

---

# Understanding Why Generalized Reweighting Does Not Improve Over ERM (Appendix)

---

Anonymous Author(s)

Affiliation

Address

email

## 1 A Related work

2 **Group Fairness.** Group fairness in machine learning was first studied in [HPS16] and [ZVGRG17],  
3 where they required the model to perform equally well over all groups. Later, [HSNL18] studied  
4 another type of group fairness called Rawlsian max-min fairness [Raw01], which does not require  
5 equal performance but rather requires high performance on the worst-off group. The subpopulation  
6 shift problem we study in this paper is most closely related to Rawlsian max-min fairness. A large  
7 body of recent work have studied how to improve this worst-group performance [DN18, OSHL19,  
8 LHC<sup>+</sup>21, ZDKR21]. Recent work however observe that these approaches, when used with modern  
9 overparameterized models, easily overfit [SKHL20, SRKL20]. Apart from group fairness, there are  
10 also other notions of fairness, such as individual fairness [DHP<sup>+</sup>12, ZWS<sup>+</sup>13] and counterfactual  
11 fairness [KLRS17], which we do not study in this work.

12 **Implicit Bias Under the Overparameterized Setting.** For overparameterized models, there could  
13 be many model parameters which all minimize the training loss. In such cases, it is of interest to study  
14 the implicit bias of specific optimization algorithms such as gradient descent i.e. to what minimizer  
15 the model parameters will converge to [DLL<sup>+</sup>19, AZLS19]. Our results use the NTK formulation of  
16 wide neural networks [JGH18], and specifically we use linearized neural networks to approximate  
17 such wide neural networks following [LXS<sup>+</sup>19]. There is some criticism of this line of work, e.g.  
18 [COB19] argued that infinitely wide neural networks fall in the “lazy training” regime and results  
19 might not be transferable to general neural networks. Nonetheless such wide neural networks are  
20 being widely studied in recent years, since they provide considerable insights into the behavior of  
21 more general neural networks, which are typically intractable to analyze otherwise.

## 22 B Extension to Multi-dimensional Regression / Multi-class Classification

23 In our results, we assume that  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  for simplicity, but our results can be very easily extended  
24 to the case where  $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$ . For most of our results, the proof consists of two major components:  
25 (i) The linearized neural network will converge to some point (interpolator, max-margin classifier,  
26 etc.); (ii) The wide fully-connected neural network can be approximated by its linearized counterpart.  
27 For both components, the extension is very simple and straightforward. For (i), the proof only  
28 relies on the smoothness of the objective function and the upper quadratic bound it entails, and  
29 the function is still smooth when its output becomes multi-dimensional; For (ii), we can prove that  
30  $\sup_t \|f(\mathbf{x}) - f_{\text{lin}}(\mathbf{x})\|_2 = O(\tilde{d}^{-1/4})$  in exactly the same way. Thus, all of our results hold for  
31 multi-dimensional regression and multi-class classification.

32 Particularly, for the multi-class cross-entropy loss, using Theorem 8 we can show that under any  
33 GRW satisfying Assumption 1, the direction of the weight of a linear classifier will converge to the

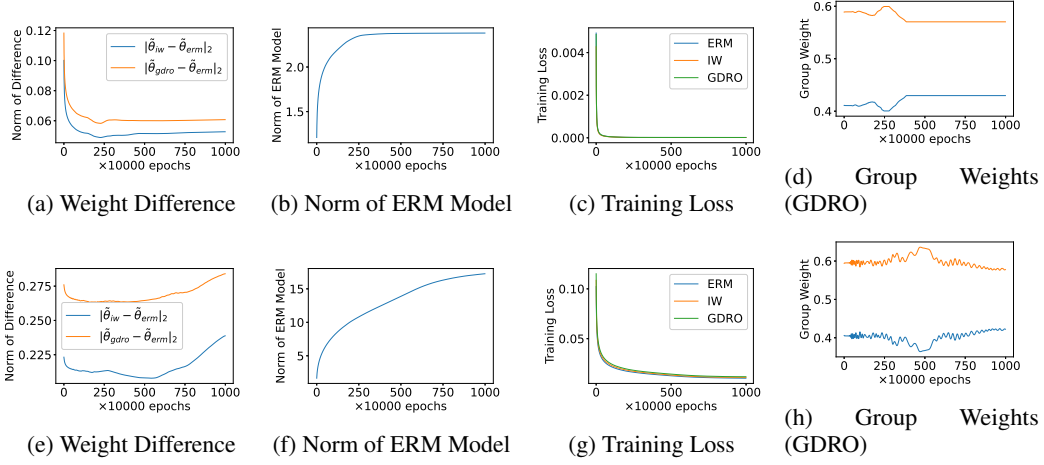


Figure 1: Experimental results of ERM, importance weighting (IW) and Group DRO (GDRO) with the logistic loss and the polynomially-tailed loss. First row: Logistic loss; Second row: Polynomially-tailed loss. All norms are  $L_2$  norms.  $\tilde{\theta}$  is a unit vector which is the direction of  $\theta$ .

34 following max-margin classifier:

$$\hat{\theta}_{\text{MM}} = \arg \min_{\theta} \left\{ \min_{i=1, \dots, n} \left[ f(\mathbf{x}_i)_{y_i} - \max_{y' \neq y_i} f(\mathbf{x}_i)_{y'} \right] : \|\theta\|_2 = 1 \right\} \quad (20)$$

35 which is still independent of  $q_i$ .

## 36 C More Experiments

37 We run ERM, importance weighting and Group DRO on the training set with 6 MNIST images which  
 38 we used in Section 4.1 with the logistic loss and the polynomially-tailed loss (Eqn. (19), with  $\alpha = 1$ ,  
 39  $\beta = 0$  and  $\ell_{\text{left}}$  being the logistic loss shifted to make the overall loss function continuous) on this  
 40 dataset for 10 million epochs (note that we run for much more epochs because the convergence is  
 41 very slow). The results are shown in Figure 1. From the plots we can see that:

- 42 • For either loss function, the training loss of each method converges to 0.
- 43 • In contrast to the theory that the norm of the ERM model will go to infinity and all models  
 44 will converge to the max-margin classifier, the weight of the ERM model gets stuck at some  
 45 point, and the norms of the gaps between the normalized model weights also get stuck.  
 46 The reason is that the training loss has got so small that it becomes zero in the floating  
 47 number representation, so the gradient also becomes zero and the training halts due to  
 48 limited computational precision.
- 49 • However, we can still observe a fundamental difference between the logistic loss and the  
 50 polynomially-tailed loss. For the logistic loss, the norm of the gap between importance  
 51 weighting (or Group DRO) and ERM will converge to around 0.06 when the training stops,  
 52 while for the polynomially-tailed loss, the norm will be larger than 0.22 and will keep  
 53 growing, which shows that for the polynomially-tailed loss the normalized model weights  
 54 do not converge to the same point.
- 55 • For either loss, the group weights of Group DRO still empirically satisfy Assumption 1.

## 56 D Proofs

57 In this paper, for any matrix  $\mathbf{A}$ , we will use  $\|\mathbf{A}\|_2$  to denote its spectral norm and  $\|\mathbf{A}\|_F$  to denote its  
 58 Frobenius norm.

## 59 D.1 Background on Smoothness

60 A first-order differentiable function  $f$  over  $\mathcal{D}$  is called  $L$ -smooth for  $L > 0$  if

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{D} \quad (21)$$

61 which is also called the *upper quadratic bound*. If  $f$  is second-order differentiable and  $\mathcal{D}$  is a convex  
62 set, then  $f$  is  $L$ -smooth is equivalent to

$$\mathbf{v}^\top \nabla^2 f(\mathbf{x}) \mathbf{v} \leq L \quad \forall \|\mathbf{v}\|_2 = 1, \forall \mathbf{x} \in \mathcal{D} \quad (22)$$

63 A classical result in convex optimization is the following:

64 **Theorem 10.** *If  $f(\mathbf{x})$  is convex and  $L$ -smooth with a unique finite minimizer  $\mathbf{x}^*$ , and is minimized  
65 by gradient descent  $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)$  starting from  $\mathbf{x}_0$  where the learning rate  $\eta \leq \frac{1}{L}$ , then we  
66 have*

$$f(\mathbf{x}_T) \leq f(\mathbf{x}^*) + \frac{1}{\eta T} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \quad (23)$$

67 which also implies that  $\mathbf{x}_T$  converges to  $\mathbf{x}^*$  as  $T \rightarrow \infty$ .

## 68 D.2 Proofs for Subsection 4.1

### 69 D.2.1 Proof of Theorem 1

70 Using the key intuition, the weight update rule (7) implies that  $\theta^{(t+1)} - \theta^{(t)} \in \text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  for  
71 all  $t$ , which further implies that  $\theta^{(t)} - \theta^{(0)} \in \text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  for all  $t$ . By Cramer's rule, in this  
72  $n$ -dimensional subspace there exists one and only one  $\theta^*$  such that  $\theta^* - \theta^{(0)} \in \text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$   
73 and  $\langle \theta^*, \mathbf{x}_i \rangle$  for all  $i$ . Then we have

$$\left\| \mathbf{X}^\top (\theta^{(t)} - \theta^*) \right\|_2 = \left\| (\mathbf{X}^\top \theta^{(t)} - \mathbf{Y}) - (\mathbf{X}^\top \theta^* - \mathbf{Y}) \right\|_2 \leq \left\| \mathbf{X}^\top \theta^{(t)} - \mathbf{Y} \right\|_2 + \left\| \mathbf{X}^\top \theta^* - \mathbf{Y} \right\|_2 \rightarrow 0 \quad (24)$$

74 because  $\left\| \mathbf{X}^\top \theta - \mathbf{Y} \right\|_2^2 = 2n\hat{\mathcal{R}}(f(\mathbf{x}; \theta))$ . On the other hand, let  $s^{\min}$  be the smallest singular value  
75 of  $\mathbf{X}$ . Since  $\mathbf{X}$  is full-rank,  $s^{\min} > 0$ , and  $\left\| \mathbf{X}^\top (\theta^{(t)} - \theta^*) \right\|_2 \geq s^{\min} \|\theta^{(t)} - \theta^*\|_2$ . This shows  
76 that  $\|\theta^{(t)} - \theta^*\|_2 \rightarrow 0$ . Thus,  $\theta^{(t)}$  converges to this unique  $\theta^*$ .  $\square$

### 77 D.2.2 Proof of Theorem 2

78 To help our readers understand the proof more easily, we will first prove the result for static GRW  
79 where  $q_i^{(t)} = q_i$  for all  $t$ , and then we will prove the result for dynamic GRW that satisfy  $q_i^{(t)} \rightarrow q_i$  as  
80  $t \rightarrow \infty$ .

81 **Static GRW.** We first prove the result for all static GRW such that  $\min_i q_i = q^* > 0$ .

82 We will use smoothness introduce in Appendix D.1. Denote  $A = \sum_{i=1}^n \|\mathbf{x}_i\|_2^2$ . The empirical risk of  
83 the linear model  $f(\mathbf{x}) = \langle \theta, \mathbf{x} \rangle$  is

$$F(\theta) = \sum_{i=1}^n q_i (\mathbf{x}_i^\top \theta - y_i)^2 \quad (25)$$

84 whose Hessian is

$$\nabla_\theta^2 F(\theta) = 2 \sum_{i=1}^n q_i \mathbf{x}_i \mathbf{x}_i^\top \quad (26)$$

85 So for any unit vector  $\mathbf{v} \in \mathbb{R}^d$ , we have (since  $q_i \in [0, 1]$ )

$$\mathbf{v}^\top \nabla_\theta^2 F(\theta) \mathbf{v} = 2 \sum_{i=1}^n q_i (\mathbf{x}_i^\top \mathbf{v})^2 \leq 2 \sum_{i=1}^n q_i \|\mathbf{x}_i\|_2^2 \leq 2A \quad (27)$$

86 which implies that  $F(\theta)$  is  $2A$ -smooth. Thus, we have the following upper quadratic bound: for any  
 87  $\theta_1, \theta_2 \in \mathbb{R}^d$ ,

$$F(\theta_2) \leq F(\theta_1) + \langle \nabla_\theta F(\theta_1), \theta_2 - \theta_1 \rangle + A \|\theta_2 - \theta_1\|_2^2 \quad (28)$$

88 Denote  $g(\theta^{(t)}) = (\mathbf{X}^\top \theta^{(t)} - \mathbf{Y}) \in \mathbb{R}^n$ . We can see that  $\|\sqrt{\mathbf{Q}}g(\theta^{(t)})\|_2^2 = F(\theta^{(t)})$ , where  
 89  $\sqrt{\mathbf{Q}} = \text{diag}(\sqrt{q_1}, \dots, \sqrt{q_n})$ . Thus,  $\nabla F(\theta^{(t)}) = 2\mathbf{X}\mathbf{Q}g(\theta^{(t)})$ . The update rule of a static GRW  
 90 with gradient descent and the squared loss is:

$$\theta^{(t+1)} = \theta^{(t)} - \eta \sum_{i=1}^n q_i \mathbf{x}_i (f^{(t)}(\mathbf{x}_i) - y_i) = \theta^{(t)} - \eta \mathbf{X}\mathbf{Q}g(\theta^{(t)}) \quad (29)$$

91 Substituting  $\theta_1$  and  $\theta_2$  in (28) with  $\theta^{(t)}$  and  $\theta^{(t+1)}$  yields

$$F(\theta^{(t+1)}) \leq F(\theta^{(t)}) - 2\eta g(\theta^{(t)})^\top \mathbf{Q}^\top \mathbf{X}^\top \mathbf{X}\mathbf{Q}g(\theta^{(t)}) + A \left\| \eta \mathbf{X}\mathbf{Q}g(\theta^{(t)}) \right\|_2^2 \quad (30)$$

92 Since  $\mathbf{x}_1, \dots, \mathbf{x}_n$  are linearly independent,  $\mathbf{X}^\top \mathbf{X}$  is a positive definite matrix. Denote the smallest  
 93 eigenvalue of  $\mathbf{X}^\top \mathbf{X}$  by  $\lambda^{\min} > 0$ . And  $\|\mathbf{Q}g(\theta^{(t)})\|_2 \geq \sqrt{q^*} \|g(\theta^{(t)})\|_2 = \sqrt{q^* F(\theta^{(t)})}$ , so we have  
 94  $g(\theta^{(t)})^\top \mathbf{Q}^\top \mathbf{X}^\top \mathbf{X}\mathbf{Q}g(\theta^{(t)}) \geq q^* \lambda^{\min} F(\theta^{(t)})$ . Thus,

$$\begin{aligned} F(\theta^{(t+1)}) &\leq F(\theta^{(t)}) - 2\eta q^* \lambda^{\min} F(\theta^{(t)}) + A\eta^2 \left\| \mathbf{X}\sqrt{\mathbf{Q}} \right\|_2^2 \left\| \sqrt{\mathbf{Q}}g(\theta^{(t)}) \right\|_2^2 \\ &\leq F(\theta^{(t)}) - 2\eta q^* \lambda^{\min} F(\theta^{(t)}) + A\eta^2 \left\| \mathbf{X}\sqrt{\mathbf{Q}} \right\|_F^2 F(\theta^{(t)}) \\ &\leq F(\theta^{(t)}) - 2\eta q^* \lambda^{\min} F(\theta^{(t)}) + A\eta^2 \|\mathbf{X}\|_F^2 F(\theta^{(t)}) \\ &= (1 - 2\eta q^* \lambda^{\min} + A^2 \eta^2) F(\theta^{(t)}) \end{aligned} \quad (31)$$

95 Let  $\eta_0 = \frac{q^* \lambda^{\min}}{A^2}$ . For any  $\eta \leq \eta_0$ , we have  $F(\theta^{(t+1)}) \leq (1 - \eta q^* \lambda^{\min}) F(\theta^{(t)})$  for all  $t$ , which  
 96 implies that  $\lim_{t \rightarrow \infty} F(\theta^{(t)}) = 0$ . This implies that the empirical training risk must converge to 0.

97 **Dynamic GRW.** Now we prove the result for all dynamic GRW satisfying Assumption 1. By  
 98 Assumption 1, for any  $\epsilon > 0$ , there exists  $t_\epsilon$  such that for all  $t \geq t_\epsilon$  and all  $i$ ,

$$q_i^{(t)} \in (q_i - \epsilon, q_i + \epsilon) \quad (32)$$

99 This is because for all  $i$ , there exists  $t_i$  such that for all  $t \geq t_i$ ,  $q_i^{(t)} \in (q_i - \epsilon, q_i + \epsilon)$ . Then, we can  
 100 define  $t_\epsilon = \max\{t_1, \dots, t_n\}$ . Denote the largest and smallest eigenvalues of  $\mathbf{X}^\top \mathbf{X}$  by  $\lambda^{\max}$  and  
 101  $\lambda^{\min}$ , and because  $\mathbf{X}$  is full-rank, we have  $\lambda^{\min} > 0$ . Define  $\epsilon = \min\{\frac{q^*}{3}, \frac{(q^* \lambda^{\min})^2}{12\lambda^{\max 2}}\}$ , and then  $t_\epsilon$  is  
 102 also fixed.

