrt{n} + \sqrt{\log(2/\delta)}))^2} = 2 - o(1), \\ \frac{\|\mathbf{B}\|_F^2}{\|\mathbf{B}\|_2^2} &\geq \frac{m \|\mathbf{b}\|_2^2/2 + \sqrt{625mn \log(2/\delta)/32}}{(\sqrt{m} \|\mathbf{b}\|_2/2 - \frac{C}{2}(\sqrt{m} + \sqrt{n} + \sqrt{\log(2/\delta)}))^2} = 2 + o(1), \\ \frac{\|\mathbf{B}\|_F^2}{\|\mathbf{B}\|_2^2} &\geq \frac{m \|\mathbf{b}\|_2^2/2 - \sqrt{625mn \log(2/\delta)/32}}{(\sqrt{m} \|\mathbf{b}\|_2/2 + \frac{C}{2}(\sqrt{m} + \sqrt{n} + \sqrt{\log(2/\delta)}))^2} = 2 - o(1). \end{aligned}$$

822 This leads to

$$\begin{aligned} \lim_{t \rightarrow \infty} \frac{\|\mathbf{W}_{+1}^{(t)}\|_F^2}{\|\mathbf{W}_{+1}^{(t)}\|_2^2} &= \frac{\|\mathbf{A}\|_F^2}{\|\mathbf{A}\|_2^2} = 2 \pm o(1), \\ \lim_{t \rightarrow \infty} \frac{\|\mathbf{W}_{-1}^{(t)}\|_F^2}{\|\mathbf{W}_{-1}^{(t)}\|_2^2} &= \frac{\|\mathbf{B}\|_F^2}{\|\mathbf{B}\|_2^2} = 2 \pm o(1). \end{aligned}$$

823 For $\mathbf{W}^{(t)}$, we have the following lower bound

$$\begin{aligned} \frac{\|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2}{(\max\{\|\mathbf{A}\|_2, \|\mathbf{B}\|_2\})^2} &\geq \frac{m(\|\mathbf{a}\|_2^2 + \|\mathbf{b}\|_2^2)/2 - \sqrt{625mn \log(2/\delta)/8}}{(\sqrt{m} \max\{\|\mathbf{a}\|_2, \|\mathbf{b}\|_2\}/2 + \frac{C}{2}(\sqrt{m} + \sqrt{n} + \sqrt{\log(2/\delta)}))^2} \\ &\geq (2 - o(1)) \cdot \frac{\|\mathbf{a}\|_2^2 + \|\mathbf{b}\|_2^2}{(\max\{\|\mathbf{a}\|_2, \|\mathbf{b}\|_2\})^2} \\ &\geq \frac{16}{9} - o(1), \end{aligned}$$

824 where the third inequality is by $(5/3)\sqrt{n} \leq \|\mathbf{a}\|_2 \leq (5/2)\sqrt{n}$ and $(5/3)\sqrt{n} \leq \|\mathbf{b}\|_2 \leq (5/2)\sqrt{n}$
825 due to $0.4m \leq |S_i^{(0)}| \leq 0.6m$. And

$$\begin{aligned} \frac{\|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2}{(\max\{\|\mathbf{A}\|_2, \|\mathbf{B}\|_2\})^2} &\leq \frac{m(\|\mathbf{a}\|_2^2 + \|\mathbf{b}\|_2^2)/2 + \sqrt{625mn \log(2/\delta)/8}}{(\sqrt{m} \max\{\|\mathbf{a}\|_2, \|\mathbf{b}\|_2\}/2 - \frac{C}{2}(\sqrt{m} + \sqrt{n} + \sqrt{\log(2/\delta)}))^2} \\ &\leq (2 + o(1)) \cdot \frac{\|\mathbf{a}\|_2^2 + \|\mathbf{b}\|_2^2}{(\max\{\|\mathbf{a}\|_2, \|\mathbf{b}\|_2\})^2} \\ &\leq 9 + o(1), \end{aligned}$$

826 where the third inequality is by $(5/3)\sqrt{n} \leq \|\mathbf{a}\|_2 \leq (5/2)\sqrt{n}$ and $(5/3)\sqrt{n} \leq \|\mathbf{b}\|_2 \leq (5/2)\sqrt{n}$
827 due to $0.4m \leq |S_i^{(0)}| \leq 0.6m$. Therefore,

$$\lim_{t \rightarrow \infty} \frac{\|\mathbf{W}\|_F^2}{\|\mathbf{W}\|_2^2} = \frac{\|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2}{(\max\{\|\mathbf{A}\|_2, \|\mathbf{B}\|_2\})^2} \in [16/9 - o(1), 9 + o(1)].$$