103 We still denote  $\mathbf{Q} = \text{diag}(q_1, \dots, q_n)$ . When  $t \geq t_\epsilon$ , the update rule of a dynamic GRW with  
 104 gradient descent and the squared loss is:

$$\theta^{(t+1)} = \theta^{(t)} - \eta \mathbf{X}\mathbf{Q}_\epsilon^{(t)} (\mathbf{X}^\top \theta^{(t)} - \mathbf{Y}) \quad (33)$$

105 where  $\mathbf{Q}_\epsilon^{(t)} = \mathbf{Q}^{(t)}$ , and we use the subscript  $\epsilon$  to indicate that  $\|\mathbf{Q}_\epsilon^{(t)} - \mathbf{Q}\|_2 < \epsilon$ . Then, note that we  
 106 can rewrite  $\mathbf{Q}_\epsilon^{(t)}$  as  $\mathbf{Q}_\epsilon^{(t)} = \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} \cdot \sqrt{\mathbf{Q}}$  as long as  $\epsilon \leq q^*/3$ . This is because  $q_i + \epsilon \leq \sqrt{(q_i + 3\epsilon)q_i}$   
 107 and  $q_i - \epsilon \geq \sqrt{(q_i - 3\epsilon)q_i}$  for all  $\epsilon \leq q_i/3$ , and  $q_i \geq q^*$ . Thus, we have

$$\theta^{(t+1)} = \theta^{(t)} - \eta \mathbf{X} \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} \sqrt{\mathbf{Q}} g(\theta^{(t)}) \quad \text{where } \mathbf{Q}_\epsilon^{(t)} = \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} \cdot \sqrt{\mathbf{Q}} \quad (34)$$

108 Again, substituting  $\theta_1$  and  $\theta_2$  in (28) with  $\theta^{(t)}$  and  $\theta^{(t+1)}$  yields

$$F(\theta^{(t+1)}) \leq F(\theta^{(t)}) - 2\eta g(\theta^{(t)})^\top \mathbf{Q}^\top \mathbf{X}^\top \mathbf{X} \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} \sqrt{\mathbf{Q}} g(\theta^{(t)}) + A \left\| \eta \mathbf{X} \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} \sqrt{\mathbf{Q}} g(\theta^{(t)}) \right\|_2^2 \quad (35)$$

109 Then, note that

$$\begin{aligned}
& \left| g(\theta^{(t)})^\top \mathbf{Q}^\top \mathbf{X}^\top \mathbf{X} \left( \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} - \sqrt{\mathbf{Q}} \right) \sqrt{\mathbf{Q}} g(\theta^{(t)}) \right| \\
& \leq \left\| \sqrt{\mathbf{Q}}^\top \mathbf{X}^\top \mathbf{X} \left( \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} - \sqrt{\mathbf{Q}} \right) \right\|_2 \left\| \sqrt{\mathbf{Q}} g(\theta^{(t)}) \right\|_2^2 \\
& \leq \left\| \sqrt{\mathbf{Q}} \right\|_2 \left\| \mathbf{X}^\top \mathbf{X} \right\|_2 \left\| \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} - \sqrt{\mathbf{Q}} \right\|_2 \left\| \sqrt{\mathbf{Q}} g(\theta^{(t)}) \right\|_2^2 \\
& \leq \lambda^{\max} \sqrt{3\epsilon} F(\theta^{(t)})
\end{aligned} \tag{36}$$

110 where the last step comes from the following fact: for all  $\epsilon < q_i/3$ ,

$$\sqrt{q_i + 3\epsilon} - \sqrt{q_i} \leq \sqrt{3\epsilon} \quad \text{and} \quad \sqrt{q_i} - \sqrt{q_i - 3\epsilon} \leq \sqrt{3\epsilon} \tag{37}$$

111 And as proved before, we also have

$$g(\theta^{(t)})^\top \mathbf{Q}^\top \mathbf{X}^\top \mathbf{X} \mathbf{Q} g(\theta^{(t)}) \geq q^* \lambda^{\min} F(\theta^{(t)}) \tag{38}$$

112 Since  $\epsilon \leq \frac{(q^* \lambda^{\min})^2}{12\lambda^{\max} 2}$ , we have

$$g(\theta^{(t)})^\top \mathbf{Q}^\top \mathbf{X}^\top \mathbf{X} \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} \sqrt{\mathbf{Q}} g(\theta^{(t)}) \geq (q^* \lambda^{\min} - \lambda^{\max} \sqrt{3\epsilon}) F(\theta^{(t)}) \geq \frac{1}{2} q^* \lambda^{\min} F(\theta^{(t)}) \tag{39}$$

113 Thus,

$$\begin{aligned}
F(\theta^{(t+1)}) & \leq F(\theta^{(t)}) - \eta q^* \lambda^{\min} F(\theta^{(t)}) + A\eta^2 \left\| \mathbf{X} \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} \right\|_2^2 \left\| \sqrt{\mathbf{Q}} g(\theta^{(t)}) \right\|_2^2 \\
& \leq (1 - \eta q^* \lambda^{\min} + A^2 \eta^2 (1 + 3\epsilon)) F(\theta^{(t)}) \\
& \leq (1 - \eta q^* \lambda^{\min} + 2A^2 \eta^2) F(\theta^{(t)})
\end{aligned} \tag{40}$$

114 for all  $\epsilon \leq 1/3$ . Let  $\eta_0 = \frac{q^* \lambda^{\min}}{4A^2}$ . For any  $\eta \leq \eta_0$ , we have  $F(\theta^{(t+1)}) \leq (1 - \eta q^* \lambda^{\min}/2) F(\theta^{(t)})$   
115 for all  $t \geq t_\epsilon$ , which implies that  $\lim_{t \rightarrow \infty} F(\theta^{(t)}) = 0$ . Thus, the empirical training risk converges to  
116 0.  $\square$

### 117 D.3 Proofs for Subsection 4.2

#### 118 D.3.1 Proof of Lemma 3

119 Note that the first  $l$  layers (except the output layer) of the original NTK formulation and our new  
120 formulation are the same, so we still have the following proposition:

121 **Proposition 11** (Proposition 1 in [JGH18]). *If  $\sigma$  is Lipschitz and  $d_l \rightarrow \infty$  for  $l = 1, \dots, L$   
122 sequentially, then for all  $l = 1, \dots, L$ , the distribution of a single element of  $\mathbf{h}^l$  converges in  
123 probability to a zero-mean Gaussian process of covariance  $\Sigma^l$  that is defined recursively by:*

$$\begin{aligned}
\Sigma^1(\mathbf{x}, \mathbf{x}') &= \frac{1}{d_0} \mathbf{x}^\top \mathbf{x}' + \beta^2 \\
\Sigma^l(\mathbf{x}, \mathbf{x}') &= \mathbb{E}_f[\sigma(f(\mathbf{x}))\sigma(f(\mathbf{x}'))] + \beta^2
\end{aligned} \tag{41}$$

124 where  $f$  is sampled from a zero-mean Gaussian process of covariance  $\Sigma^{(l-1)}$ .

125 Now we show that for an infinitely wide neural network with  $L \geq 1$  hidden layers,  $\Theta^{(0)}$  converges in  
126 probability to the following non-degenerated deterministic limiting kernel

$$\Theta = \mathbb{E}_{f \sim \Sigma^L}[\sigma(f(\mathbf{x}))\sigma(f(\mathbf{x}'))] + \beta^2 \tag{42}$$

127 Consider the output layer  $\mathbf{h}^{L+1} = \frac{W^L}{\sqrt{d}} \sigma(\mathbf{h}^L) + \beta \mathbf{b}^L$ . We can see that for any parameter  $\theta_i$  before  
128 the output layer,

$$\nabla_{\theta_i} \mathbf{h}^{L+1} = \text{diag}(\dot{\sigma}(\mathbf{h}^L)) \frac{W^{L\top}}{\sqrt{d_L}} \nabla_{\theta_i} \mathbf{h}^L = 0 \tag{43}$$

129 And for  $W^L$  and  $\mathbf{b}^L$ , we have

$$\nabla_{W^L} \mathbf{h}^{L+1} = \frac{1}{\sqrt{d_L}} \sigma(\mathbf{h}^L) \quad \text{and} \quad \nabla_{\mathbf{b}^L} \mathbf{h}^{L+1} = \beta \quad (44)$$

130 Then we can achieve (42) by the law of large numbers.  $\square$

### 131 D.3.2 Proof of Lemma 5

132 We will use the following short-hand in the proof:

$$\begin{cases} g(\theta^{(t)}) = f^{(t)}(\mathbf{X}) - \mathbf{Y} \\ J(\theta^{(t)}) = \nabla_{\theta} f(\mathbf{X}; \theta^{(t)}) \in \mathbb{R}^{p \times n} \\ \Theta^{(t)} = J(\theta^{(t)})^{\top} J(\theta^{(t)}) \end{cases} \quad (45)$$

133 For any  $\epsilon > 0$ , there exists  $t_{\epsilon}$  such that for all  $t \geq t_{\epsilon}$  and all  $i, q_i^{(t)} \in (q_i - \epsilon, q_i + \epsilon)$ . Like what we  
134 have done in (34), we can rewrite  $\mathbf{Q}^{(t)} = \mathbf{Q}_{3\epsilon}^{(t)} = \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} \cdot \sqrt{\mathbf{Q}}$ , where  $\mathbf{Q} = \text{diag}(q_1, \dots, q_n)$ .

135 The update rule of a GRW with gradient descent and the squared loss for the wide neural network is:

$$\theta^{(t+1)} = \theta^{(t)} - \eta J(\theta^{(t)}) \mathbf{Q}^{(t)} g(\theta^{(t)}) \quad (46)$$

136 and for  $t \geq t_{\epsilon}$ , it can be rewritten as

$$\theta^{(t+1)} = \theta^{(t)} - \eta J(\theta^{(t)}) \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} \left[ \sqrt{\mathbf{Q}} g(\theta^{(t)}) \right] \quad (47)$$

137 First, we will prove the following theorem:

138 **Theorem 12.** *There exist constants  $M > 0$  and  $\epsilon_0 > 0$  such that for all  $\epsilon \in (0, \epsilon_0]$ ,  $\eta \leq \eta^*$  and any  
139  $\delta > 0$ , there exist  $R_0 > 0$ ,  $\tilde{D} > 0$  and  $B > 1$  such that for any  $\tilde{d} \geq \tilde{D}$ , the following (i) and (ii)  
140 hold with probability at least  $(1 - \delta)$  over random initialization when applying gradient descent with  
141 learning rate  $\eta$ :*

142 (i) For all  $t \leq t_{\epsilon}$ , there is

$$\left\| g(\theta^{(t)}) \right\|_2 \leq B^t R_0 \quad (48)$$

$$\sum_{j=1}^t \left\| \theta^{(j)} - \theta^{(j-1)} \right\|_2 \leq \eta M R_0 \sum_{j=1}^t B^{j-1} < \frac{M B^{t_{\epsilon}} R_0}{B - 1} \quad (49)$$

143 (ii) For all  $t \geq t_{\epsilon}$ , we have

$$\left\| \sqrt{\mathbf{Q}} g(\theta^{(t)}) \right\|_2 \leq \left( 1 - \frac{\eta q^* \lambda^{\min}}{3} \right)^{t-t_{\epsilon}} B^{t_{\epsilon}} R_0 \quad (50)$$

$$\begin{aligned} \sum_{j=t_{\epsilon}+1}^t \left\| \theta^{(j)} - \theta^{(j-1)} \right\|_2 &\leq \eta \sqrt{1 + 3\epsilon} M B^{t_{\epsilon}} R_0 \sum_{j=t_{\epsilon}+1}^t \left( 1 - \frac{\eta q^* \lambda^{\min}}{3} \right)^{j-t_{\epsilon}} \\ &< \frac{3\sqrt{1 + 3\epsilon} M B^{t_{\epsilon}} R_0}{q^* \lambda^{\min}} \end{aligned} \quad (51)$$

144 *Proof.* The proof is based on the following lemma:

145 **Lemma 13** (Local Lipschitzness of the Jacobian). *Under Assumption 2, there is a constant  $M > 0$   
146 such that for any  $C_0 > 0$  and any  $\delta > 0$ , there exists a  $\tilde{D}$  such that: If  $\tilde{d} \geq \tilde{D}$ , then with probability  
147 at least  $(1 - \delta)$  over random initialization, for any  $\mathbf{x}$  such that  $\|\mathbf{x}\|_2 \leq 1$ ,*

$$\left\{ \begin{aligned} \left\| \nabla_{\theta} f(\mathbf{x}; \theta) - \nabla_{\theta} f(\mathbf{x}; \tilde{\theta}) \right\|_2 &\leq \frac{M}{\sqrt[4]{\tilde{d}}} \left\| \theta - \tilde{\theta} \right\|_2 \\ \left\| \nabla_{\theta} f(\mathbf{x}; \theta) \right\|_2 &\leq M \\ \left\| J(\theta) - J(\tilde{\theta}) \right\|_F &\leq \frac{M}{\sqrt[4]{\tilde{d}}} \left\| \theta - \tilde{\theta} \right\|_2 \\ \left\| J(\theta) \right\|_F &\leq M \end{aligned} \right., \quad \forall \theta, \tilde{\theta} \in B(\theta^{(0)}, C_0) \quad (52)$$

148 where  $B(\theta^{(0)}, R) = \{\theta : \|\theta - \theta^{(0)}\|_2 < R\}$ .

149 The proof of this lemma can be found in Appendix D.3.3. Note that for any  $\mathbf{x}$ ,  $f^{(0)}(\mathbf{x}) = \beta \mathbf{b}^L$  where  
 150  $\mathbf{b}^L$  is sampled from the standard Gaussian distribution. Thus, for any  $\delta > 0$ , there exists a constant  
 151  $R_0$  such that with probability at least  $(1 - \delta/3)$  over random initialization,

$$\|g(\theta^{(0)})\|_2 < R_0 \quad (53)$$

152 And by Proposition 3, there exists  $D_2 \geq 0$  such that for any  $\tilde{d} \geq D_2$ , with probability at least  
 153  $(1 - \delta/3)$ ,

$$\|\Theta - \Theta^{(0)}\|_F \leq \frac{q^* \lambda^{\min}}{3} \quad (54)$$

154 Let  $M$  be the constant in Lemma 13. Let  $\epsilon_0 = \frac{(q^* \lambda^{\min})^2}{108M^4}$ . Let  $B = 1 + \eta^* M^2$ , and  $C_0 =$   
 155  $\frac{MB^{t_\epsilon} R_0}{B-1} + \frac{3\sqrt{1+3\epsilon} MB^{t_\epsilon} R_0}{q^* \lambda^{\min}}$ . By Lemma 13, there exists  $D_1 > 0$  such that with probability at least  
 156  $(1 - \delta/3)$ , for any  $\tilde{d} \geq D_1$ , (52) is true for all  $\theta, \tilde{\theta} \in B(\theta^{(0)}, C_0)$ .

157 By union bound, with probability at least  $(1 - \delta)$ , (52), (53) and (54) are all true. Now we assume  
 158 that all of them are true, and prove (48) and (49) by induction. (48) is true for  $t = 0$  due to (53), and  
 159 (49) is always true for  $t = 0$ . Suppose (48) and (49) are true for  $t$ , then for  $t + 1$  we have

$$\begin{aligned} \|\theta^{(t+1)} - \theta^{(t)}\|_2 &\leq \eta \|J(\theta^{(t)}) \mathbf{Q}^{(t)}\|_2 \|g(\theta^{(t)})\|_2 \leq \eta \|J(\theta^{(t)}) \mathbf{Q}^{(t)}\|_F \|g(\theta^{(t)})\|_2 \\ &\leq \eta \|J(\theta^{(t)})\|_F \|g(\theta^{(t)})\|_2 \leq M \eta B^t R_0 \end{aligned} \quad (55)$$

160 So (49) is also true for  $t + 1$ . And we also have

$$\begin{aligned} \|g(\theta^{(t+1)})\|_2 &= \|g(\theta^{(t+1)}) - g(\theta^{(t)}) + g(\theta^{(t)})\|_2 \\ &= \|J(\tilde{\theta}^{(t)})^\top (\theta^{(t+1)} - \theta^{(t)}) + g(\theta^{(t)})\|_2 \\ &= \|- \eta J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)}) \mathbf{Q}^{(t)} g(\theta^{(t)}) + g(\theta^{(t)})\|_2 \\ &\leq \|\mathbf{I} - \eta J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)}) \mathbf{Q}^{(t)}\|_2 \|g(\theta^{(t)})\|_2 \\ &\leq \left(1 + \|\eta J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)}) \mathbf{Q}^{(t)}\|_2\right) \|g(\theta^{(t)})\|_2 \\ &\leq \left(1 + \eta \|J(\tilde{\theta}^{(t)})\|_F \|J(\theta^{(t)})\|_F\right) \|g(\theta^{(t)})\|_2 \\ &\leq (1 + \eta^* M^2) \|g(\theta^{(t)})\|_2 \leq B^{t+1} R_0 \end{aligned} \quad (56)$$

161 Therefore, (48) and (49) are true for all  $t \leq t_\epsilon$ , which implies that  $\|\sqrt{\mathbf{Q}}g(\theta^{(t_\epsilon)})\|_2 \leq \|g(\theta^{(t_\epsilon)})\|_2 \leq$   
 162  $B^{t_\epsilon} R_0$ , so (50) is true for  $t = t_\epsilon$ . And (51) is obviously true for  $t = t_\epsilon$ . Now, let us prove (ii) by  
 163 induction. Note that when  $t \geq t_\epsilon$ , we have the alternative update rule (47). If (50) and (51) are true  
 164 for  $t$ , then for  $t + 1$ , there is

$$\begin{aligned} \|\theta^{(t+1)} - \theta^{(t)}\|_2 &\leq \eta \|J(\theta^{(t)}) \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}}\|_2 \|\sqrt{\mathbf{Q}}g(\theta^{(t)})\|_2 \leq \eta \|J(\theta^{(t)}) \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}}\|_F \|\sqrt{\mathbf{Q}}g(\theta^{(t)})\|_2 \\ &\leq \eta \sqrt{1+3\epsilon} \|J(\theta^{(t)})\|_F \|\sqrt{\mathbf{Q}}g(\theta^{(t)})\|_2 \leq M \eta \sqrt{1+3\epsilon} \left(1 - \frac{\eta q^* \lambda^{\min}}{3}\right)^{t-t_\epsilon} B^{t_\epsilon} R_0 \end{aligned} \quad (57)$$

165 So (51) is true for  $t + 1$ . And we also have

$$\begin{aligned}
\left\| \sqrt{\mathbf{Q}}g(\theta^{(t+1)}) \right\|_2 &= \left\| \sqrt{\mathbf{Q}}g(\theta^{(t+1)}) - \sqrt{\mathbf{Q}}g(\theta^{(t)}) + \sqrt{\mathbf{Q}}g(\theta^{(t)}) \right\|_2 \\
&= \left\| \sqrt{\mathbf{Q}}J(\tilde{\theta}^{(t)})^\top (\theta^{(t+1)} - \theta^{(t)}) + \sqrt{\mathbf{Q}}g(\theta^{(t)}) \right\|_2 \\
&= \left\| -\eta\sqrt{\mathbf{Q}}J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)})\mathbf{Q}^{(t)}g(\theta^{(t)}) + \sqrt{\mathbf{Q}}g(\theta^{(t)}) \right\|_2 \\
&\leq \left\| \mathbf{I} - \eta\sqrt{\mathbf{Q}}J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)})\sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} \right\|_2 \left\| \sqrt{\mathbf{Q}}g(\theta^{(t)}) \right\|_2 \\
&\leq \left\| \mathbf{I} - \eta\sqrt{\mathbf{Q}}J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)})\sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} \right\|_2 \left( 1 - \frac{\eta q^* \lambda^{\min}}{3} \right)^t R_0
\end{aligned} \tag{58}$$

166 where  $\tilde{\theta}^{(t)}$  is some linear interpolation between  $\theta^{(t)}$  and  $\theta^{(t+1)}$ . Now we prove that

$$\left\| \mathbf{I} - \eta\sqrt{\mathbf{Q}}J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)})\sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} \right\|_2 \leq 1 - \frac{\eta q^* \lambda^{\min}}{3} \tag{59}$$

167 For any unit vector  $\mathbf{v} \in \mathbb{R}^n$ , we have

$$\mathbf{v}^\top (\mathbf{I} - \eta\sqrt{\mathbf{Q}}\Theta\sqrt{\mathbf{Q}})\mathbf{v} = 1 - \eta\mathbf{v}^\top \sqrt{\mathbf{Q}}\Theta\sqrt{\mathbf{Q}}\mathbf{v} \tag{60}$$

168  $\left\| \sqrt{\mathbf{Q}}\mathbf{v} \right\|_2 \in [\sqrt{q^*}, 1]$ , so for any  $\eta \leq \eta^*$ ,  $\mathbf{v}^\top (\mathbf{I} - \eta\sqrt{\mathbf{Q}}\Theta\sqrt{\mathbf{Q}})\mathbf{v} \in [0, 1 - \eta\lambda^{\min}q^*]$ , which implies  
169 that  $\left\| \mathbf{I} - \eta\sqrt{\mathbf{Q}}\Theta\sqrt{\mathbf{Q}} \right\|_2 \leq 1 - \eta\lambda^{\min}q^*$ . Thus,

$$\begin{aligned}
&\left\| \mathbf{I} - \eta\sqrt{\mathbf{Q}}J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)})\sqrt{\mathbf{Q}} \right\|_2 \\
&\leq \left\| \mathbf{I} - \eta\sqrt{\mathbf{Q}}\Theta\sqrt{\mathbf{Q}} \right\|_2 + \eta \left\| \sqrt{\mathbf{Q}}(\Theta - \Theta^{(0)})\sqrt{\mathbf{Q}} \right\|_2 + \eta \left\| \sqrt{\mathbf{Q}}(J(\theta^{(0)})^\top J(\theta^{(0)}) - J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)}))\sqrt{\mathbf{Q}} \right\|_2 \\
&\leq 1 - \eta\lambda^{\min}q^* + \eta \left\| \sqrt{\mathbf{Q}}(\Theta - \Theta^{(0)})\sqrt{\mathbf{Q}} \right\|_F + \eta \left\| \sqrt{\mathbf{Q}}(J(\theta^{(0)})^\top J(\theta^{(0)}) - J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)}))\sqrt{\mathbf{Q}} \right\|_F \\
&\leq 1 - \eta\lambda^{\min}q^* + \eta \left\| \Theta - \Theta^{(0)} \right\|_F + \eta \left\| J(\theta^{(0)})^\top J(\theta^{(0)}) - J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)}) \right\|_F \\
&\leq 1 - \eta\lambda^{\min}q^* + \frac{\eta q^* \lambda^{\min}}{3} + \frac{\eta M^2}{\sqrt[4]{d}} \left( \left\| \theta^{(t)} - \theta^{(0)} \right\|_2 + \left\| \tilde{\theta}^{(t)} - \theta^{(0)} \right\|_2 \right) \leq 1 - \frac{\eta q^* \lambda^{\min}}{2}
\end{aligned} \tag{61}$$

170 for all  $\tilde{d} \geq \max \left\{ D_1, D_2, \left( \frac{12M^2C_0}{q^*\lambda^{\min}} \right)^4 \right\}$ , which implies that

$$\begin{aligned}
&\left\| \mathbf{I} - \eta\sqrt{\mathbf{Q}}J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)})\sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} \right\|_2 \\
&\leq 1 - \frac{\eta q^* \lambda^{\min}}{2} + \left\| \eta\sqrt{\mathbf{Q}}J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)}) \left( \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} - \sqrt{\mathbf{Q}} \right) \right\|_2 \\
&\leq 1 - \frac{\eta q^* \lambda^{\min}}{2} + \eta M^2 \sqrt{3\epsilon} \leq 1 - \frac{\eta q^* \lambda^{\min}}{3} \quad (\text{due to (37)})
\end{aligned} \tag{62}$$

171 for all  $\epsilon \leq \epsilon_0$ . Thus, (50) is also true for  $t + 1$ . In conclusion, (50) and (51) are true with probability  
172 at least  $(1 - \delta)$  for all  $\tilde{d} \geq \tilde{D} = \max \left\{ D_1, D_2, \left( \frac{12M^2C_0}{q^*\lambda^{\min}} \right)^4 \right\}$ .  $\square$

173 Returning back to the proof of Lemma 5. Choose and fix an  $\epsilon$  such that  $\epsilon <$   
174  $\min\{\epsilon_0, \frac{1}{3} \left( \frac{q^* \lambda^{\min}}{3\lambda_{\max} + q^* \lambda^{\min}} \right)^2\}$ , where  $\epsilon_0$  is defined by Theorem 12. Then,  $t_\epsilon$  is also fixed. There  
175 exists  $\tilde{D} \geq 0$  such that for any  $\tilde{d} \geq \tilde{D}$ , with probability at least  $(1 - \delta)$ , Theorem 12 and Lemma 13  
176 are true and

$$\left\| \Theta - \Theta^{(0)} \right\|_F \leq \frac{q^* \lambda^{\min}}{3} \tag{63}$$



177 which immediately implies that

$$\left\| \Theta^{(0)} \right\|_2 \leq \left\| \Theta \right\|_2 + \left\| \Theta - \Theta^{(0)} \right\|_F \leq \lambda^{\max} + \frac{q^* \lambda^{\min}}{3} \quad (64)$$

178 We still denote  $B = 1 + \eta^* M^2$  and  $C_0 = \frac{MB^{t_\epsilon} R_0}{B-1} + \frac{3\sqrt{1+3\epsilon} MB^{t_\epsilon} R_0}{q^* \lambda^{\min}}$ . Theorem 12 ensures that for  
 179 all  $t, \theta^{(t)} \in B(\theta^{(0)}, C_0)$ . Then we have

$$\begin{aligned} \left\| \mathbf{I} - \eta \sqrt{\mathbf{Q}} \Theta^{(0)} \sqrt{\mathbf{Q}} \right\|_2 &\leq \left\| \mathbf{I} - \eta \sqrt{\mathbf{Q}} \Theta \sqrt{\mathbf{Q}} \right\|_2 + \eta \left\| \sqrt{\mathbf{Q}} (\Theta - \Theta^{(0)}) \sqrt{\mathbf{Q}} \right\|_2 \\ &\leq 1 - \eta \lambda^{\min} q^* + \frac{\eta q^* \lambda^{\min}}{3} = 1 - \frac{2\eta q^* \lambda^{\min}}{3} \end{aligned} \quad (65)$$

180 so it follows that

$$\begin{aligned} \left\| \mathbf{I} - \eta \sqrt{\mathbf{Q}} \Theta^{(0)} \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} \right\|_2 &\leq \left\| \mathbf{I} - \eta \sqrt{\mathbf{Q}} \Theta^{(0)} \sqrt{\mathbf{Q}} \right\|_2 + \left\| \eta \sqrt{\mathbf{Q}} \Theta^{(0)} \left( \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} - \sqrt{\mathbf{Q}} \right) \right\|_2 \\ &\leq 1 - \frac{2\eta q^* \lambda^{\min}}{3} + \eta \left( \lambda^{\max} + \frac{q^* \lambda^{\min}}{3} \right) \sqrt{3\epsilon} \end{aligned} \quad (66)$$

181 Thus, for all  $\epsilon < \frac{1}{3} \left( \frac{q^* \lambda^{\min}}{3\lambda^{\max} + q^* \lambda^{\min}} \right)^2$ , there is

$$\left\| \mathbf{I} - \eta \sqrt{\mathbf{Q}} \Theta^{(0)} \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} \right\|_2 \leq 1 - \frac{\eta q^* \lambda^{\min}}{3} \quad (67)$$

182 The update rule of the GRW for the linearized neural network is:

$$\theta_{\text{lin}}^{(t+1)} = \theta_{\text{lin}}^{(t)} - \eta J(\theta^{(0)}) \mathbf{Q}^{(t)} g_{\text{lin}}(\theta^{(t)}) \quad (68)$$

183 where we use the subscript “lin” to denote the linearized neural network, and with a slight abuse of  
 184 notion denote  $g_{\text{lin}}(\theta^{(t)}) = g(\theta_{\text{lin}}^{(t)})$ .

185 First, let us consider the training data  $\mathbf{X}$ . Denote  $\Delta_t = g_{\text{lin}}(\theta^{(t)}) - g(\theta^{(t)})$ . We have

$$\begin{cases} g_{\text{lin}}(\theta^{(t+1)}) - g_{\text{lin}}(\theta^{(t)}) = -\eta J(\theta^{(0)})^\top J(\theta^{(0)}) \mathbf{Q}^{(t)} g_{\text{lin}}(\theta^{(t)}) \\ g(\theta^{(t+1)}) - g(\theta^{(t)}) = -\eta J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)}) \mathbf{Q}^{(t)} g(\theta^{(t)}) \end{cases} \quad (69)$$

186 where  $\tilde{\theta}^{(t)}$  is some linear interpolation between  $\theta^{(t)}$  and  $\theta^{(t+1)}$ . Thus,

$$\begin{aligned} \Delta_{t+1} - \Delta_t &= \eta \left[ J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)}) - J(\theta^{(0)})^\top J(\theta^{(0)}) \right] \mathbf{Q}^{(t)} g(\theta^{(t)}) \\ &\quad - \eta J(\theta^{(0)})^\top J(\theta^{(0)}) \mathbf{Q}^{(t)} \Delta_t \end{aligned} \quad (70)$$

187 By Lemma 13, we have

$$\begin{aligned} &\left\| J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)}) - J(\theta^{(0)})^\top J(\theta^{(0)}) \right\|_F \\ &\leq \left\| \left( J(\tilde{\theta}^{(t)}) - J(\theta^{(0)}) \right)^\top J(\theta^{(t)}) \right\|_F + \left\| J(\theta^{(0)})^\top \left( J(\theta^{(t)}) - J(\theta^{(0)}) \right) \right\|_F \\ &\leq 2M^2 C_0 \tilde{d}^{-1/4} \end{aligned} \quad (71)$$

188 which implies that for all  $t < t_\epsilon$ ,

$$\begin{aligned} \|\Delta_{t+1}\|_2 &\leq \left\| \left[ \mathbf{I} - \eta J(\theta^{(0)})^\top J(\theta^{(0)}) \mathbf{Q}^{(t)} \right] \Delta_t \right\|_2 + \left\| \eta \left[ J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)}) - J(\theta^{(0)})^\top J(\theta^{(0)}) \right] \mathbf{Q}^{(t)} g(\theta^{(t)}) \right\|_2 \\ &\leq \left\| \mathbf{I} - \eta J(\theta^{(0)})^\top J(\theta^{(0)}) \mathbf{Q}^{(t)} \right\|_F \|\Delta_t\|_2 + \eta \left\| J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)}) - J(\theta^{(0)})^\top J(\theta^{(0)}) \right\|_F \|g(\theta^{(t)})\|_2 \\ &\leq (1 + \eta M^2) \|\Delta_t\|_2 + 2\eta M^2 C_0 B^t R_0 \tilde{d}^{-1/4} \\ &\leq B \|\Delta_t\|_2 + 2\eta M^2 C_0 B^t R_0 \tilde{d}^{-1/4} \end{aligned} \quad (72)$$

189 Therefore, we have

$$B^{-(t+1)} \|\Delta_{t+1}\|_2 \leq B^{-t} \|\Delta_t\|_2 + 2\eta M^2 C_0 B^{-1} R_0 \tilde{d}^{-1/4} \quad (73)$$

190 Since  $\Delta_0 = 0$ , it follows that for all  $t \leq t_\epsilon$ ,

$$\|\Delta_t\|_2 \leq 2t\eta M^2 C_0 B^{t-1} R_0 \tilde{d}^{-1/4} \quad (74)$$

191 and particularly we have

$$\left\| \sqrt{\mathbf{Q}} \Delta_{t_\epsilon} \right\|_2 \leq \|\Delta_{t_\epsilon}\|_2 \leq 2t_\epsilon \eta M^2 C_0 B^{t_\epsilon-1} R_0 \tilde{d}^{-1/4} \quad (75)$$

192 For  $t \geq t_\epsilon$ , we have the alternative update rule (47). Thus,

$$\begin{aligned} \sqrt{\mathbf{Q}} \Delta_{t+1} - \sqrt{\mathbf{Q}} \Delta_t &= \eta \sqrt{\mathbf{Q}} \left[ J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)}) - J(\theta^{(0)})^\top J(\theta^{(0)}) \right] \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} \left[ \sqrt{\mathbf{Q}} g(\theta^{(t)}) \right] \\ &\quad - \eta \sqrt{\mathbf{Q}} J(\theta^{(0)})^\top J(\theta^{(0)}) \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} \left[ \sqrt{\mathbf{Q}} \Delta_t \right] \end{aligned} \quad (76)$$

193 Let  $\mathbf{A} = \mathbf{I} - \eta \sqrt{\mathbf{Q}} J(\theta^{(0)})^\top J(\theta^{(0)}) \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} = \mathbf{I} - \eta \sqrt{\mathbf{Q}} \Theta^{(0)} \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}}$ . Then, we have

$$\sqrt{\mathbf{Q}} \Delta_{t+1} = \mathbf{A} \sqrt{\mathbf{Q}} \Delta_t + \eta \sqrt{\mathbf{Q}} \left[ J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)}) - J(\theta^{(0)})^\top J(\theta^{(0)}) \right] \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} \left( \sqrt{\mathbf{Q}} g(\theta^{(t)}) \right) \quad (77)$$

194 Let  $\gamma = 1 - \frac{\eta q^* \lambda^{\min}}{3} < 1$ . Combining with Theorem 12 and (67), the above leads to

$$\begin{aligned} \left\| \sqrt{\mathbf{Q}} \Delta_{t+1} \right\|_2 &\leq \|\mathbf{A}\|_2 \left\| \sqrt{\mathbf{Q}} \Delta_t \right\|_2 + \eta \left\| \sqrt{\mathbf{Q}} \left[ J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)}) - J(\theta^{(0)})^\top J(\theta^{(0)}) \right] \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}} \right\|_2 \left\| \sqrt{\mathbf{Q}} g(\theta^{(t)}) \right\|_2 \\ &\leq \gamma \left\| \sqrt{\mathbf{Q}} \Delta_t \right\|_2 + \eta \left\| J(\tilde{\theta}^{(t)})^\top J(\theta^{(t)}) - J(\theta^{(0)})^\top J(\theta^{(0)}) \right\|_F \sqrt{1 + 3\epsilon} \gamma^{t-t_\epsilon} B^{t_\epsilon} R_0 \\ &\leq \gamma \left\| \sqrt{\mathbf{Q}} \Delta_t \right\|_2 + 2\eta M^2 C_0 \sqrt{1 + 3\epsilon} \gamma^{t-t_\epsilon} B^{t_\epsilon} R_0 \tilde{d}^{-1/4} \end{aligned} \quad (78)$$

195 This implies that

$$\gamma^{-(t+1)} \left\| \sqrt{\mathbf{Q}} \Delta_{t+1} \right\|_2 \leq \gamma^{-t} \left\| \sqrt{\mathbf{Q}} \Delta_t \right\|_2 + 2\eta M^2 C_0 \sqrt{1 + 3\epsilon} \gamma^{-1-t_\epsilon} B^{t_\epsilon} R_0 \tilde{d}^{-1/4} \quad (79)$$

196 Combining with (75), it implies that for all  $t \geq t_\epsilon$ ,

$$\left\| \sqrt{\mathbf{Q}} \Delta_t \right\|_2 \leq 2\gamma^{t-t_\epsilon} \eta M^2 C_0 B^{t_\epsilon} R_0 \left[ t_\epsilon B^{-1} + \sqrt{1 + 3\epsilon} \gamma^{-1} (t - t_\epsilon) \right] \tilde{d}^{-1/4} \quad (80)$$

197 Next, we consider an arbitrary test point  $\mathbf{x}$  such that  $\|\mathbf{x}\|_2 \leq 1$ . Denote  $\delta_t = f_{\text{lin}}^{(t)}(\mathbf{x}) - f^{(t)}(\mathbf{x})$ .

198 Then we have

$$\begin{cases} f_{\text{lin}}^{(t+1)}(\mathbf{x}) - f_{\text{lin}}^{(t)}(\mathbf{x}) = -\eta \nabla_\theta f(\mathbf{x}; \theta^{(0)})^\top J(\theta^{(0)}) \mathbf{Q}^{(t)} g_{\text{lin}}(\theta^{(t)}) \\ f^{(t+1)}(\mathbf{x}) - f^{(t)}(\mathbf{x}) = -\eta \nabla_\theta f(\mathbf{x}; \tilde{\theta}^{(t)})^\top J(\theta^{(t)}) \mathbf{Q}^{(t)} g(\theta^{(t)}) \end{cases} \quad (81)$$

199 which yields

$$\begin{aligned} \delta_{t+1} - \delta_t &= \eta \left[ \nabla_\theta f(\mathbf{x}; \tilde{\theta}^{(t)})^\top J(\theta^{(t)}) - \nabla_\theta f(\mathbf{x}; \theta^{(0)})^\top J(\theta^{(0)}) \right] \mathbf{Q}^{(t)} g(\theta^{(t)}) \\ &\quad - \eta \nabla_\theta f(\mathbf{x}; \theta^{(0)})^\top J(\theta^{(0)}) \mathbf{Q}^{(t)} \Delta_t \end{aligned} \quad (82)$$

200 For  $t \leq t_\epsilon$ , we have

$$\begin{aligned}
\|\delta_t\|_2 &\leq \eta \sum_{s=0}^{t-1} \left\| \left[ \nabla_\theta f(\mathbf{x}; \tilde{\theta}^{(s)})^\top J(\theta^{(s)}) - \nabla_\theta f(\mathbf{x}; \theta^{(0)})^\top J(\theta^{(0)}) \right] \mathbf{Q}^{(s)} \right\|_2 \|g(\theta^{(s)})\|_2 \\
&\quad + \eta \sum_{s=0}^{t-1} \left\| \nabla_\theta f(\mathbf{x}; \theta^{(0)})^\top J(\theta^{(0)}) \mathbf{Q}^{(s)} \right\|_2 \|\Delta_s\|_2 \\
&\leq \eta \sum_{s=0}^{t-1} \left\| \nabla_\theta f(\mathbf{x}; \tilde{\theta}^{(s)})^\top J(\theta^{(s)}) - \nabla_\theta f(\mathbf{x}; \theta^{(0)})^\top J(\theta^{(0)}) \right\|_F \|g(\theta^{(s)})\|_2 \\
&\quad + \eta \sum_{s=0}^{t-1} \left\| \nabla_\theta f(\mathbf{x}; \theta^{(0)}) \right\|_2 \|J(\theta^{(0)})\|_F \|\Delta_s\|_2 \\
&\leq 2\eta M^2 C_0 \tilde{d}^{-1/4} \sum_{s=0}^{t-1} B^s R_0 + \eta M^2 \sum_{s=0}^{t-1} (2s\eta M^2 C_0 B^{s-1} R_0 \tilde{d}^{-1/4})
\end{aligned} \tag{83}$$

201 So we can see that there exists a constant  $C_1$  such that  $\|\delta_{t_\epsilon}\|_2 \leq C_1 \tilde{d}^{-1/4}$ . Then, for  $t > t_\epsilon$ , we have

$$\begin{aligned}
\|\delta_t\|_2 - \|\delta_{t_\epsilon}\|_2 &\leq \eta \sum_{s=t_\epsilon}^{t-1} \left\| \left[ \nabla_\theta f(\mathbf{x}; \tilde{\theta}^{(s)})^\top J(\theta^{(s)}) - \nabla_\theta f(\mathbf{x}; \theta^{(0)})^\top J(\theta^{(0)}) \right] \sqrt{\mathbf{Q}_{3\epsilon}^{(s)}} \right\|_2 \left\| \sqrt{\mathbf{Q}} g(\theta^{(s)}) \right\|_2 \\
&\quad + \eta \sum_{s=t_\epsilon}^{t-1} \left\| \nabla_\theta f(\mathbf{x}; \theta^{(0)})^\top J(\theta^{(0)}) \sqrt{\mathbf{Q}_{3\epsilon}^{(s)}} \right\|_2 \left\| \sqrt{\mathbf{Q}} \Delta_s \right\|_2 \\
&\leq 2\eta M^2 C_0 \tilde{d}^{-1/4} \sqrt{1+3\epsilon} \sum_{s=t_\epsilon}^{t-1} \gamma^{s-t_\epsilon} B^{t_\epsilon} R_0 \\
&\quad + \eta M^2 \sqrt{1+3\epsilon} \sum_{s=t_\epsilon}^{t-1} \left( 2\gamma^{s-t_\epsilon} \eta M^2 C_0 B^{t_\epsilon} R_0 [t_\epsilon B^{-1} + \sqrt{1+3\epsilon} \gamma^{-1}(s-t_\epsilon)] \tilde{d}^{-1/4} \right)
\end{aligned} \tag{84}$$

202 Note that  $\sum_{t=0}^\infty t\gamma^t$  is finite as long as  $\gamma \in (0, 1)$ . Therefore, there is a constant  $C$  such that for any  $t$ ,  
203  $\|\delta_t\|_2 \leq C \tilde{d}^{-1/4}$  with probability at least  $(1-\delta)$  for any  $\tilde{d} \geq \tilde{D}$ .  $\square$

### 204 D.3.3 Proof of Lemma 13

205 We will use the following theorem regarding the eigenvalues of random Gaussian matrices:

206 **Theorem 14** (Corollary 5.35 in [Ver10]). *If  $\mathbf{A} \in \mathbb{R}^{p \times q}$  is a random matrix whose entries are*  
207 *independent standard normal random variables, then for every  $t \geq 0$ , with probability at least*  
208  *$1 - 2\exp(-t^2/2)$ ,*

$$\sqrt{p} - \sqrt{q} - t \leq \lambda^{\min}(\mathbf{A}) \leq \lambda^{\max}(\mathbf{A}) \leq \sqrt{p} + \sqrt{q} + t \tag{85}$$

209 By this theorem, and also note that  $W^L$  is a vector, we can see that for any  $\delta$ , there exist  $\tilde{D} > 0$  and  
210  $M_1 > 0$  such that if  $\tilde{d} \geq \tilde{D}$ , then with probability at least  $(1-\delta)$ , for all  $\theta \in B(\theta^{(0)}, C_0)$ , we have

$$\|W^l\|_2 \leq 3\sqrt{\tilde{d}} \quad (\forall 0 \leq l \leq L-1) \quad \text{and} \quad \|W^L\|_2 \leq C_0 \leq 3\sqrt[4]{\tilde{d}} \tag{86}$$

211 as well as

$$\|\beta \mathbf{b}^l\|_2 \leq M_1 \sqrt{\tilde{d}} \quad (\forall l = 0, \dots, L) \tag{87}$$

212 Now we assume that (86) and (87) are true. Then, for any  $\mathbf{x}$  such that  $\|\mathbf{x}\|_2 \leq 1$ ,

$$\begin{aligned}
\|\mathbf{h}^1\|_2 &= \left\| \frac{1}{\sqrt{d_0}} W^0 \mathbf{x} + \beta \mathbf{b}^0 \right\|_2 \leq \frac{1}{\sqrt{d_0}} \|W^0\|_2 \|\mathbf{x}\|_2 + \|\beta \mathbf{b}^0\|_2 \leq \left( \frac{3}{\sqrt{d_0}} + M_1 \right) \sqrt{\tilde{d}} \\
\|\mathbf{h}^{l+1}\|_2 &= \left\| \frac{1}{\sqrt{\tilde{d}}} W^l \mathbf{x}^l + \beta \mathbf{b}^l \right\|_2 \leq \frac{1}{\sqrt{\tilde{d}}} \|W^l\|_2 \|\mathbf{x}^l\|_2 + \|\beta \mathbf{b}^l\|_2 \quad (\forall l \geq 1) \\
\|\mathbf{x}^l\|_2 &= \|\sigma(\mathbf{h}^l) - \sigma(\mathbf{0}^l) + \sigma(\mathbf{0}^l)\|_2 \leq L_0 \|\mathbf{h}^l\|_2 + \sigma(0) \sqrt{\tilde{d}} \quad (\forall l \geq 1)
\end{aligned} \tag{88}$$

213 where  $L_0$  is the Lipschitz constant of  $\sigma$  and  $\sigma(\mathbf{0}^l) = (\sigma(0), \dots, \sigma(0)) \in \mathbb{R}^{d_l}$ . By induction, there  
 214 exists an  $M_2 > 0$  such that  $\|\mathbf{x}^l\|_2 \leq M_2 \sqrt{\tilde{d}}$  and  $\|\mathbf{h}^l\|_2 \leq M_2 \sqrt{\tilde{d}}$  for all  $l = 1, \dots, L$ .