828

□

829 G Loss Convergence of ReLU Network

830 We have the following convergence rate of the loss function.

831 **Lemma G.1.** For ReLU neural network defined in (E.1), for any $t \geq 1$, we have

$$L_S(\mathbf{W}^{(t)}) \leq \frac{e^{2\beta}}{n} \sum_{i=1}^n \left(1 + \frac{\eta |S_i^{(0)}| \|\mathbf{x}_i\|_2^2 e^{-\beta}}{2nm^2} \cdot t \right)^{-\frac{1-c_1^{-1}(1+c')R_{\min}^{-2}pn}{c'^2(1+c_1^{-1}R_{\min}^{-2}pn)}},$$

832 where c_1 is the same constant defined in Lemma E.1 and c' is a constant such that $m \leq c'|S_i^{(0)}|, \forall i \in$
833 $[n]$. This indicates that the training loss will converge with rate $O(t^{-\alpha})$ where $\alpha = c'^{-2}(1 - c_1^{-1}(1 +$
834 $c')R_{\min}^{-2}pn)/(1 + c_1^{-1}R_{\min}^{-2}pn)$ is a positive constant.

835 *Proof of Lemma G.1.* To establish an upper bound for $L_S(\mathbf{W}^{(t)})$, we need to first determine a lower
836 bound for the margin $y_i f(\mathbf{W}^{(t)}, \mathbf{x}_i)$. This can be accomplished by leveraging the increasing rate of
837 $\rho_{y_i, r, i}^{(t)}$ and the constant lower bound for $\rho_{y_i, r, i}^{(t)}/|\rho_{j, r', i'}^{(t)}|$ given in Lemma E.1 as follows:

$$\begin{aligned} &y_i f(\mathbf{W}^{(t)}, \mathbf{x}_i) \\ &= F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \mathbf{x}_i) - F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i) \\ &\geq \frac{1}{m} \sum_{r \in S_i^{(0)}} \langle \mathbf{w}_{y_i, r}^{(t)}, \mathbf{x}_i \rangle - F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i) \\ &= \frac{1}{m} \sum_{r \in S_i^{(0)}} \left[\langle \mathbf{w}_{y_i, r}^{(0)}, \mathbf{x}_i \rangle + \sum_{i'=1}^n \rho_{y_i, r, i'}^{(t)} \|\mathbf{x}_{i'}\|_2^{-2} \cdot \langle \mathbf{x}_{i'}, \mathbf{x}_i \rangle \right] - F_{-y_i}(\mathbf{W}_{-y_i}^{(t)}, \mathbf{x}_i) \\ &\geq \frac{1}{m} \sum_{r \in S_i^{(0)}} \rho_{y_i, r, i}^{(t)} - \frac{1}{m} \sum_{r \in S_i^{(0)}} \sum_{i' \neq i} |\rho_{y_i, r, i'}^{(t)}| R_{\min}^{-2} p - \frac{|S_i^{(0)}|}{m} \beta - \frac{1}{m} \sum_{r=1}^m \sum_{i' \neq i} |\rho_{-y_i, r, i'}^{(t)}| R_{\min}^{-2} p - \beta \end{aligned}$$

$$\begin{aligned}
&\geq \frac{1 - c_1^{-1} R_{\min}^{-2} p n - c_1^{-1} R_{\min}^{-2} p n m / |S_i^{(0)}|}{m} \sum_{r \in S_i^{(0)}} \rho_{y_i, r, i}^{(t)} - \left(\frac{|S_i^{(0)}|}{m} + 1 \right) \beta \\
&\geq \frac{1 - c_1^{-1} (1 + c') R_{\min}^{-2} p n}{m} \sum_{r \in S_i^{(0)}} \rho_{y_i, r, i}^{(t)} - 2\beta \\
&\geq \frac{c_2 (1 - c_1^{-1} (1 + c') R_{\min}^{-2} p n) |S_i^{(0)}|}{m} \log \left(1 + \frac{\eta |S_i^{(0)}| \|\mathbf{x}_i\|_2^2 e^{-\beta}}{2 n m^2} \cdot t \right) - 2\beta \\
&\geq \frac{1 - c_1^{-1} (1 + c') R_{\min}^{-2} p n}{c'^2 (1 + c_1^{-1} R_{\min}^{-2} p n)} \log \left(1 + \frac{\eta |S_i^{(0)}| \|\mathbf{x}_i\|_2^2 e^{-\beta}}{2 n m^2} \cdot t \right) - 2\beta, \tag{G.1}
\end{aligned}$$