215 Denote  $\boldsymbol{\alpha}^l = \nabla_{\mathbf{h}^l} f(\mathbf{x}) = \nabla_{\mathbf{h}^l} \mathbf{h}^{L+1}$ . For all  $l = 1, \dots, L$ , we have  $\boldsymbol{\alpha}^l = \text{diag}(\dot{\sigma}(\mathbf{h}^l)) \frac{W^{l\top}}{\sqrt{\tilde{d}}} \boldsymbol{\alpha}^{l+1}$   
 216 where  $\dot{\sigma}(x) \leq L_0$  for all  $x \in \mathbb{R}$  since  $\sigma$  is  $L_0$ -Lipschitz,  $\boldsymbol{\alpha}^{L+1} = \mathbf{1}$  and  $\|\boldsymbol{\alpha}^L\|_2 =$   
 217  $\left\| \text{diag}(\dot{\sigma}(\mathbf{h}^L)) \frac{W^{L\top}}{\sqrt{\tilde{d}}} \right\|_2 \leq \frac{3}{\sqrt[4]{\tilde{d}}} L_0$ . Then, we can easily prove by induction that there exists an  
 218  $M_3 > 1$  such that  $\|\boldsymbol{\alpha}^l\|_2 \leq M_3 / \sqrt[4]{\tilde{d}}$  for all  $l = 1, \dots, L$  (note that this is not true for  $L+1$  because  
 219  $\boldsymbol{\alpha}^{L+1} = \mathbf{1}$ ).

220 For  $l = 0$ ,  $\nabla_{W^0} f(\mathbf{x}) = \frac{1}{\sqrt{d_0}} \mathbf{x}^0 \boldsymbol{\alpha}^{1\top}$ , so  $\|\nabla_{W^0} f(\mathbf{x})\|_2 \leq \frac{1}{\sqrt{d_0}} \|\mathbf{x}^0\|_2 \|\boldsymbol{\alpha}^1\|_2 \leq \frac{1}{\sqrt{d_0}} M_3 / \sqrt[4]{\tilde{d}}$ . And  
 221 for any  $l = 1, \dots, L$ ,  $\nabla_{W^l} f(\mathbf{x}) = \frac{1}{\sqrt{\tilde{d}}} \mathbf{x}^l \boldsymbol{\alpha}^{l+1\top}$ , so  $\|\nabla_{W^l} f(\mathbf{x})\|_2 \leq \frac{1}{\sqrt{\tilde{d}}} \|\mathbf{x}^l\|_2 \|\boldsymbol{\alpha}^{l+1}\|_2 \leq M_2 M_3$ .  
 222 (Note that if  $M_3 > 1$ , then  $\|\boldsymbol{\alpha}^{L+1}\|_2 \leq M_3$ ; and since  $\tilde{d} \geq 1$ , there is  $\|\boldsymbol{\alpha}^l\|_2 \leq M_3$  for  $l \leq L$ .)  
 223 Moreover, for  $l = 0, \dots, L$ ,  $\nabla_{\mathbf{b}^l} f(\mathbf{x}) = \beta \boldsymbol{\alpha}^{l+1}$ , so  $\|\nabla_{\mathbf{b}^l} f(\mathbf{x})\|_2 \leq \beta M_3$ . Thus, if (86) and (87) are  
 224 true, then there exists an  $M_4 > 0$ , such that  $\|\nabla_{\theta} f(\mathbf{x})\|_2 \leq M_4 / \sqrt{n}$ . And since  $\|\mathbf{x}_i\|_2 \leq 1$  for all  $i$ ,  
 225 so  $\|J(\theta)\|_F \leq M_4$ .

226 Next, we consider the difference in  $\nabla_{\theta} f(\mathbf{x})$  between  $\theta$  and  $\tilde{\theta}$ . Let  $\tilde{f}, \tilde{W}, \tilde{\mathbf{b}}, \tilde{\mathbf{x}}, \tilde{\mathbf{h}}, \tilde{\boldsymbol{\alpha}}$  be the function  
 227 and the values corresponding to  $\tilde{\theta}$ . There is

$$\begin{aligned} \|\mathbf{h}^1 - \tilde{\mathbf{h}}^1\|_2 &= \left\| \frac{1}{\sqrt{d_0}} (W^0 - \tilde{W}^0) \mathbf{x} + \beta (\mathbf{b}^0 - \tilde{\mathbf{b}}^0) \right\|_2 \\ &\leq \frac{1}{\sqrt{d_0}} \|W^0 - \tilde{W}^0\|_2 \|\mathbf{x}\|_2 + \beta \|\mathbf{b}^0 - \tilde{\mathbf{b}}^0\|_2 \leq \left( \frac{1}{\sqrt{d_0}} + \beta \right) \|\theta - \tilde{\theta}\|_2 \\ \|\mathbf{h}^{l+1} - \tilde{\mathbf{h}}^{l+1}\|_2 &= \left\| \frac{1}{\sqrt{\tilde{d}}} W^l (\mathbf{x}^l - \tilde{\mathbf{x}}^l) + \frac{1}{\sqrt{\tilde{d}}} (W^l - \tilde{W}^l) \tilde{\mathbf{x}}^l + \beta (\mathbf{b}^l - \tilde{\mathbf{b}}^l) \right\|_2 \\ &\leq \frac{1}{\sqrt{\tilde{d}}} \|W^l\|_2 \|\mathbf{x}^l - \tilde{\mathbf{x}}^l\|_2 + \frac{1}{\sqrt{\tilde{d}}} \|W^l - \tilde{W}^l\|_2 \|\tilde{\mathbf{x}}^l\|_2 + \beta \|\mathbf{b}^l - \tilde{\mathbf{b}}^l\|_2 \\ &\leq 3 \|\mathbf{x}^l - \tilde{\mathbf{x}}^l\|_2 + (M_2 + \beta) \|\theta - \tilde{\theta}\|_2 \quad (\forall l \geq 1) \\ \|\mathbf{x}^l - \tilde{\mathbf{x}}^l\|_2 &= \|\sigma(\mathbf{h}^l) - \sigma(\tilde{\mathbf{h}}^l)\|_2 \leq L_0 \|\mathbf{h}^l - \tilde{\mathbf{h}}^l\|_2 \quad (\forall l \geq 1) \end{aligned} \quad (89)$$

228 By induction, there exists an  $M_5 > 0$  such that  $\|\mathbf{x}^l - \tilde{\mathbf{x}}^l\|_2 \leq M_5 \|\theta - \tilde{\theta}\|_2$  for all  $l$ .

229 For  $\boldsymbol{\alpha}^l$ , we have  $\boldsymbol{\alpha}^{L+1} = \tilde{\boldsymbol{\alpha}}^{L+1} = \mathbf{1}$ , and for all  $l \geq 1$ ,

$$\begin{aligned} \|\boldsymbol{\alpha}^l - \tilde{\boldsymbol{\alpha}}^l\|_2 &= \left\| \text{diag}(\dot{\sigma}(\mathbf{h}^l)) \frac{W^{l\top}}{\sqrt{\tilde{d}}} \boldsymbol{\alpha}^{l+1} - \text{diag}(\dot{\sigma}(\tilde{\mathbf{h}}^l)) \frac{\tilde{W}^{l\top}}{\sqrt{\tilde{d}}} \tilde{\boldsymbol{\alpha}}^{l+1} \right\|_2 \\ &\leq \left\| \text{diag}(\dot{\sigma}(\mathbf{h}^l)) \frac{W^{l\top}}{\sqrt{\tilde{d}}} (\boldsymbol{\alpha}^{l+1} - \tilde{\boldsymbol{\alpha}}^{l+1}) \right\|_2 + \left\| \text{diag}(\dot{\sigma}(\mathbf{h}^l)) \frac{(W^l - \tilde{W}^l)^\top}{\sqrt{\tilde{d}}} \tilde{\boldsymbol{\alpha}}^{l+1} \right\|_2 \\ &\quad + \left\| \text{diag}((\dot{\sigma}(\mathbf{h}^l) - \dot{\sigma}(\tilde{\mathbf{h}}^l))) \frac{\tilde{W}^{l\top}}{\sqrt{\tilde{d}}} \tilde{\boldsymbol{\alpha}}^{l+1} \right\|_2 \\ &\leq 3L_0 \|\boldsymbol{\alpha}^{l+1} - \tilde{\boldsymbol{\alpha}}^{l+1}\|_2 + \left( M_3 L_0 \tilde{d}^{-1/2} + 3M_3 M_5 L_1 \tilde{d}^{-1/4} \right) \|\theta - \tilde{\theta}\|_2 \end{aligned} \quad (90)$$

230 where  $L_1$  is the Lipschitz constant of  $\dot{\sigma}$ . Particularly, for  $l = L$ , though  $\tilde{\boldsymbol{\alpha}}^{L+1} = \mathbf{1}$ , since  $\|\tilde{W}^L\|_2 \leq$   
 231  $3\tilde{d}^{1/4}$ , (90) is still true. By induction, there exists an  $M_6 > 0$  such that  $\|\boldsymbol{\alpha}^l - \tilde{\boldsymbol{\alpha}}^l\|_2 \leq \frac{M_6}{\sqrt[4]{\tilde{d}}} \|\theta - \tilde{\theta}\|_2$   
 232 for all  $l \geq 1$  (note that this is also true for  $l = L+1$ ).

Thus, if (86) and (87) are true, then for all  $\theta, \tilde{\theta} \in B(\theta^{(0)}, C_0)$ , any  $\mathbf{x}$  such that  $\|\mathbf{x}\|_2 \leq 1$ , we have

$$\begin{aligned} \left\| \nabla_{W^0} f(\mathbf{x}) - \nabla_{\tilde{W}^0} \tilde{f}(\mathbf{x}) \right\|_2 &= \frac{1}{\sqrt{d_0}} \left\| \mathbf{x} \boldsymbol{\alpha}^{1\top} - \mathbf{x} \tilde{\boldsymbol{\alpha}}^{1\top} \right\|_2 \\ &\leq \frac{1}{\sqrt{d_0}} \left\| \boldsymbol{\alpha}^1 - \tilde{\boldsymbol{\alpha}}^1 \right\|_2 \\ &\leq \frac{1}{\sqrt{d_0}} \frac{M_6}{\sqrt[4]{\tilde{d}}} \left\| \theta - \tilde{\theta} \right\|_2 \end{aligned} \quad (91)$$

and for  $l = 1, \dots, L$ , we have

$$\begin{aligned} \left\| \nabla_{W^l} f(\mathbf{x}) - \nabla_{\tilde{W}^l} \tilde{f}(\mathbf{x}) \right\|_2 &= \frac{1}{\sqrt{\tilde{d}}} \left\| \mathbf{x}^l \boldsymbol{\alpha}^{l+1\top} - \tilde{\mathbf{x}}^l \tilde{\boldsymbol{\alpha}}^{l+1\top} \right\|_2 \\ &\leq \frac{1}{\sqrt{\tilde{d}}} (\left\| \mathbf{x}^l \right\|_2 \left\| \boldsymbol{\alpha}^{l+1} - \tilde{\boldsymbol{\alpha}}^{l+1} \right\|_2 + \left\| \mathbf{x}^l - \tilde{\mathbf{x}}^l \right\|_2 \left\| \tilde{\boldsymbol{\alpha}}^{l+1} \right\|_2) \\ &\leq \left( \frac{M_2 M_6}{\sqrt[4]{\tilde{d}}} + \frac{M_5 M_3}{\sqrt{\tilde{d}}} \right) \left\| \theta - \tilde{\theta} \right\|_2 \end{aligned} \quad (92)$$

Moreover, for any  $l = 0, \dots, L$ , there is

$$\left\| \nabla_{b^l} f(\mathbf{x}) - \nabla_{\tilde{b}^l} \tilde{f}(\mathbf{x}) \right\|_2 = \beta \left\| \boldsymbol{\alpha}^{l+1} - \tilde{\boldsymbol{\alpha}}^{l+1} \right\|_2 \leq \frac{\beta M_6}{\sqrt[4]{\tilde{d}}} \left\| \theta - \tilde{\theta} \right\|_2 \quad (93)$$

Overall, we can see that there exists a constant  $M_7 > 0$  such that  $\left\| \nabla_{\theta} f(\mathbf{x}) - \nabla_{\tilde{\theta}} \tilde{f}(\mathbf{x}) \right\|_2 \leq \frac{M_7}{\sqrt{n} \cdot \sqrt[4]{\tilde{d}}} \left\| \theta - \tilde{\theta} \right\|_2$ , so that  $\left\| J(\theta) - J(\tilde{\theta}) \right\|_F \leq \frac{M_7}{\sqrt[4]{\tilde{d}}} \left\| \theta - \tilde{\theta} \right\|_2$ .  $\square$

#### D.3.4 Proof of Theorem 4

First of all, for a linearized neural network (11), if we view  $\{\nabla_{\theta} f^{(0)}(\mathbf{x}_i)\}_{i=1}^n$  as the inputs and  $\{y_i - f^{(0)}(\mathbf{x}_i) + \langle \theta^{(0)}, \nabla_{\theta} f^{(0)}(\mathbf{x}_i) \rangle\}_{i=1}^n$  as the targets, then the model becomes a linear model. So by Theorem 2 we have the following corollary:

**Corollary 15.** *If  $\nabla_{\theta} f^{(0)}(\mathbf{x}_1), \dots, \nabla_{\theta} f^{(0)}(\mathbf{x}_n)$  are linearly independent, then there exists  $\eta_0 > 0$  such that for any GRW satisfying Assumption 1, and any  $\eta \leq \eta_0$ ,  $\theta^{(t)}$  converges to the same interpolator  $\theta^*$  that does not depend on  $q_i$ .*

Let  $\eta_1 = \min\{\eta_0, \eta^*\}$ , where  $\eta_0$  is defined in Corollary 15 and  $\eta^*$  is defined in Lemma 5. Let  $f_{\text{lin}}^{(t)}(\mathbf{x})$  and  $f_{\text{linERM}}^{(t)}(\mathbf{x})$  be the linearized neural networks of  $f^{(t)}(\mathbf{x})$  and  $f_{\text{ERM}}^{(t)}(\mathbf{x})$ , respectively. By Lemma 5, for any  $\delta > 0$ , there exists  $\tilde{D} > 0$  and a constant  $C$  such that

$$\begin{cases} \sup_{t \geq 0} \left| f_{\text{lin}}^{(t)}(\mathbf{x}) - f^{(t)}(\mathbf{x}) \right| \leq C \tilde{d}^{-1/4} \\ \sup_{t \geq 0} \left| f_{\text{linERM}}^{(t)}(\mathbf{x}) - f_{\text{ERM}}^{(t)}(\mathbf{x}) \right| \leq C \tilde{d}^{-1/4} \end{cases} \quad (94)$$

By Corollary 15, we have

$$\lim_{t \rightarrow \infty} \left| f_{\text{lin}}^{(t)}(\mathbf{x}) - f_{\text{linERM}}^{(t)}(\mathbf{x}) \right| = 0 \quad (95)$$

Summing the above yields

$$\limsup_{t \rightarrow \infty} \left| f^{(t)}(\mathbf{x}) - f_{\text{ERM}}^{(t)}(\mathbf{x}) \right| \leq 2C \tilde{d}^{-1/4} \quad (96)$$

which is the result we want.  $\square$

## 251 D.4 Proofs for Subsection 4.3

### 252 D.4.1 A New Approximation Theorem

253 **Lemma 16** (Approximation Theorem for Regularized GRW). *For a wide fully-connected neural*  
 254 *network  $f$ , denote  $J(\theta) = \nabla_{\theta} f(\mathbf{X}; \theta) \in \mathbb{R}^{p \times n}$  and  $g(\theta) = \nabla_{\hat{y}} \ell(f(\mathbf{X}; \theta), \mathbf{Y}) \in \mathbb{R}^n$ . Given that the*  
 255 *loss function  $\ell$  satisfies:  $\nabla_{\theta} g(\theta) = J(\theta)U(\theta)$  for any  $\theta$ , and  $U(\theta)$  is a positive semi-definite diagonal*  
 256 *matrix whose elements are uniformly bounded, we have: for any GRW that minimizes the regularized*  
 257 *weighted empirical risk (13) with a sufficiently small learning rate  $\eta$ , there is: for a sufficiently large*  
 258  *$\tilde{d}$ , with high probability over random initialization, on any test point  $\mathbf{x}$  such that  $\|\mathbf{x}\|_2 \leq 1$ ,*

$$\sup_{t \geq 0} \left| f_{\text{linreg}}^{(t)}(\mathbf{x}) - f_{\text{reg}}^{(t)}(\mathbf{x}) \right| \leq C \tilde{d}^{-1/4} \quad (97)$$

259 *where both  $f_{\text{linreg}}^{(t)}$  and  $f_{\text{reg}}^{(t)}$  are trained by the same regularized GRW and start from the same initial*  
 260 *point.*

261 First of all, with some simple linear algebra analysis, we can prove the following proposition:

262 **Proposition 17.** *For any positive definite symmetric matrix  $\mathbf{H} \in \mathbb{R}^{n \times n}$ , denote its largest and*  
 263 *smallest eigenvalues by  $\lambda^{\max}$  and  $\lambda^{\min}$ . Then, for any positive semi-definite diagonal matrix*  
 264  *$\mathbf{Q} = \text{diag}(q_1, \dots, q_n)$ ,  $\mathbf{H}\mathbf{Q}$  has  $n$  eigenvalues that all lie in  $[\min_i q_i \cdot \lambda^{\min}, \max_i q_i \cdot \lambda^{\max}]$ .*

265 *Proof.*  $\mathbf{H}$  is a positive definite symmetric matrix, so there exists  $\mathbf{A} \in \mathbb{R}^{n \times n}$  such that  $\mathbf{H} = \mathbf{A}^{\top} \mathbf{A}$ ,  
 266 and  $\mathbf{A}$  is full-rank. First, any eigenvalue of  $\mathbf{A}\mathbf{Q}\mathbf{A}^{\top}$  is also an eigenvalue of  $\mathbf{A}^{\top} \mathbf{A}\mathbf{Q}$ , because for  
 267 any eigenvalue  $\lambda$  of  $\mathbf{A}\mathbf{Q}\mathbf{A}^{\top}$  we have some  $\mathbf{v} \neq 0$  such that  $\mathbf{A}\mathbf{Q}\mathbf{A}^{\top} \mathbf{v} = \lambda \mathbf{v}$ . Multiplying both sides  
 268 by  $\mathbf{A}^{\top}$  on the left yields  $\mathbf{A}^{\top} \mathbf{A}\mathbf{Q}(\mathbf{A}^{\top} \mathbf{v}) = \lambda(\mathbf{A}^{\top} \mathbf{v})$  which implies that  $\lambda$  is also an eigenvalue of  
 269  $\mathbf{A}^{\top} \mathbf{A}\mathbf{Q}$  because  $\mathbf{A}^{\top} \mathbf{v} \neq 0$  as  $\lambda \mathbf{v} \neq 0$ .