838 where the first inequality is by the activation pattern $S_i^{(0)} \subseteq S_i^{(t)}$ given in Lemma E.1; the
 839 second inequality is by triangle inequality and (E.16); the third inequality is due to $\rho_{y_i, r, i}^{(t)} \geq$
 840 $c_1 |\rho_{j, r', i'}^{(t)}|$ from Lemma E.1 and hence $|S_i^{(0)}| \cdot |\rho_{y_i, r, i'}^{(t)}| \leq c_1^{-1} \sum_{r \in S_i^{(0)}} \rho_{y_i, r, i}^{(t)} |S_i^{(0)}| \cdot |\rho_{-y_i, r, i'}^{(t)}| \leq$
 841 $c_1^{-1} \sum_{r \in S_i^{(0)}} \rho_{y_i, r, i}^{(t)}$; the fourth inequality is by $|S_i^{(0)}| \leq m \leq c' |S_i^{(0)}|$; the second last inequality
 842 is by the second bullet of Lemma E.1; the last inequality is by $m \leq c' |S_i^{(0)}|$ and taking c_2 as
 843 $\frac{1}{c' (1 + R_{\min}^{-2} p n / c_1)}$ from the proof of Lemma E.1. Having obtained a lower bound for the margin, we
 844 can now use it to derive an upper bound for the loss function as follows:

$$\begin{aligned}
L_S(\mathbf{W}^{(t)}) &= \frac{1}{n} \sum_{i=1}^n \ell(y_i f(\mathbf{W}^{(t)}, \mathbf{x}_i)) \\
&= \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i f(\mathbf{W}^{(t)}, \mathbf{x}_i))) \\
&\leq \frac{1}{n} \sum_{i=1}^n \exp(-y_i f(\mathbf{W}^{(t)}, \mathbf{x}_i)) \\
&\leq \frac{1}{n} \sum_{i=1}^n \exp \left(- \frac{1 - c_1^{-1} (1 + c') R_{\min}^{-2} p n}{c'^2 (1 + R_{\min}^{-2} p n / c_1)} \log \left(1 + \frac{\eta |S_i^{(0)}| \|\mathbf{x}_i\|_2^2 e^{-\beta}}{2 n m^2} \cdot t \right) + 2\beta \right) \\
&= \frac{e^{2\beta}}{n} \sum_{i=1}^n \left(1 + \frac{\eta |S_i^{(0)}| \|\mathbf{x}_i\|_2^2 e^{-\beta}}{2 n m^2} \cdot t \right)^{-\frac{1 - c_1^{-1} (1 + c') R_{\min}^{-2} p n}{c'^2 (1 + c_1^{-1} R_{\min}^{-2} p n)}} \\
&= O \left(t^{-c'^{-2} (1 - c_1^{-1} (1 + c') R_{\min}^{-2} p n) / (1 + c_1^{-1} R_{\min}^{-2} p n)} \right),
\end{aligned}$$

845 where the first inequality is by $\log(1 + z) \leq z$; the second inequality is by (G.1). This completes the
 846 proof. \square

847 H Auxiliary Lemmas

848 **Lemma H.1.** Let $\{x_t\}_{t=0}^\infty$ be a non-negative sequence satisfying the following inequality:

$$x_{t+1} - x_t \leq C \cdot e^{-x_t}, \forall t \geq 0$$

849 then we have

$$x_t \leq \log(e^{x_0} + C e^C \cdot t).$$

850 *Proof of Lemma H.1.* Given the inequality $x_{t+1} - x_t \leq C \cdot e^{-x_t}$ for all $t \geq 0$, we want to prove that
 851 $x_T \leq \log(e^{x_0} + C e^C \cdot T)$ for $T \geq 0$. We start by manipulating the inequality as follows:

$$\begin{aligned}
&x_{t+1} - x_t \leq C \cdot e^{-x_t} \\
\Rightarrow &x_{t+1} - x_t \leq C \cdot e^{-x_{t+1} + C} \quad (\text{using } x_{t+1} \text{ instead of } x_t)
\end{aligned}$$

$$\implies e^{x_{t+1}}(x_{t+1} - x_t) \leq Ce^C \quad (\text{multiplying both sides by } e^{x_{t+1}}).$$