270 Second, by condition we know that the eigenvalues of  $\mathbf{A}^{\top} \mathbf{A}$  are all in  $[\lambda^{\min}, \lambda^{\max}]$  where  $\lambda^{\min} > 0$ ,  
 271 which implies for any unit vector  $\mathbf{v}$ ,  $\mathbf{v}^{\top} \mathbf{A}^{\top} \mathbf{A} \mathbf{v} \in [\lambda^{\min}, \lambda^{\max}]$ , which is equivalent to  $\|\mathbf{A}\mathbf{v}\|_2 \in$   
 272  $[\sqrt{\lambda^{\min}}, \sqrt{\lambda^{\max}}]$ . Thus, we have  $\mathbf{v}^{\top} \mathbf{A}^{\top} \mathbf{Q} \mathbf{A} \mathbf{v} \in [\lambda^{\min} \min_i q_i, \lambda^{\max} \max_i q_i]$ , which implies that  
 273 the eigenvalues of  $\mathbf{A}^{\top} \mathbf{Q} \mathbf{A}$  are all in  $[\lambda^{\min} \min_i q_i, \lambda^{\max} \max_i q_i]$ .

274 Thus, the eigenvalues of  $\mathbf{H}\mathbf{Q} = \mathbf{A}^{\top} \mathbf{A}\mathbf{Q}$  are all in  $[\lambda^{\min} \min_i q_i, \lambda^{\max} \max_i q_i]$ .  $\square$

275 **Proof of Lemma 16** By the condition  $\ell$  satisfies, without loss of generality, assume that the elements  
 276 of  $U(\theta)$  are in  $[0, 1]$  for all  $\theta$ . Then, let  $\eta \leq (\mu + \lambda^{\min} + \lambda^{\max})^{-1}$ . (If the elements of  $U(\theta)$  are  
 277 bounded by  $[0, C]$ , then we can let  $\eta \leq (\mu + C\lambda^{\min} + C\lambda^{\max})^{-1}$  and prove the result in the same  
 278 way.)

279 With  $L_2$  penalty, the update rule of the GRW for the neural network is:

$$\theta^{(t+1)} = \theta^{(t)} - \eta J(\theta^{(t)}) \mathbf{Q}^{(t)} g(\theta^{(t)}) - \eta \mu (\theta^{(t)} - \theta^{(0)}) \quad (98)$$

280 And the update rule for the linearized neural network is:

$$\theta_{\text{lin}}^{(t+1)} = \theta_{\text{lin}}^{(t)} - \eta J(\theta^{(0)}) \mathbf{Q}^{(t)} g(\theta_{\text{lin}}^{(t)}) - \eta \mu (\theta_{\text{lin}}^{(t)} - \theta^{(0)}) \quad (99)$$

281 By Proposition 11,  $f(\mathbf{x}; \theta)$  converges in probability to a zero-mean Gaussian process. Thus, for any  
 282  $\delta > 0$ , there exists a constant  $R_0 > 0$  such that with probability at least  $(1 - \delta/3)$ ,  $\|g(\theta^{(0)})\|_2 < R_0$ .  
 283 Let  $M$  be as defined in Lemma 13. Denote  $A = \eta M R_0$ , and let  $C_0 = \frac{4A}{\eta \mu}$  in Lemma 13<sup>1</sup>. By Lemma  
 284 13, there exists  $D_1$  such that for all  $\tilde{d} \geq D_1$ , with probability at least  $(1 - \delta/3)$ , (52) is true.

285 Similar to the proof of Proposition 17, we can show that for arbitrary  $\tilde{\theta}$ , all non-zero eigenval-  
 286 ues of  $J(\theta^{(0)}) \mathbf{Q}^{(t)} U(\tilde{\theta}) J(\theta^{(0)})^{\top}$  are eigenvalues of  $J(\theta^{(0)})^{\top} J(\theta^{(0)}) \mathbf{Q}^{(t)} U(\tilde{\theta})$ . This is because for  
 287 any  $\lambda \neq 0$ , if  $J(\theta^{(0)}) \mathbf{Q}^{(t)} U(\tilde{\theta}) J(\theta^{(0)})^{\top} \mathbf{v} = \lambda \mathbf{v}$ , then  $J(\theta^{(0)})^{\top} J(\theta^{(0)}) \mathbf{Q}^{(t)} U(\tilde{\theta}) (J(\theta^{(0)})^{\top} \mathbf{v}) =$   
 288  $\lambda (J(\theta^{(0)})^{\top} \mathbf{v})$ , and  $J(\theta^{(0)})^{\top} \mathbf{v} \neq 0$  since  $\lambda \mathbf{v} \neq 0$ , so  $\lambda$  is also an eigenvalue of

<sup>1</sup>Note that Lemma 13 only depends on the network structure and does not depend on the update rule, so we can use this lemma here.

289  $J(\theta^{(0)})^\top J(\theta^{(0)})\mathbf{Q}^{(t)}U(\tilde{\theta})$ . On the other hand, by Proposition 3,  $J(\theta^{(0)})^\top J(\theta^{(0)})\mathbf{Q}^{(t)}U(\tilde{\theta})$  con-  
 290 verges in probability to  $\Theta\mathbf{Q}^{(t)}U(\tilde{\theta})$  whose eigenvalues are all in  $[0, \lambda^{\max}]$  by Proposition 17. So  
 291 there exists  $D_2$  such that for all  $\tilde{d} \geq D_2$ , with probability at least  $(1 - \delta/3)$ , the eigenvalues of  
 292  $J(\theta^{(0)})\mathbf{Q}^{(t)}U(\tilde{\theta})J(\theta^{(0)})^\top$  are all in  $[0, \lambda^{\max} + \lambda^{\min}]$  for all  $t$ .

293 By union bound, with probability at least  $(1 - \delta)$ , all three above are true, which we will assume in  
 294 the rest of this proof.

295 First, we need to prove that there exists  $D_0$  such that for all  $\tilde{d} \geq D_0$ ,  $\sup_{t \geq 0} \|\theta^{(t)} - \theta^{(0)}\|_2$  is  
 296 bounded with high probability. Denote  $a_t = \theta^{(t)} - \theta^{(0)}$ . By (98) we have

$$\begin{aligned} a_{t+1} = & (1 - \eta\mu)a_t - \eta[J(\theta^{(t)}) - J(\theta^{(0)})]\mathbf{Q}^{(t)}g(\theta^{(t)}) \\ & - \eta J(\theta^{(0)})\mathbf{Q}^{(t)}[g(\theta^{(t)}) - g(\theta^{(0)})] - \eta J(\theta^{(0)})\mathbf{Q}^{(t)}g(\theta^{(0)}) \end{aligned} \quad (100)$$

297 which implies

$$\begin{aligned} \|a_{t+1}\|_2 \leq & \left\| (1 - \eta\mu)\mathbf{I} - \eta J(\theta^{(0)})\mathbf{Q}^{(t)}U(\tilde{\theta}^{(t)})J(\tilde{\theta}^{(t)})^\top \right\|_2 \|a_t\|_2 \\ & + \eta \left\| J(\theta^{(t)}) - J(\theta^{(0)}) \right\|_F \|g(\theta^{(t)})\|_2 + \eta \left\| J(\theta^{(0)}) \right\|_F \|g(\theta^{(0)})\|_2 \end{aligned} \quad (101)$$

298 where  $\tilde{\theta}^{(t)}$  is some linear interpolation between  $\theta^{(t)}$  and  $\theta^{(0)}$ . Our choice of  $\eta$  ensures that  $\eta\mu < 1$ .

299 Now we prove by induction that  $\|a_t\|_2 < C_0$ . It is true for  $t = 0$ , so we need to prove that if  
 300  $\|a_t\|_2 < C_0$ , then  $\|a_{t+1}\|_2 < C_0$ .

301 For the first term on the right-hand side of (101), we have

$$\begin{aligned} \left\| (1 - \eta\mu)\mathbf{I} - \eta J(\theta^{(0)})\mathbf{Q}^{(t)}U(\tilde{\theta}^{(t)})J(\tilde{\theta}^{(t)})^\top \right\|_2 \leq & (1 - \eta\mu) \left\| \mathbf{I} - \frac{\eta}{1 - \eta\mu} J(\theta^{(0)})\mathbf{Q}^{(t)}U(\tilde{\theta}^{(t)})J(\theta^{(0)})^\top \right\|_2 \\ & + \eta \left\| J(\theta^{(0)}) \right\|_F \left\| J(\tilde{\theta}^{(t)}) - J(\theta^{(0)}) \right\|_F \end{aligned} \quad (102)$$

302 Since  $\eta/(1 - \eta\mu) \leq (\lambda^{\min} + \lambda^{\max})^{-1}$  by our choice of  $\eta$ , we have

$$\left\| \mathbf{I} - \frac{\eta}{1 - \eta\mu} J(\theta^{(0)})\mathbf{Q}^{(t)}U(\tilde{\theta}^{(t)})J(\theta^{(0)})^\top \right\|_2 \leq 1 \quad (103)$$

303 On the other hand, we can use (52) since  $\|a_t\|_2 < C_0$ , so  $\|J(\theta^{(0)})\|_F \|J(\tilde{\theta}^{(t)}) - J(\theta^{(0)})\|_F \leq$   
 304  $\frac{M^2}{\sqrt[4]{d}}C_0$ . Therefore, there exists  $D_3$  such that for all  $\tilde{d} \geq D_3$ ,

$$\left\| (1 - \eta\mu)\mathbf{I} - \eta J(\theta^{(0)})\mathbf{Q}^{(t)}U(\tilde{\theta}^{(t)})J(\tilde{\theta}^{(t)})^\top \right\|_2 \leq 1 - \frac{\eta\mu}{2} \quad (104)$$

305 For the second term, we have

$$\begin{aligned} \|g(\theta^{(t)})\|_2 & \leq \|g(\theta^{(t)}) - g(\theta^{(0)})\|_2 + \|g(\theta^{(0)})\|_2 \\ & \leq \|J(\tilde{\theta}^{(t)})\|_2 \|U(\tilde{\theta}^{(t)})\|_2 \|\theta^{(t)} - \theta^{(0)}\|_2 + R_0 \leq MC_0 + R_0 \end{aligned} \quad (105)$$

306 And for the third term, we have

$$\eta \left\| J(\theta^{(0)}) \right\|_F \|g(\theta^{(0)})\|_2 \leq \eta MR_0 = A \quad (106)$$

307 Thus, we have

$$\|a_{t+1}\|_2 \leq \left(1 - \frac{\eta\mu}{2}\right) \|a_t\|_2 + \frac{\eta M(MC_0 + R_0)}{\sqrt[4]{d}} + A \quad (107)$$

308 So there exists  $D_4$  such that for all  $\tilde{d} \geq D_4$ ,  $\|a_{t+1}\|_2 \leq (1 - \frac{\eta\mu}{2}) \|a_t\|_2 + 2A$ . This shows that if  
 309  $\|a_t\|_2 < C_0$  is true, then  $\|a_{t+1}\|_2 < C_0$  will also be true.

310 In conclusion, for all  $\tilde{d} \geq D_0 = \max\{D_1, D_2, D_3, D_4\}$ ,  $\|\theta^{(t)} - \theta^{(0)}\|_2 < C_0$  is true for all  $t$ . This  
 311 also implies that for  $C_1 = MC_0 + R_0$ , we have  $\|g(\theta^{(t)})\|_2 \leq C_1$  for all  $t$  by (105). Similarly, we  
 312 can prove that  $\|\theta_{\text{lin}}^{(t)} - \theta^{(0)}\|_2 < C_0$  for all  $t$ .

313 Second, let  $\Delta_t = \theta_{\text{lin}}^{(t)} - \theta^{(t)}$ . Then we have

$$\Delta_{t+1} - \Delta_t = \eta(J(\theta^{(t)})\mathbf{Q}^{(t)}g(\theta^{(t)}) - J(\theta^{(0)})\mathbf{Q}^{(t)}g(\theta_{\text{lin}}^{(t)}) - \mu\Delta_t) \quad (108)$$

314 which implies

$$\Delta_{t+1} = \left[ (1 - \eta\mu)\mathbf{I} - \eta J(\theta^{(0)})\mathbf{Q}^{(t)}U(\tilde{\theta}^{(t)})J(\tilde{\theta}^{(t)})^\top \right] \Delta_t + \eta(J(\theta^{(t)}) - J(\theta^{(0)}))\mathbf{Q}^{(t)}g(\theta^{(t)}) \quad (109)$$

315 where  $\tilde{\theta}^{(t)}$  is some linear interpolation between  $\theta^{(t)}$  and  $\theta_{\text{lin}}^{(t)}$ . By (104), with probability at least  
 316  $(1 - \delta)$  for all  $\tilde{d} \geq D_0$ , we have

$$\begin{aligned} \|\Delta_{t+1}\|_2 &\leq \left\| (1 - \eta\mu)\mathbf{I} - \eta J(\theta^{(0)})\mathbf{Q}^{(t)}U(\tilde{\theta}^{(t)})J(\tilde{\theta}^{(t)})^\top \right\|_2 \|\Delta_t\|_2 + \eta \left\| J(\theta^{(t)}) - J(\theta^{(0)}) \right\|_F \|g(\theta^{(t)})\|_2 \\ &\leq \left( 1 - \frac{\eta\mu}{2} \right) \|\Delta_t\|_2 + \eta \frac{M}{\sqrt[4]{\tilde{d}}} C_0 C_1 \end{aligned} \quad (110)$$

317 Again, as  $\Delta_0 = 0$ , we can prove by induction that for all  $t$ ,

$$\|\Delta_t\|_2 < \frac{2MC_0C_1}{\mu} \tilde{d}^{-1/4} \quad (111)$$

318 For any test point  $\mathbf{x}$  such that  $\|\mathbf{x}\|_2 \leq 1$ , we have

$$\begin{aligned} \left| f_{\text{reg}}^{(t)}(\mathbf{x}) - f_{\text{linreg}}^{(t)}(\mathbf{x}) \right| &= \left| f(\mathbf{x}; \theta^{(t)}) - f_{\text{lin}}(\mathbf{x}; \theta_{\text{lin}}^{(t)}) \right| \\ &\leq \left| f(\mathbf{x}; \theta^{(t)}) - f_{\text{lin}}(\mathbf{x}; \theta^{(t)}) \right| + \left| f_{\text{lin}}(\mathbf{x}; \theta^{(t)}) - f_{\text{lin}}(\mathbf{x}; \theta_{\text{lin}}^{(t)}) \right| \\ &\leq \left| f(\mathbf{x}; \theta^{(t)}) - f_{\text{lin}}(\mathbf{x}; \theta^{(t)}) \right| + \left\| \nabla_{\theta} f(\mathbf{x}; \theta^{(0)}) \right\|_2 \left\| \theta^{(t)} - \theta_{\text{lin}}^{(t)} \right\|_2 \\ &\leq \left| f(\mathbf{x}; \theta^{(t)}) - f_{\text{lin}}(\mathbf{x}; \theta^{(t)}) \right| + M \|\Delta_t\|_2 \end{aligned} \quad (112)$$

319 For the first term, note that

$$\begin{cases} f(\mathbf{x}; \theta^{(t)}) - f(\mathbf{x}; \theta^{(0)}) = \nabla_{\theta} f(\mathbf{x}; \tilde{\theta}^{(t)}) (\theta^{(t)} - \theta^{(0)}) \\ f_{\text{lin}}(\mathbf{x}; \theta^{(t)}) - f_{\text{lin}}(\mathbf{x}; \theta^{(0)}) = \nabla_{\theta} f(\mathbf{x}; \theta^{(0)}) (\theta^{(t)} - \theta^{(0)}) \end{cases} \quad (113)$$

320 where  $\tilde{\theta}^{(t)}$  is some linear interpolation between  $\theta^{(t)}$  and  $\theta^{(0)}$ . Since  $f(\mathbf{x}; \theta^{(0)}) = f_{\text{lin}}(\mathbf{x}; \theta^{(0)})$ ,

$$\left| f(\mathbf{x}; \theta^{(t)}) - f_{\text{lin}}(\mathbf{x}; \theta^{(t)}) \right| \leq \left\| \nabla_{\theta} f(\mathbf{x}; \tilde{\theta}^{(t)}) - \nabla_{\theta} f(\mathbf{x}; \theta^{(0)}) \right\|_2 \left\| \theta^{(t)} - \theta^{(0)} \right\|_2 \leq \frac{M}{\sqrt[4]{\tilde{d}}} C_0^2 \quad (114)$$

321 Thus, we have shown that for all  $\tilde{d} \geq D_0$ , with probability at least  $(1 - \delta)$  for all  $t$  and all  $\mathbf{x}$ ,

$$\left| f_{\text{reg}}^{(t)}(\mathbf{x}) - f_{\text{linreg}}^{(t)}(\mathbf{x}) \right| \leq \left( MC_0^2 + \frac{2M^2C_0C_1}{\mu} \right) \tilde{d}^{-1/4} = O(\tilde{d}^{-1/4}) \quad (115)$$

322 which is the result we need.  $\square$

#### 323 D.4.2 Result for Linearized Neural Networks

324 **Lemma 18.** Suppose there exists  $M_0 > 0$  such that  $\|\nabla_{\theta} f^{(0)}(\mathbf{x})\|_2 \leq M_0$  for all test point  $\mathbf{x}$ . If the  
 325 gradients  $\nabla_{\theta} f^{(0)}(\mathbf{x}_1), \dots, \nabla_{\theta} f^{(0)}(\mathbf{x}_n)$  are linearly independent, and the empirical training risk of  
 326  $f_{\text{linreg}}^{(t)}$  satisfies

$$\limsup_{t \rightarrow \infty} \hat{\mathcal{R}}(f_{\text{linreg}}^{(t)}) < \epsilon, \quad (116)$$

327 for some  $\epsilon > 0$ , then for  $\mathbf{x}$  such that  $\|\mathbf{x}\|_2 \leq 1$  we have

$$\limsup_{t \rightarrow \infty} \left| f_{\text{linreg}}^{(t)}(\mathbf{x}) - f_{\text{linERM}}^{(t)}(\mathbf{x}) \right| = O(\sqrt{\epsilon}). \quad (117)$$



First, we can see that under the new weight update rule,  $\theta^{(t)} - \theta^{(0)} \in \text{span}\{\nabla_{\theta} f^{(0)}(\mathbf{x}_1), \dots, \nabla_{\theta} f^{(0)}(\mathbf{x}_n)\}$  is still true for all  $t$ . Let  $\theta^*$  be the interpolator in  $\text{span}(\nabla_{\theta} f^{(0)}(\mathbf{x}_1), \dots, \nabla_{\theta} f^{(0)}(\mathbf{x}_n))$ , then the empirical risk of  $\theta$  is  $\frac{1}{2n} \sum_{i=1}^n \langle \theta - \theta^*, \nabla_{\theta} f^{(0)}(\mathbf{x}_i) \rangle^2 = \frac{1}{2n} \|\nabla_{\theta} f^{(0)}(\mathbf{X})^{\top} (\theta - \theta^*)\|_2^2$ . Thus, there exists  $T > 0$  such that for any  $t \geq T$ ,

$$\left\| \nabla_{\theta} f^{(0)}(\mathbf{X})^{\top} (\theta^{(t)} - \theta^*) \right\|_2^2 \leq 2n\epsilon \quad (118)$$

Let the smallest singular value of  $\frac{1}{\sqrt{n}} \nabla_{\theta} f^{(0)}(\mathbf{X})$  be  $s^{\min}$ , and we have  $s^{\min} > 0$ . Note that the column space of  $\nabla_{\theta} f^{(0)}(\mathbf{X})$  is exactly  $\text{span}(\nabla_{\theta} f^{(0)}(\mathbf{x}_1), \dots, \nabla_{\theta} f^{(0)}(\mathbf{x}_n))$ . Define  $\mathbf{H} \in \mathbb{R}^{p \times n}$  such that its columns form an orthonormal basis of this subspace, then there exists  $\mathbf{G} \in \mathbb{R}^{n \times n}$  such that  $\nabla_{\theta} f^{(0)}(\mathbf{X}) = \mathbf{H}\mathbf{G}$ , and the smallest singular value of  $\frac{1}{\sqrt{n}}\mathbf{G}$  is also  $s^{\min}$ . Since  $\theta^{(t)} - \theta^{(0)}$  is also in this subspace, there exists  $\mathbf{v} \in \mathbb{R}^n$  such that  $\theta^{(t)} - \theta^* = \mathbf{H}\mathbf{v}$ . Then we have  $\sqrt{2n\epsilon} \geq \|\mathbf{G}^{\top} \mathbf{H}^{\top} \mathbf{H} \mathbf{v}\|_2 = \|\mathbf{G}^{\top} \mathbf{v}\|_2$ . Thus,  $\|\mathbf{v}\|_2 \leq \frac{\sqrt{2\epsilon}}{s^{\min}}$ , which implies

$$\left\| \theta^{(t)} - \theta^* \right\|_2 \leq \frac{\sqrt{2\epsilon}}{s^{\min}} \quad (119)$$

We have already proved in previous results that if we minimize the unregularized risk with ERM, then  $\theta$  always converges to the interpolator  $\theta^*$ . So for any  $t \geq T$  and any test point  $\mathbf{x}$  such that  $\|\mathbf{x}\|_2 \leq 1$ , we have

$$|f_{\text{linreg}}^{(t)}(\mathbf{x}) - f_{\text{linERM}}^{(t)}(\mathbf{x})| = |\langle \theta^{(t)} - \theta^*, \nabla_{\theta} f^{(0)}(\mathbf{x}) \rangle| \leq \frac{M_0 \sqrt{2\epsilon}}{s^{\min}} \quad (120)$$

which implies (117).  $\square$

#### D.4.3 Proof of Theorem 6

Given that  $\hat{\mathcal{R}}(f_{\text{linreg}}^{(t)}) < \epsilon$  for sufficiently large  $t$ , Lemma 16 implies that

$$\left| \hat{\mathcal{R}}(f_{\text{linreg}}^{(t)}) - \hat{\mathcal{R}}(f_{\text{reg}}^{(t)}) \right| = O(\tilde{d}^{-1/4} \sqrt{\epsilon} + \tilde{d}^{-1/2}) \quad (121)$$

So for a fixed  $\epsilon$ , there exists  $D > 0$  such that for all  $d \geq D$ , for sufficiently large  $t$ ,

$$\hat{\mathcal{R}}(f_{\text{reg}}^{(t)}) < \epsilon \Rightarrow \hat{\mathcal{R}}(f_{\text{linreg}}^{(t)}) < 2\epsilon \quad (122)$$

By Lemma 5 and Lemma 16, we have

$$\begin{cases} \sup_{t \geq 0} \left| f_{\text{linERM}}^{(t)}(\mathbf{x}) - f_{\text{ERM}}^{(t)}(\mathbf{x}) \right| = O(\tilde{d}^{-1/4}) \\ \sup_{t \geq 0} \left| f_{\text{linreg}}^{(t)}(\mathbf{x}) - f_{\text{reg}}^{(t)}(\mathbf{x}) \right| = O(\tilde{d}^{-1/4}) \end{cases} \quad (123)$$

Combining Lemma 18 with (123) derives

$$\limsup_{t \rightarrow \infty} \left| f_{\text{reg}}^{(t)}(\mathbf{x}) - f_{\text{ERM}}^{(t)}(\mathbf{x}) \right| = O(\tilde{d}^{-1/4} + \sqrt{\epsilon}) \quad (124)$$

Letting  $\tilde{d} \rightarrow \infty$  leads to the result we need.  $\square$

**Remark.** One might wonder whether  $\|\nabla_{\theta} f^{(0)}(\mathbf{x})\|_2$  will diverge as  $\tilde{d} \rightarrow \infty$ . In fact, in Lemma 13, we have proved that there exists a constant  $M$  such that with high probability, for any  $\tilde{d}$  there is  $\|\nabla_{\theta} f^{(0)}(\mathbf{x})\|_2 \leq M$  for any  $\mathbf{x}$  such that  $\|\mathbf{x}\|_2 \leq 1$ . Therefore, it is fine to suppose that there exists such an  $M_0$ .

## 353 D.5 Proofs for Subsection 5.1

### 354 D.5.1 Proof of Theorem 7

355 First we need to show that  $\hat{\theta}_{\text{MM}}$  is unique. Suppose both  $\theta_1$  and  $\theta_2$  maximize  $\min_{i=1,\dots,n} y_i \cdot$   
 356  $\langle \theta, \mathbf{x}_i \rangle$  and  $\theta_1 \neq \theta_2$ ,  $\|\theta_1\|_2 = \|\theta_2\|_2 = 1$ . Then consider  $\theta_0 = \theta / \|\theta\|_2$  where  $\theta = (\theta_1 + \theta_2)/2$ .  
 357 Obviously,  $\|\theta\|_2 < 1$ , and for any  $i$ ,  $y_i \cdot \langle \theta, \mathbf{x}_i \rangle = (y_i \cdot \langle \theta_1, \mathbf{x}_i \rangle + y_i \cdot \langle \theta_2, \mathbf{x}_i \rangle)/2$ , so  $y_i \cdot \langle \theta_0, \mathbf{x}_i \rangle >$   
 358  $\min\{y_i \cdot \langle \theta_1, \mathbf{x}_i \rangle, y_i \cdot \langle \theta_2, \mathbf{x}_i \rangle\}$ , which implies that  $\min_{i=1,\dots,n} y_i \cdot \langle \theta_0, \mathbf{x}_i \rangle > \min\{\min_{i=1,\dots,n} y_i \cdot$   
 359  $\langle \theta_1, \mathbf{x}_i \rangle, \min_{i=1,\dots,n} y_i \cdot \langle \theta_2, \mathbf{x}_i \rangle\}$ , contradiction!

360 Now we start proving the result. Without loss of generality, let  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$  be the samples  
 361 with the smallest margin to  $\mathbf{u}$ , i.e.

$$\arg \min_{1 \leq i \leq n} y_i \cdot \langle \mathbf{u}, \mathbf{x}_i \rangle = \{1, \dots, m\} \quad (125)$$

362 And denote  $y_1 \cdot \langle \mathbf{u}, \mathbf{x}_1 \rangle = \dots = y_m \cdot \langle \mathbf{u}, \mathbf{x}_m \rangle = \gamma_{\mathbf{u}}$ . Since the training error converges to 0,  $\gamma_{\mathbf{u}} > 0$ .  
 363 Note that for the logistic loss, if  $y_i \cdot \langle \theta, \mathbf{x}_i \rangle < y_j \cdot \langle \theta, \mathbf{x}_j \rangle$ , then for any  $M > 0$ , there exists an  
 364  $R_M > 0$  such that for all  $R \geq R_M$ ,

$$\frac{\nabla_{\theta} \ell(\langle R\theta, \mathbf{x}_i \rangle, y_i)}{\nabla_{\theta} \ell(\langle R\theta, \mathbf{x}_j \rangle, y_j)} > M \quad (126)$$

365 which can be shown with some simple calculation. And because the training error converges to 0,  
 366 we must have  $\|\theta^{(t)}\| \rightarrow \infty$ . Then, by Assumption 3 this means that when  $t$  gets sufficiently large,  
 367 the impact of  $(\mathbf{x}_j, y_j)$  to  $\theta^{(t)}$  where  $j > m$  is an infinitesimal compared to  $(\mathbf{x}_i, y_i)$  where  $i \leq m$   
 368 (because there exists a positive constant  $\delta$  such that  $q_i^{(t)} > \delta$  for all sufficiently large  $t$  by Assumption  
 369 3). Thus, we must have  $\mathbf{u} \in \text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ .

370 Let  $\mathbf{u} = \alpha_1 y_1 \mathbf{x}_1 + \dots + \alpha_m y_m \mathbf{x}_m$ . Now we show that  $\alpha_i \geq 0$  for all  $i = 1, \dots, m$ . This is  
 371 because when  $t$  is sufficiently large such that the impact of  $(\mathbf{x}_j, y_j)$  to  $\theta^{(t)}$  where  $j > m$  becomes  
 372 infinitesimal, we have

$$\theta^{(t+1)} - \theta^{(t)} \approx \eta \frac{q_i^{(t)} \exp(y_i \cdot \langle \theta^{(t)}, \mathbf{x}_i \rangle)}{1 + \exp(y_i \cdot \langle \theta^{(t)}, \mathbf{x}_i \rangle)} y_i \mathbf{x}_i \quad (127)$$

373 and since  $\|\theta^{(t)}\| \rightarrow \infty$  as  $t \rightarrow \infty$ , we have

$$\alpha_i \propto \lim_{T \rightarrow \infty} \sum_{t=T_0}^T \frac{q_i^{(t)} \exp(y_i \cdot \langle \theta^{(t)}, \mathbf{x}_i \rangle)}{1 + \exp(y_i \cdot \langle \theta^{(t)}, \mathbf{x}_i \rangle)} := \lim_{T \rightarrow \infty} \alpha_i(T) \quad (128)$$

374 where  $T_0$  is sufficiently large. Here the notion  $\alpha_i \propto \lim_{T \rightarrow \infty} \alpha_i(T)$  means that  $\lim_{T \rightarrow \infty} \frac{\alpha_i(T)}{\alpha_j(T)} = \frac{\alpha_i}{\alpha_j}$   
 375 for any pair of  $i, j$  and  $\alpha_j \neq 0$ . Note that each term in the sum is non-negative. This implies that all  
 376  $\alpha_1, \dots, \alpha_m$  have the same sign (or equal to 0). On the other hand,

$$\sum_{i=1}^m \alpha_i \gamma_{\mathbf{u}} = \sum_{i=1}^m \alpha_i y_i \cdot \langle \mathbf{u}, \mathbf{x}_i \rangle = \langle \mathbf{u}, \mathbf{u} \rangle > 0 \quad (129)$$

377 Thus,  $\alpha_i \geq 0$  for all  $i$  and at least one of them is positive. Now suppose  $\mathbf{u} \neq \hat{\theta}_{\text{MM}}$ , which means that  
 378  $\gamma_{\mathbf{u}}$  is smaller than the margin of  $\hat{\theta}_{\text{MM}}$ . Then, for all  $i = 1, \dots, m$ , there is  $y_i \cdot \langle \mathbf{u}, \mathbf{x}_i \rangle < y_i \cdot \langle \hat{\theta}_{\text{MM}}, \mathbf{x}_i \rangle$ .  
 379 This implies that

$$\langle \mathbf{u}, \mathbf{u} \rangle = \sum_{i=1}^m \alpha_i y_i \cdot \langle \mathbf{u}, \mathbf{x}_i \rangle < \sum_{i=1}^m \alpha_i y_i \cdot \langle \hat{\theta}_{\text{MM}}, \mathbf{x}_i \rangle = \langle \hat{\theta}_{\text{MM}}, \mathbf{u} \rangle \quad (130)$$

380 which is a contradiction. Thus, we must have  $\mathbf{u} = \hat{\theta}_{\text{MM}}$ . □

### 381 D.5.2 Proof of Theorem 8

382 Denote the largest and smallest eigenvalues of  $\mathbf{X}^\top \mathbf{X}$  by  $\lambda^{\max}$  and  $\lambda^{\min}$ , and by condition we have  
 383  $\lambda^{\min} > 0$ . Let  $\epsilon = \min\{\frac{q^*}{3}, \frac{(q^* \lambda^{\min})^2}{192 \lambda^{\max 2}}\}$ . Then similar to the proof in Appendix D.2.2, there exists  $t_\epsilon$

such that for all  $t \geq t_\epsilon$  and all  $i, q_i^{(t)} \in (q_i - \epsilon, q_i + \epsilon)$ . Denote  $\mathbf{Q} = \text{diag}(q_1, \dots, q_n)$ , then for all  $t \geq t_\epsilon$ ,  $\mathbf{Q}^{(t)} := \mathbf{Q}_\epsilon^{(t)} = \sqrt{\mathbf{Q}} \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}}$ , where we use the subscript  $\epsilon$  to indicate that  $\|\mathbf{Q}_\epsilon^{(t)} - \mathbf{Q}\|_2 < \epsilon$ . First, we prove that  $F(\theta)$  is  $L$ -smooth as long as  $\|\mathbf{x}_i\|_2 \leq 1$  for all  $i$ . The gradient of  $F$  is

$$\nabla F(\theta) = \sum_{i=1}^n q_i \nabla_{\hat{y}} \ell(\langle \theta, \mathbf{x}_i \rangle, y_i) \mathbf{x}_i \quad (131)$$

Since  $\ell(\hat{y}, y)$  is  $L$ -smooth in  $\hat{y}$ , we have for any  $\theta_1, \theta_2$  and any  $i$ ,

$$\begin{aligned} & \ell(\langle \theta_2, \mathbf{x}_i \rangle, y_i) - \ell(\langle \theta_1, \mathbf{x}_i \rangle, y_i) \\ & \leq \nabla_{\hat{y}} \ell(\langle \theta_1, \mathbf{x}_i \rangle, y_i) \cdot (\langle \theta_2, \mathbf{x}_i \rangle - \langle \theta_1, \mathbf{x}_i \rangle) + \frac{L}{2} (\langle \theta_2, \mathbf{x}_i \rangle - \langle \theta_1, \mathbf{x}_i \rangle)^2 \\ & = \langle \nabla_{\hat{y}} \ell(\langle \theta_1, \mathbf{x}_i \rangle, y_i) \cdot \mathbf{x}_i, \theta_2 - \theta_1 \rangle + \frac{L}{2} (\langle \theta_2 - \theta_1, \mathbf{x}_i \rangle)^2 \\ & \leq \langle \nabla_{\hat{y}} \ell(\langle \theta_1, \mathbf{x}_i \rangle, y_i) \cdot \mathbf{x}_i, \theta_2 - \theta_1 \rangle + \frac{L}{2} \|\theta_2 - \theta_1\|_2^2 \end{aligned} \quad (132)$$

Thus, we have

$$\begin{aligned} F(\theta_2) - F(\theta_1) &= \sum_{i=1}^n q_i [\ell(\langle \theta_2, \mathbf{x}_i \rangle, y_i) - \ell(\langle \theta_1, \mathbf{x}_i \rangle, y_i)] \\ &\leq \sum_{i=1}^n q_i \langle \nabla_{\hat{y}} \ell(\langle \theta_1, \mathbf{x}_i \rangle, y_i) \cdot \mathbf{x}_i, \theta_2 - \theta_1 \rangle + \frac{L}{2} \sum_{i=1}^n q_i \|\theta_2 - \theta_1\|_2^2 \\ &= \langle \nabla F(\theta_1), \theta_2 - \theta_1 \rangle + \frac{L}{2} \|\theta_2 - \theta_1\|_2^2 \end{aligned} \quad (133)$$

which implies that  $F(\theta)$  is  $L$ -smooth.

Denote  $\tilde{g}(\theta) = \nabla_{\hat{y}} \ell(f(\mathbf{X}; \theta), \mathbf{Y}) \in \mathbb{R}^n$ , then  $\nabla F(\theta^{(t)}) = \mathbf{X} \mathbf{Q} \tilde{g}(\theta^{(t)})$ , and the update rule is

$$\theta^{(t+1)} = \theta^{(t)} - \eta \mathbf{X} \mathbf{Q}^{(t)} \tilde{g}(\theta^{(t)}) \quad (134)$$

So by the upper quadratic bound, we have

$$F(\theta^{(t+1)}) \leq F(\theta^{(t)}) - \eta \langle \mathbf{X} \mathbf{Q} \tilde{g}(\theta^{(t)}), \mathbf{X} \mathbf{Q}^{(t)} \tilde{g}(\theta^{(t)}) \rangle + \frac{\eta^2 L}{2} \|\mathbf{X} \mathbf{Q}^{(t)} \tilde{g}(\theta^{(t)})\|_2^2 \quad (135)$$

Let  $\eta_1 = \frac{q^* \lambda^{\min}}{2L(1+3\epsilon)\lambda^{\max}}$ . Similar to what we did in Appendix D.2.2 (Eqn. (40)), we can prove that for all  $\eta \leq \eta_1$ , (135) implies that for all  $t \geq t_\epsilon$ , there is

$$\begin{aligned} F(\theta^{(t+1)}) &\leq F(\theta^{(t)}) - \frac{\eta q^* \lambda^{\min}}{2} \|\sqrt{\mathbf{Q}} \tilde{g}(\theta^{(t)})\|_2^2 + \frac{\eta^2 L}{2} \|\mathbf{X} \sqrt{\mathbf{Q}_{3\epsilon}^{(t)}}\|_2^2 \|\sqrt{\mathbf{Q}} \tilde{g}(\theta^{(t)})\|_2^2 \\ &\leq F(\theta^{(t)}) - \frac{\eta q^* \lambda^{\min}}{2} \|\sqrt{\mathbf{Q}} \tilde{g}(\theta^{(t)})\|_2^2 + \frac{\eta^2 L}{2} \|\mathbf{X}\|_2^2 (1+3\epsilon) \|\sqrt{\mathbf{Q}} \tilde{g}(\theta^{(t)})\|_2^2 \\ &\leq F(\theta^{(t)}) - \frac{\eta q^* \lambda^{\min}}{4} \|\sqrt{\mathbf{Q}} \tilde{g}(\theta^{(t)})\|_2^2 \\ &\leq F(\theta^{(t)}) - \frac{\eta q^{*2} \lambda^{\min}}{4} \|\tilde{g}(\theta^{(t)})\|_2^2 \end{aligned} \quad (136)$$

This shows that  $F(\theta^{(t)})$  is monotonically non-increasing. Since  $F(\theta) \geq 0$ ,  $F(\theta^{(t)})$  must converge as  $t \rightarrow \infty$ , and we need to prove that it converges to 0. Suppose that  $F(\theta^{(t)})$  does not converge to 0, then there exists a constant  $C > 0$  such that  $F(\theta^{(t)}) \geq 2C$  for all  $t$ . On the other hand, it is easy to see that there exists  $\theta^*$  such that  $\ell(\langle \theta^*, \mathbf{x}_i \rangle, y_i) < C$  for all  $i$ . (136) also implies that  $\|\tilde{g}(\theta^{(t)})\|_2 \rightarrow 0$  as  $t \rightarrow \infty$  because we must have  $F(\theta^{(t)}) - F(\theta^{(t+1)}) \rightarrow 0$ .