852 Summing the inequality from $t = 0$ to $t = T - 1$, we get:

$$\sum_{t=0}^{T-1} e^{x_{t+1}}(x_{t+1} - x_t) \leq Ce^C \cdot T.$$

853 Since e^x is a monotone increasing function, we can approximate the above sum with an integral:

$$\int_{x_0}^{x_T} e^x dx \leq Ce^C \cdot T.$$

854 Evaluating the integral, we get:

$$e^{x_T} - e^{x_0} \leq Ce^C \cdot T.$$

855 Rearranging the inequality, we get:

$$e^{x_T} \leq e^{x_0} + Ce^C \cdot T.$$

856 Taking the natural logarithm of both sides, we get:

$$x_T \leq \log(e^{x_0} + Ce^C \cdot T).$$

857 Therefore, we have shown that $x_T \leq \log(e^{x_0} + Ce^C \cdot T)$, as required. \square

858 **Lemma H.2.** Let $\{x_t\}_{t=0}^{\infty}$ be an sequence satisfying the following inequality:

$$x_{t+1} - x_t \geq C \cdot e^{-x_t}, \forall t \geq 0$$

859 then we have

$$x_t \geq \log(e^{x_0} + C \cdot t).$$

860 *Proof of Lemma H.2.* Given the inequality $x_{t+1} - x_t \geq C \cdot e^{-x_t}$ for all $t \geq 0$, we want to prove that
861 $x_T \geq \log(e^{x_0} + C \cdot T)$ for $T \geq 0$. We start by manipulating the inequality as follows:

$$\begin{aligned} x_{t+1} - x_t &\geq C \cdot e^{-x_t} \\ \implies e^{x_t}(x_{t+1} - x_t) &\geq C \quad (\text{multiplying both sides by } e^{x_t}). \end{aligned}$$

862 Summing the inequality from $t = 0$ to $t = T - 1$, we get:

$$\sum_{t=0}^{T-1} e^{x_t}(x_{t+1} - x_t) \geq C \cdot T.$$

863 Since e^x is a monotone increasing function, we can approximate the above sum with an integral:

$$\int_{x_0}^{x_T} e^x dx \geq C \cdot T.$$

864 Evaluating the integral, we get:

$$e^{x_T} - e^{x_0} \geq C \cdot T.$$

865 Rearranging the inequality, we get:

$$e^{x_T} \geq e^{x_0} + C \cdot T.$$

866 Taking the natural logarithm of both sides, we get:

$$x_T \geq \log(e^{x_0} + C \cdot T).$$

867 Therefore, we have shown that $x_T \geq \log(e^{x_0} + C \cdot T)$, as required. \square

868 **Lemma H.3** (Theorem 4.4.5 in Vershynin (2018)). Let \mathbf{A} be an $m \times n$ random matrix whose entries
869 a_{ij} are independent, mean zero, sub-gaussian random variables. Then for any $t > 0$ we have

$$\|\mathbf{A}\|_2 \leq CK(\sqrt{m} + \sqrt{n} + t)$$

870 with probability at least $1 - 2 \exp(-t^2)$. Here $K = \max_{i,j} \|a_{ij}\|_{\phi_2}$ where $\|\cdot\|_{\phi_2}$ is the sub-gaussian
871 norm.

872 **Lemma H.4.** For $t \geq s > 0$, we have

$$\frac{\log(1+at)}{\log(1+bt)} \geq \frac{\log(1+as)}{\log(1+bs)},$$

873 if $b > a > 0$.

874 *Proof of Lemma H.4.* Let $f(t) = \log(1+at)/\log(1+bt)$, and we want to prove that $f'(t) > 0$ for
875 all $t > 0$. To find the derivative of $f(t)$, we use the quotient rule:

$$\begin{aligned} f'(t) &= \frac{(\log(1+bt)) \frac{d}{dt}(\log(1+at)) - (\log(1+at)) \frac{d}{dt}(\log(1+bt))}{(\log(1+bt))^2} \\ &= \frac{(\log(1+bt)) \frac{a}{1+at} - (\log(1+at)) \frac{b}{1+bt}}{(\log(1+bt))^2} \\ &= \frac{a(1+bt) \log(1+bt) - b(1+at) \log(1+at)}{(1+at)(1+bt)(\log(1+bt))^2}. \end{aligned}$$

876 Next, we define the function $g(t) = (\frac{1}{b} + t) \log(1+bt) - (\frac{1}{a} + t) \log(1+at)$, and we aim to show
877 that $g'(t) > 0$ for all $t > 0$. We start by computing the derivative of $g(t)$:

$$g'(t) = \log(1+bt) - \log(1+at).$$

878 Since $b > a$ and $t > 0$, we have $1+bt > 1+at$, which implies that $\log(1+bt) > \log(1+at)$.
879 Therefore, we have $g'(t) > 0$ for all $t > 0$. Note that $g(0) = 0$, we then have $g(t) > 0$ for all $t > 0$.
880 Therefore, we have $a(1+bt) \log(1+bt) - b(1+at) \log(1+at) > 0$ for all $t > 0$, which in turn
881 implies that $f'(t) > 0$ for all $t > 0$. Thus, we have shown that $f(t)$ is increasing for $t > 0$ and hence
882 $f(t) > f(s)$, which completes the proof. \square

883 I Additional Experiments

884 In this section, we conduct additional experiments on the MINIST dataset. Our focus is the train-
885 ing of a two-layer feed-forward neural network, as discussed in Section 3, utilizing either ReLU
886 or leaky-ReLU activation functions. We examine different widths, specifically choosing from
887 $\{10, 50, 100, 500, 1000\}$.

888 The network initialization process follows a Gaussian distribution, with a variance of $\sigma_0 = 0.00001$.
889 Training is executed using stochastic gradient descent, a batch size of 64, and a learning rate of 0.1,
890 for a total of 10 epochs. As discerned from Figures 4 and 5, the stable rank of networks utilizing
891 either ReLU or leaky ReLU is weakly influenced by the width. For an exceedingly small width such
892 as 10, the weight matrix is low rank with a correspondingly small stable rank. However, this also
893 results in low test accuracy as the network cannot effectively learn all necessary features. As the
894 width increases, the test accuracy and final stable rank will increase. However, for sufficiently large
895 widths, an increase in width no longer corresponds to stable rank or test accuracy increases.

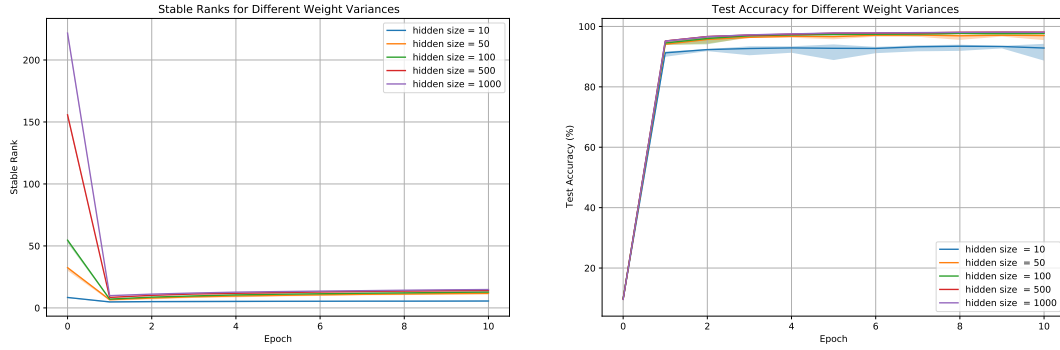


Figure 4: Stable ranks and test errors for different width across multiple runs (ReLU Activation Function). Each line represents the mean stable rank or test error for a given weight variance, while the shaded regions indicate the variability of the values (± 3 times the standard deviation) across the 5 runs.

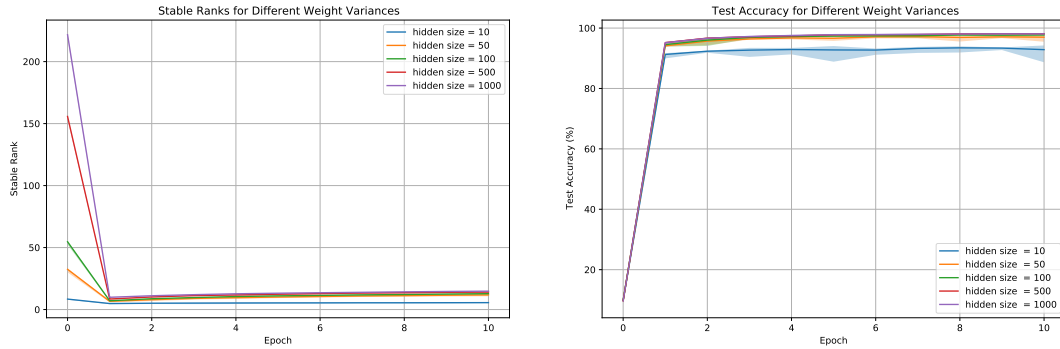


Figure 5: Stable ranks and test errors for different width across multiple runs (leaky-ReLU Activation Function). Each line represents the mean stable rank or test error for a given weight variance, while the shaded regions indicate the variability of the values (± 3 times the standard deviation) across the 5 runs.