Note that from (134) we have

$$\|\theta^{(t+1)} - \theta^*\|_2^2 = \|\theta^{(t)} - \theta^*\|_2^2 + 2\eta \langle \mathbf{X} \mathbf{Q}^{(t)} \tilde{g}(\theta^{(t)}), \theta^* - \theta^{(t)} \rangle + \eta^2 \|\mathbf{X} \mathbf{Q}^{(t)} \tilde{g}(\theta^{(t)})\|_2^2 \quad (137)$$

400 Denote

$$F_t(\theta) = \sum_{i=1}^n q_i^{(t)} \ell(\langle \theta, \mathbf{x}_i \rangle, y_i) \quad (138)$$

401 Then  $F_t$  is convex because  $\ell$  is convex and  $q_i^{(t)}$  are non-negative, and  $\nabla F_t(\theta^{(t)}) = \mathbf{X} \mathbf{Q}^{(t)} \tilde{g}(\theta^{(t)})$ .  
 402 By the lower linear bound  $F_t(\mathbf{y}) \geq F_t(\mathbf{x}) + \langle \nabla F_t(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$ , we have for all  $t$ ,

$$\langle \mathbf{X} \mathbf{Q}^{(t)} \tilde{g}(\theta^{(t)}), \theta^* - \theta^{(t)} \rangle \leq F_t(\theta^*) - F_t(\theta^{(t)}) \leq F_t(\theta^*) - \frac{2}{3} F(\theta^{(t)}) \leq C - \frac{4C}{3} = -\frac{C}{3} \quad (139)$$

403 because  $q_i^{(t)} \geq q_i - \epsilon \geq \frac{2}{3} q_i$  and  $\sum_{i=1}^n q_i^{(t)} = 1$ . Since  $\|\tilde{g}(\theta^{(t)})\|_2 \rightarrow 0$ , there exists  $T > 0$  such that  
 404 for all  $t \geq T$  and all  $\eta \leq \eta_0$ ,

$$\|\theta^{(t+1)} - \theta^*\|_2^2 \leq \|\theta^{(t)} - \theta^*\|_2^2 - \frac{\eta C}{3} \quad (140)$$

405 which means that  $\|\theta^{(t)} - \theta^*\|_2^2 \rightarrow -\infty$  because  $\frac{\eta C}{3}$  is a positive constant. This is a contradiction!  
 406 Thus,  $F(\theta^{(t)})$  must converge to 0, which is result (i).

407 (i) immediately implies (ii) because  $\ell$  is strictly decreasing to 0 by condition.

408 Now let's prove (iii). First of all, the uniqueness of  $\theta_R$  can be easily proved from the convexity  
 409 of  $F(\theta)$ . The condition implies that  $y_i \langle \theta_R, \mathbf{x}_i \rangle > 0$ , i.e.  $\theta_R$  must classify all training samples  
 410 correctly. If there are two different minimizers  $\theta_R$  and  $\theta'_R$  in whose norm is at most  $R$ , then consider  
 411  $\theta''_R = \frac{1}{2}(\theta_R + \theta'_R)$ . By the convexity of  $F$ , we know that  $\theta''_R$  must also be a minimizer, and  $\|\theta''_R\|_2 < R$ .  
 412 Thus,  $F(\frac{R}{\|\theta''_R\|_2} \theta''_R) < F(\theta''_R)$  and  $\|\frac{R}{\|\theta''_R\|_2} \theta''_R\|_2 = R$ , which contradicts with the fact that  $\theta''_R$  is a  
 413 minimizer.

414 To prove the rest of (iii), the key is to consider (135). On one hand, similar to (36) we can prove that  
 415 for all  $t \geq t_\epsilon$ , there is

$$\left| \langle \mathbf{X} \mathbf{Q}^{(t)} \tilde{g}(\theta^{(t)}), \mathbf{X}(\mathbf{Q}^{(t)} - \mathbf{Q}) \tilde{g}(\theta^{(t)}) \rangle \right| \leq \lambda^{\max} \sqrt{3\epsilon} \left\| \sqrt{\mathbf{Q}^{(t)}} \tilde{g}(\theta^{(t)}) \right\|_2^2 \quad (141)$$

416 Since we choose  $\epsilon = \min\{\frac{q^*}{3}, \frac{(q^* \lambda^{\min})^2}{192 \lambda^{\max} 2}\}$ , this inequality implies that

$$\begin{aligned} \left\| \nabla F_t(\theta^{(t)}) \right\|_2^2 &= \left\| \mathbf{X} \mathbf{Q}^{(t)} \tilde{g}(\theta^{(t)}) \right\|_2^2 \geq \lambda^{\min} \left\| \mathbf{Q}^{(t)} \tilde{g}(\theta^{(t)}) \right\|_2^2 \geq \lambda^{\min} (q^* - \epsilon) \left\| \sqrt{\mathbf{Q}^{(t)}} \tilde{g}(\theta^{(t)}) \right\|_2^2 \\ &\geq \frac{\lambda^{\min} q^*}{2} \left\| \sqrt{\mathbf{Q}^{(t)}} \tilde{g}(\theta^{(t)}) \right\|_2^2 \geq 4 \left| \langle \mathbf{X} \mathbf{Q}^{(t)} \tilde{g}(\theta^{(t)}), \mathbf{X}(\mathbf{Q}^{(t)} - \mathbf{Q}) \tilde{g}(\theta^{(t)}) \rangle \right| \end{aligned} \quad (142)$$

417 On the other hand, if  $\eta \leq \eta_2 = \frac{1}{2L}$ , we will have

$$\frac{\eta^2 L}{2} \left\| \mathbf{X} \mathbf{Q}^{(t)} \tilde{g}(\theta^{(t)}) \right\|_2^2 \leq \frac{\eta}{4} \left\| \nabla F_t(\theta^{(t)}) \right\|_2^2 \quad (143)$$

418 Combining all the above with (135) yields

$$F(\theta^{(t+1)}) - F(\theta^{(t)}) \leq -\frac{\eta}{2} \left\| \nabla F_t(\theta^{(t)}) \right\|_2^2 \quad (144)$$

419 Denote  $\mathbf{u} = \lim_{R \rightarrow \infty} \frac{\theta_R}{\|\theta_R\|_2}$ . Similar to Lemma 9 in [JDST20], we can prove that: for any  $\alpha > 0$ ,  
 420 there exists a constant  $\rho(\alpha) > 0$  such that for any  $\theta$  subject to  $\|\theta\|_2 \geq \rho(\alpha)$ , there is

$$F_t((1 + \alpha)\|\theta\|_2 \mathbf{u}) \leq F_t(\theta) \quad (145)$$

421 for any  $t$ . Let  $t_\alpha \geq t_\epsilon$  satisfy that for all  $t \geq t_\alpha$ ,  $\|\theta^{(t)}\|_2 \geq \max\{\rho(\alpha), 1\}$ . By the convexity of  $F_t$ ,  
 422 for all  $t \geq t_\alpha$ ,

$$\langle \nabla F_t(\theta^{(t)}), \theta^{(t)} - (1 + \alpha)\|\theta^{(t)}\|_2 \mathbf{u} \rangle \geq F_t(\theta^{(t)}) - F_t((1 + \alpha)\|\theta^{(t)}\|_2 \mathbf{u}) \geq 0 \quad (146)$$

Thus, we have

$$\begin{aligned}
\langle \theta^{(t+1)} - \theta^{(t)}, \mathbf{u} \rangle &= \langle -\eta \nabla F_t(\theta^{(t)}), \mathbf{u} \rangle \\
&\geq \langle -\eta \nabla F_t(\theta^{(t)}), \theta^{(t)} \rangle \frac{1}{(1+\alpha)\|\theta^{(t)}\|_2} \\
&= \langle \theta^{(t+1)} - \theta^{(t)}, \theta^{(t)} \rangle \frac{1}{(1+\alpha)\|\theta^{(t)}\|_2} \\
&= \left( \frac{1}{2} \|\theta^{(t+1)}\|_2^2 - \frac{1}{2} \|\theta^{(t)}\|_2^2 - \frac{1}{2} \|\theta^{(t+1)} - \theta^{(t)}\|_2^2 \right) \frac{1}{(1+\alpha)\|\theta^{(t)}\|_2}
\end{aligned} \tag{147}$$

By  $\frac{1}{2}(\|\theta^{(t+1)}\|_2 - \|\theta^{(t)}\|_2)^2 \geq 0$ , we have  $(\frac{1}{2} \|\theta^{(t+1)}\|_2^2 - \frac{1}{2} \|\theta^{(t)}\|_2^2) / \|\theta^{(t)}\|_2 \geq \|\theta^{(t+1)}\|_2 - \|\theta^{(t)}\|_2$ .  
Moreover, by (144) we have

$$\frac{\|\theta^{(t+1)} - \theta^{(t)}\|_2^2}{2(1+\alpha)\|\theta^{(t)}\|_2} \leq \frac{\|\theta^{(t+1)} - \theta^{(t)}\|_2^2}{2} = \frac{\eta^2 \|\nabla F_t(\theta^{(t)})\|_2^2}{2} \leq \eta (F(\theta^{(t)}) - F(\theta^{(t+1)})) \tag{148}$$

Summing up (147) from  $t = t_\alpha$  to  $t - 1$ , we have

$$\langle \theta^{(t)} - \theta^{(t_\alpha)}, \mathbf{u} \rangle \geq \frac{\|\theta^{(t)}\|_2 - \|\theta^{(t_\alpha)}\|_2}{1+\alpha} + \eta (F(\theta^{(t)}) - F(\theta^{(t_\alpha)})) \geq \frac{\|\theta^{(t)}\|_2 - \|\theta^{(t_\alpha)}\|_2}{1+\alpha} - \eta F(\theta^{(t_\alpha)}) \tag{149}$$

which implies that

$$\left\langle \frac{\theta^{(t)}}{\|\theta^{(t)}\|_2}, \mathbf{u} \right\rangle \geq \frac{1}{1+\alpha} + \frac{1}{\|\theta^{(t)}\|_2} \left( \langle \theta^{(t_\alpha)}, \mathbf{u} \rangle - \frac{\|\theta^{(t_\alpha)}\|_2}{1+\alpha} - \eta F(\theta^{(t_\alpha)}) \right) \tag{150}$$

Since  $\lim_{t \rightarrow \infty} \|\theta^{(t)}\|_2 = \infty$ , we have

$$\liminf_{t \rightarrow \infty} \left\langle \frac{\theta^{(t)}}{\|\theta^{(t)}\|_2}, \mathbf{u} \right\rangle \geq \frac{1}{1+\alpha} \tag{151}$$

Since  $\alpha$  is arbitrary, we must have  $\lim_{t \rightarrow \infty} \frac{\theta^{(t)}}{\|\theta^{(t)}\|_2} = \mathbf{u}$  as long as  $\eta \leq \min\{\eta_1, \eta_2\}$ .  $\square$

### D.5.3 Corollary of Theorem 8

We can show that for the logistic loss, it satisfies all conditions of Theorem 8 and  $\lim_{R \rightarrow \infty} \frac{\theta_R}{R} = \hat{\theta}_{\text{MM}}$ .

First of all, for the logistic loss we have  $\nabla_y^2 \ell(\hat{y}, y) = \frac{y^2}{e^{y\hat{y}} + e^{-y\hat{y}} + 2} \leq \max_i \frac{y_i^2}{4}$ , so  $\ell$  is smooth.

Then, we prove that  $\lim_{R \rightarrow \infty} \frac{\theta_R}{R}$  exists and is equal to  $\hat{\theta}_{\text{MM}}$ . For the logistic loss, it is easy to show that for any  $\hat{\theta}' \neq \hat{\theta}_{\text{MM}}$ , there exists an  $R(\hat{\theta}') > 0$  and an  $\delta(\hat{\theta}') > 0$  such that  $F(R \cdot \theta) > F(R \cdot \hat{\theta}_{\text{MM}})$  for all  $R \geq R(\hat{\theta}')$  and  $\theta \in B(\hat{\theta}', \delta(\hat{\theta}'))$ .

Let  $S = \{\theta : \|\theta\|_2 = 1\}$ . For any  $\epsilon > 0$ ,  $S - B(\hat{\theta}_{\text{MM}}, \epsilon)$  is a compact set. And for any  $\theta \in S - B(\hat{\theta}_{\text{MM}}, \epsilon)$ , there exist  $R(\theta)$  and  $\delta(\theta)$  as defined above. Thus, there must exist  $\theta_1, \dots, \theta_m \in S - B(\hat{\theta}_{\text{MM}}, \epsilon)$  such that  $S - B(\hat{\theta}_{\text{MM}}, \epsilon) \subseteq \cup_{i=1}^m B(\theta_i, \delta(\theta_i))$ . Let  $R(\epsilon) = \max\{R(\theta_1), \dots, R(\theta_m)\}$ , then for all  $R \geq R(\epsilon)$  and all  $\theta \in S - B(\hat{\theta}_{\text{MM}}, \epsilon)$ ,  $F(R \cdot \theta) > F(R \cdot \hat{\theta}_{\text{MM}})$ , which means that  $\frac{\theta_R}{R} \in B(\hat{\theta}_{\text{MM}}, \epsilon)$  for all  $R \geq R(\epsilon)$ . Therefore,  $\lim_{R \rightarrow \infty} \frac{\theta_R}{R}$  exists and is equal to  $\hat{\theta}_{\text{MM}}$ .

Therefore, by Theorem 8, any GRW satisfying Assumption 1 makes a linear model converge to the max-margin classifier under the logistic loss.

### D.6 Proof of Theorem 9

We first consider the regularized linearized neural network  $f_{\text{linreg}}^{(t)}$ . Since by Proposition 11  $f^{(0)}(\mathbf{x})$  is sampled from a zero-mean Gaussian process, there exists a constant  $M > 0$  such that  $|f^{(0)}(\mathbf{x}_i)| < M$

446 for all  $i$  with high probability. Define

$$F(\theta) = \sum_{i=1}^n q_i \ell(\langle \theta, \nabla_{\theta} f^{(0)}(\mathbf{x}_i) \rangle + f^{(0)}(\mathbf{x}_i), y_i) \quad (152)$$

447 Denote  $\tilde{\theta}_R = \arg \min_{\theta} \{F(R \cdot \theta) : \|\theta\|_2 \leq 1\}$ . when the linearized neural network is trained  
 448 by a GRW satisfying Assumption 1 with regularization, since this is convex optimization and the  
 449 objective function is smooth, we can prove that with a sufficiently small learning rate, as  $t \rightarrow \infty$ ,  
 450  $\theta^{(t)} \rightarrow R \cdot \tilde{\theta}_R + \theta^{(0)}$  where  $R = \lim_{t \rightarrow \infty} \|\theta^{(t)} - \theta^{(0)}\|_2$  (which is the minimizer). And define

$$\gamma = \min_{i=1, \dots, n} y_i \cdot \langle \hat{\theta}_{\text{MM}}, \nabla_{\theta} f^{(0)}(\mathbf{x}_i) \rangle \quad (153)$$

451 First, we derive the lower bound of  $R$ . By Theorem 16, with a sufficiently large  $\tilde{d}$ , with high  
 452 probability  $\hat{\mathcal{R}}(f_{\text{reg}}^{(t)}) < \epsilon$  implies  $\hat{\mathcal{R}}(f_{\text{linreg}}^{(t)}) < 2\epsilon$ . By the convexity of  $\ell$ , we have

$$\begin{aligned} 2\epsilon &> \frac{1}{n} \sum_{i=1}^n \ell(\langle R\tilde{\theta}_R, \mathbf{x}_i \rangle + f^{(0)}(\mathbf{x}_i), y_i) \geq \log \left( 1 + \exp \left( -\frac{1}{n} \sum_{i=1}^n (\langle R\tilde{\theta}_R, \mathbf{x}_i \rangle + f^{(0)}(\mathbf{x}_i)) y_i \right) \right) \\ &\geq \log \left( 1 + \exp \left( -\frac{1}{n} \sum_{i=1}^n R \langle \tilde{\theta}_R, \mathbf{x}_i \rangle y_i - M \right) \right) \end{aligned} \quad (154)$$

453 which implies that  $R = \Omega(-\log 2\epsilon)$  for all  $\epsilon \in (0, \frac{1}{4})$ .

454 Denote  $\delta = \|\hat{\theta}_{\text{MM}} - \tilde{\theta}_R\|_2$ . Let  $\theta' = \frac{\hat{\theta}_{\text{MM}} + \tilde{\theta}_R}{2}$ , then we can see that  $\|\theta'\|_2 = \sqrt{1 - \frac{\delta^2}{4}}$ . Let  $\tilde{\theta}' = \frac{\theta'}{\|\theta'\|_2}$ .  
 455 By the definition of  $\hat{\theta}_{\text{MM}}$ , there exists  $j$  such that  $y_j \cdot \langle \tilde{\theta}', \nabla_{\theta} f^{(0)}(\mathbf{x}_j) \rangle \leq \gamma$ , which implies

$$y_j \cdot \left\langle \frac{\hat{\theta}_{\text{MM}} + \tilde{\theta}_R}{2} \frac{1}{\sqrt{1 - \frac{\delta^2}{4}}}, \nabla_{\theta} f^{(0)}(\mathbf{x}_j) \right\rangle \leq \gamma \quad (155)$$

456 Thus, we have

$$\begin{aligned} y_j \cdot \langle \tilde{\theta}_R, \nabla_{\theta} f^{(0)}(\mathbf{x}_j) \rangle &\leq 2\sqrt{1 - \frac{\delta^2}{4}} \gamma - y_j \cdot \langle \hat{\theta}_{\text{MM}}, \nabla_{\theta} f^{(0)}(\mathbf{x}_j) \rangle \\ &\leq \left( 2\sqrt{1 - \frac{\delta^2}{4}} - 1 \right) \gamma \\ &\leq \left( 2\left(1 - \frac{\delta^2}{8}\right) - 1 \right) \gamma \quad (\text{since } \sqrt{1-x} \leq 1 - \frac{x}{2}) \\ &= \left(1 - \frac{\delta^2}{4}\right) \gamma \end{aligned} \quad (156)$$

457 On the other hand, we have

$$\begin{aligned} q_j \log(1 + \exp(-y_j \cdot \langle R \cdot \tilde{\theta}_R, \nabla_{\theta} f^{(0)}(\mathbf{x}_j) \rangle - M)) &\leq F(R \cdot \tilde{\theta}_R) \\ &\leq F(R \cdot \hat{\theta}_{\text{MM}}) \leq \log(1 + \exp(-R\gamma + M)) \end{aligned} \quad (157)$$

458 which implies that

$$q^* \log \left( 1 + \exp \left( -\left(1 - \frac{\delta^2}{4}\right) R\gamma - M \right) \right) \leq \log(1 + \exp(-R\gamma + M)) \quad (158)$$

459 and this leads to

$$1 + \exp(-R\gamma + M) \geq \left( 1 + \exp \left( -\left(1 - \frac{\delta^2}{4}\right) R\gamma - M \right) \right)^{q^*} \geq 1 + q^* \exp \left( -\left(1 - \frac{\delta^2}{4}\right) R\gamma - M \right) \quad (159)$$

460 which is equivalent to

$$-R\gamma + M \geq -(1 - \frac{\delta^2}{4})R\gamma - M + \log(q^*) \quad (160)$$

461 Thus, we have

$$\delta = O(R^{-1/2}) = O((- \log 2\epsilon)^{-1/2}) \quad (161)$$

462 So for any test point  $\mathbf{x}$ , since  $\|\nabla_{\theta} f^{(0)}(\mathbf{x})\|_2 \leq M_0$ , we have

$$|\langle \hat{\theta}_{\text{MM}} - \tilde{\theta}_R, \nabla_{\theta} f^{(0)}(\mathbf{x}) \rangle| \leq \delta M_0 = O((- \log 2\epsilon)^{-1/2}) \quad (162)$$

463 Combined with Theorem 16, we have: with high probability,

$$\limsup_{t \rightarrow \infty} |R \cdot f_{\text{MM}}(\mathbf{x}) - f_{\text{reg}}^{(t)}(\mathbf{x})| = O(R \cdot (- \log 2\epsilon)^{-1/2} + \tilde{d}^{-1/4}) \quad (163)$$

464 So there exists a constant  $C > 0$  such that: As  $\tilde{d} \rightarrow \infty$ , with high probability, for all  $\epsilon \in (0, \frac{1}{4})$ , if  
 465  $|f_{\text{MM}}(\mathbf{x})| > C \cdot (- \log 2\epsilon)^{-1/2}$ , then  $f_{\text{reg}}^{(t)}(\mathbf{x})$  will have the same sign as  $f_{\text{MM}}(\mathbf{x})$  for a sufficiently  
 466 large  $t$ . Note that this  $C$  only depends on  $n, q^*, \gamma, M$  and  $M_0$ , so it is a constant independent of  
 467  $\epsilon$ .  $\square$

468 **Remark.** Note that Theorem 9 requires Assumption 1 while Theorem 6 does not due to the  
 469 fundamental difference between the classification and regression. In regression the model converges  
 470 to a finite point. However, in classification, the training loss converging to zero implies that either (i)  
 471 The direction of the weight is close to the max-margin classifier or (ii) The norm of the weight is  
 472 very large. Assumption 1 is used to eliminate the possibility of (ii). If the regularization parameter  $\mu$   
 473 is sufficiently large, then a small empirical risk could imply a small weight norm. However, in our  
 474 theorem we do not assume anything on  $\mu$ , so Assumption 1 is necessary.

## 475 E A Note on the Proofs in [LXS<sup>+</sup>19]

476 We have mentioned that the proofs in [LXS<sup>+</sup>19], particularly the proofs of their Theorem 2.1 and  
 477 Lemma 1 in their Appendix G, are flawed. In order to fix their proof, we change the network  
 478 initialization to (9). In this section, we will demonstrate what goes wrong in the proofs in [LXS<sup>+</sup>19],  
 479 and how we manage to fix the proof. For clarity, we are referring to the following version of the  
 480 paper: <https://arxiv.org/pdf/1902.06720v4.pdf>.

481 To avoid confusion, in this section we will still use the notations used in our paper.

### 482 E.1 Their Problems

483 [LXS<sup>+</sup>19] claimed in their Theorem 2.1 that under the conditions of our Lemma 5, for any  $\delta > 0$ ,  
 484 there exist  $\tilde{D} > 0$  and a constant  $C$  such that for any  $\tilde{d} \geq \tilde{D}$ , with probability at least  $(1 - \delta)$ , the  
 485 gap between the output of a sufficiently wide fully-connected neural network and the output of its  
 486 linearized neural network at any test point  $\mathbf{x}$  can be uniformly bounded by

$$\sup_{t \geq 0} |f^{(t)}(\mathbf{x}) - f_{\text{lin}}^{(t)}(\mathbf{x})| \leq C\tilde{d}^{-1/2} \quad (\text{claimed}) \quad (164)$$

487 where they used the original NTK formulation and initialization in [JGH18]:

$$\begin{cases} \mathbf{h}^{l+1} = \frac{W^l}{\sqrt{d_l}} \mathbf{x}^l + \beta \mathbf{b}^l \\ \mathbf{x}^{l+1} = \sigma(\mathbf{h}^{l+1}) \end{cases} \quad \text{and} \quad \begin{cases} W_{i,j}^{l(0)} \sim \mathcal{N}(0, 1) \\ b_i^{l(0)} \sim \mathcal{N}(0, 1) \end{cases} \quad (\forall l = 0, \dots, L) \quad (165)$$

488 where  $\mathbf{x}_0 = \mathbf{x}$  and  $f(\mathbf{x}) = h^{L+1}$ . However, in their proof in their Appendix G, they did not directly  
 489 prove their result for the NTK formulation, but instead they proved another result for the following  
 490 formulation which they called the *standard formulation*:

$$\begin{cases} \mathbf{h}^{l+1} = W^l \mathbf{x}^l + \beta \mathbf{b}^l \\ \mathbf{x}^{l+1} = \sigma(\mathbf{h}^{l+1}) \end{cases} \quad \text{and} \quad \begin{cases} W_{i,j}^{l(0)} \sim \mathcal{N}(0, \frac{1}{d_l}) \\ b_i^{l(0)} \sim \mathcal{N}(0, 1) \end{cases} \quad (\forall l = 0, \dots, L) \quad (166)$$

See their Appendix F for the definition of their standard formulation. In the original formulation, they also included two constants  $\sigma_w$  and  $\sigma_b$  for standard deviations, and for simplicity we omit these constants here. Note that the outputs of the NTK formulation and the standard formulation at initialization are actually the same. The only difference is that the norm of the weight  $W^l$  and the gradient of the model output with respect to  $W^l$  are different for all  $l$ .

In their Appendix G, they claimed that if a network with the standard formulation is trained by minimizing the squared loss with gradient descent and learning rate  $\eta' = \eta/\tilde{d}$ , where  $\eta$  is our learning rate in Lemma 5 and also their learning rate in their Theorem 2.1, then (164) is true for this network, so it is also true for a network with the NTK formulation because the two formulations have the same network output. And then they claimed in their equation (S37) that applying learning rate  $\eta'$  to the standard formulation is equivalent to applying the following learning rates

$$\eta_W^l = \frac{d_l}{d_{\max}}\eta \quad \text{and} \quad \eta_b^l = \frac{1}{d_{\max}}\eta \quad (167)$$

to  $W^l$  and  $b^l$  of the NTK formulation, where  $d_{\max} = \max\{d_0, \dots, d_L\}$ .

To avoid confusion, in the following discussions we will still use the NTK formulation and initialization if not stated otherwise.

**Problem 1.** Claim (167) is true, but it leads to two problems. The first problem is that  $\eta_b^l = O(d_{\max}^{-1})$  since  $\eta = O(1)$ , while their Theorem 2.1 needs the learning rate to be  $O(1)$ . Nevertheless, this problem can be simply fixed by modifying their standard formulation as  $\mathbf{h}^{l+1} = W^l \mathbf{x}^l + \beta \sqrt{d_l} \mathbf{b}^l$  where  $b_i^{l(0)} \sim \mathcal{N}(0, d_l^{-1})$ . The real problem that is non-trivial to fix is that by (167), there is  $\eta_W^0 = \frac{d_0}{d_{\max}}\eta$ . However, note that  $d_0$  is a constant since it is the dimension of the input space, while  $d_{\max}$  goes to infinity. Consequently, in (167) they were essentially using a very small learning rate for the first layer  $W^0$  but a normal learning rate for the rest of the layers, which definitely does not match with their claim in their Theorem 2.1.

**Problem 2.** Another big problem is that the proof of their Lemma 1 in their Appendix G is erroneous, and consequently their Theorem 2.1 is unsound as it heavily depends on their Lemma 1. In their Lemma 1, they claimed that for some constant  $M > 0$ , for any two models with the parameters  $\theta$  and  $\tilde{\theta}$  such that  $\theta, \tilde{\theta} \in B(\theta^{(0)}, C_0)$  for some constant  $C_0$ , there is

$$\|J(\theta) - J(\tilde{\theta})\|_F \leq \frac{M}{\sqrt{\tilde{d}}} \|\theta - \tilde{\theta}\|_2 \quad (\text{claimed}) \quad (168)$$

Note that the original claim in their paper was  $\|J(\theta) - J(\tilde{\theta})\|_F \leq M\sqrt{\tilde{d}} \|\theta - \tilde{\theta}\|_2$ . This is because they were proving this result for their standard formulation. Compared to the standard formulation, in the NTK formulation  $\theta$  is  $\sqrt{\tilde{d}}$  times larger, while the Jacobian  $J(\theta)$  is  $\sqrt{\tilde{d}}$  times smaller. This is also why here we have  $\theta, \tilde{\theta} \in B(\theta^{(0)}, C_0)$  instead of  $\theta, \tilde{\theta} \in B(\theta^{(0)}, C_0\tilde{d}^{-1/2})$  for the NTK formulation. Therefore, equivalently they were claiming (168) for the NTK formulation.

However, their proof of (168) is incorrect. Specifically, the right-hand side of their inequality (S86) is incorrect. Using the notations in our Appendix D.3.3, their (S86) essentially claimed that

$$\|\alpha^l - \tilde{\alpha}^l\|_2 \leq \frac{M}{\sqrt{\tilde{d}}} \|\theta - \tilde{\theta}\|_2 \quad (\text{claimed}) \quad (169)$$

for any  $\theta, \tilde{\theta} \in B(\theta^{(0)}, C_0)$ , where  $\alpha^l = \nabla_{\mathbf{h}^l} \mathbf{h}^{L+1}$  and  $\tilde{\alpha}^l$  is the same gradient for the second model. Note that their (S86) does not have the  $\sqrt{\tilde{d}}$  in the denominator which appears in (169). This is because for their standard formulation,  $\theta$  is  $\sqrt{\tilde{d}}$  times smaller than the original NTK formulation, while  $\|\alpha^l\|_2$  has the same order in the two formulations because all  $\mathbf{h}^l$  are the same.

However, it is actually impossible to prove (169). Consider the following counterexample: Since  $\theta$  and  $\tilde{\theta}$  are arbitrarily chosen, we can choose them such that they only differ in  $b_1^l$  for some  $1 \leq l < L$ . Then,  $\|\theta - \tilde{\theta}\|_2 = |b_1^l - \tilde{b}_1^l|$ . We can see that  $\mathbf{h}^{l+1}$  and  $\tilde{\mathbf{h}}^{l+1}$  only differ in the first element, and



531  $|h_1^{l+1} - \tilde{h}_1^{l+1}| = |\beta(b_1^l - \tilde{b}_1^l)|$ . Moreover, we have  $W^{l+1} = \tilde{W}^{l+1}$ , so there is

$$\begin{aligned} \alpha^{l+1} - \tilde{\alpha}^{l+1} &= \text{diag}(\dot{\sigma}(\mathbf{h}^{l+1})) \frac{W^{l+1\top}}{\sqrt{\tilde{d}}} \alpha^{l+2} - \text{diag}(\dot{\sigma}(\tilde{\mathbf{h}}^{l+1})) \frac{\tilde{W}^{l+1\top}}{\sqrt{\tilde{d}}} \tilde{\alpha}^{l+2} \\ &= \left[ \text{diag}(\dot{\sigma}(\mathbf{h}^{l+1})) - \text{diag}(\dot{\sigma}(\tilde{\mathbf{h}}^{l+1})) \right] \frac{W^{l+1\top}}{\sqrt{\tilde{d}}} \alpha^{l+2} \\ &\quad + \text{diag}(\dot{\sigma}(\tilde{\mathbf{h}}^{l+1})) \frac{W^{l+1\top}}{\sqrt{\tilde{d}}} (\alpha^{l+2} - \tilde{\alpha}^{l+2}) \end{aligned} \quad (170)$$

532 Then we can lower bound  $\|\alpha^{l+1} - \tilde{\alpha}^{l+1}\|_2$  by

$$\begin{aligned} \|\alpha^{l+1} - \tilde{\alpha}^{l+1}\|_2 &\geq \left\| \left[ \text{diag}(\dot{\sigma}(\mathbf{h}^{l+1})) - \text{diag}(\dot{\sigma}(\tilde{\mathbf{h}}^{l+1})) \right] \frac{W^{l+1\top}}{\sqrt{\tilde{d}}} \alpha^{l+2} \right\|_2 \\ &\quad - \left\| \text{diag}(\dot{\sigma}(\tilde{\mathbf{h}}^{l+1})) \frac{W^{l+1\top}}{\sqrt{\tilde{d}}} (\alpha^{l+2} - \tilde{\alpha}^{l+2}) \right\|_2 \end{aligned} \quad (171)$$

533 The first term on the right-hand side is equal to  $\left| \left[ \dot{\sigma}(h_1^{l+1}) - \dot{\sigma}(\tilde{h}_1^{l+1}) \right] \langle W_1^{l+1} / \sqrt{\tilde{d}}, \alpha^{l+2} \rangle \right|$  where  
 534  $W_1^{l+1}$  is the first row of  $W^{l+1}$ . We know that  $\|W_1^{l+1}\|_2 = \Theta(\sqrt{\tilde{d}})$  with high probability as its  
 535 elements are sampled from  $\mathcal{N}(0, 1)$ , and in their (S85) they claimed that  $\|\alpha^{l+2}\|_2 = O(1)$ , which is  
 536 true. In addition, they assumed that  $\dot{\sigma}$  is Lipschitz. Hence, we can see that

$$\left\| \left[ \text{diag}(\dot{\sigma}(\mathbf{h}^{l+1})) - \text{diag}(\dot{\sigma}(\tilde{\mathbf{h}}^{l+1})) \right] \frac{W^{l+1\top}}{\sqrt{\tilde{d}}} \alpha^{l+2} \right\|_2 = O\left(|h_1^{l+1} - \tilde{h}_1^{l+1}|\right) = O\left(\|\theta - \tilde{\theta}\|_2\right) \quad (172)$$

537 On the other hand, suppose that claim (169) is true, then  $\|\alpha^{l+2} - \tilde{\alpha}^{l+2}\|_2 = O\left(\tilde{d}^{-1/2} \|\theta - \tilde{\theta}\|_2\right)$ .

538 Then we can see that the second term on the right-hand side is  $O\left(\tilde{d}^{-1/2} \|\theta - \tilde{\theta}\|_2\right)$  because

539  $\|W^{l+1}\|_2 = O(\sqrt{\tilde{d}})$  and  $\dot{\sigma}(x)$  is bounded by a constant as  $\sigma$  is Lipschitz. Thus, for a very large  $\tilde{d}$ ,  
 540 the second-term is an infinitesimal compared to the first term, so we can only prove that

$$\|\alpha^{l+1} - \tilde{\alpha}^{l+1}\|_2 = O\left(\|\theta - \tilde{\theta}\|_2\right) \quad (173)$$

541 which is different from (169) because it lacks a critical  $\tilde{d}^{-1/2}$  and thus leads to a contradiction. Hence,  
 542 we cannot prove (169) with the  $\tilde{d}^{-1/2}$  factor, and consequently we cannot prove (168) with the  $\sqrt{\tilde{d}}$   
 543 in the denominator on the right-hand side. As a result, their Lemma 1 and Theorem 2.1 cannot be  
 544 proved without this critical  $\tilde{d}^{-1/2}$ . Similarly, we can also construct a counterexample where  $\theta$  and  $\tilde{\theta}$   
 545 only differ in the first row of some  $W^l$ .

## 546 E.2 Our Fixes

547 Regarding Problem 1, we can still use an  $O(1)$  learning rate for the first layer in the NTK formulation  
 548 given that  $\|\mathbf{x}\|_2 \leq 1$ . This is because for the first layer, we have

$$\nabla_{W^0} f(\mathbf{x}) = \frac{1}{\sqrt{d_0}} \mathbf{x}^0 \alpha^{1\top} = \frac{1}{\sqrt{d_0}} \mathbf{x} \alpha^{1\top} \quad (174)$$

549 For all  $l \geq 1$ , we have  $\|\mathbf{x}^l\|_2 = O(\tilde{d}^{1/2})$ . However, for  $l = 0$ , we instead have  $\|\mathbf{x}^0\|_2 = O(1)$ . Thus,  
 550 we can prove that the norm of  $\nabla_{W^0} f(\mathbf{x})$  has the same order as the gradient with respect to any other  
 551 layer, so there is no need to use a smaller learning rate for the first layer.

552 Regarding Problem 2, in our formulation (8) and initialization (9), the initialization of the last layer  
 553 of the NTK formulation is changed from the Gaussian initialization  $W_{i,j}^{L(0)} \sim \mathcal{N}(0, 1)$  to the zero  
 554 initialization  $W_{i,j}^{L(0)} = 0$ . Now we show how this modification solves Problem 2.

The main consequence of changing the initialization of the last layer is that (86) becomes different: instead of  $\|W^L\|_2 \leq 3\sqrt{\tilde{d}}$ , we now have  $\|W^L\|_2 \leq C_0 \leq 3\sqrt[4]{\tilde{d}}$ . In fact, for any  $r \in (0, 1/2)$ , we can prove that  $\|W^L\|_2 \leq 3\tilde{d}^r$  for sufficiently large  $\tilde{d}$ . In our proof we choose  $r = 1/4$ .

Consequently, instead of  $\|\alpha^l\|_2 \leq M_3$ , we can now prove that  $\|\alpha^l\|_2 \leq M_3\tilde{d}^{r-1/2}$  for all  $l \leq L$  by induction. So now we can prove  $\|\alpha^l - \tilde{\alpha}^l\|_2 = O(\tilde{d}^{r-1/2}\|\theta - \tilde{\theta}\|_2)$  instead of  $O(\|\theta - \tilde{\theta}\|_2)$ , because

- For  $l < L$ , we now have  $\|\alpha^{l+1}\|_2 = O(\tilde{d}^{r-1/2})$  instead of  $O(1)$ , so we can have the additional  $\tilde{d}^{r-1/2}$  factor in the bound.
- For  $l = L$ , although  $\|\alpha^{L+1}\|_2 = 1$ , note that  $\|W^L\|_2$  now becomes  $O(\tilde{d}^r)$  instead of  $O(\tilde{d}^{1/2})$ , so again we can decrease the bound by a factor of  $\tilde{d}^{r-1/2}$ .

Then, with this critical  $\tilde{d}^{r-1/2}$ , we can prove the approximation theorem with the form

$$\sup_{t \geq 0} \left| f^{(t)}(\mathbf{x}) - f_{\text{lin}}^{(t)}(\mathbf{x}) \right| \leq C\tilde{d}^{r-1/2} \quad (175)$$

for any  $r \in (0, 1/2)$ , though we cannot really prove the  $O(\tilde{d}^{-1/2})$  bound as originally claimed in (164). So this is how we solve Problem 2.

One caveat of changing the initialization to zero initialization is whether we can still safely assume that  $\lambda^{\min} > 0$  where  $\lambda^{\min}$  is the smallest eigenvalue of  $\Theta$ , the kernel matrix of our new formulation. The answer is yes. In fact, in our Proposition 3 we proved that  $\Theta$  is non-degenerated (which means that  $\Theta(\mathbf{x}, \mathbf{x}')$  still depends on  $\mathbf{x}$  and  $\mathbf{x}'$ ), and under the overparameterized setting where  $d_L \gg n$ , chances are high that  $\Theta$  is full-rank. Hence, we can still assume that  $\lambda^{\min} > 0$ .

As a final remark, one key reason why we need to initialize  $W^L$  as zero is that the dimension of the output space (i.e. the dimension of  $\mathbf{h}^{L+1}$ ) is finite, and in our case it is 1. Suppose we allow the dimension of  $\mathbf{h}^{L+1}$  to be  $\tilde{d}$  which goes to infinity, then using the same proof techniques, for the NTK formulation we can prove that  $\sup_t \left\| \mathbf{h}^{L+1(t)} - \mathbf{h}_{\text{lin}}^{L+1(t)} \right\|_2 \leq C$ , i.e. the gap between two vectors of infinite dimension is always bounded by a finite constant. This is the approximation theorem we need for the infinite-dimensional output space. However, when the dimension of the output space is finite,  $\sup_t \left\| \mathbf{h}^{L+1(t)} - \mathbf{h}_{\text{lin}}^{L+1(t)} \right\|_2 \leq C$  no longer suffices, so we need to decrease the order of the norm of  $W^L$  in order to obtain a smaller upper bound.

## References

- [AZLS19] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In *International Conference on Machine Learning*, pages 242–252. PMLR, 2019.
- [COB19] Lénaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [DHP<sup>+</sup>12] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*, pages 214–226, 2012.
- [DLL<sup>+</sup>19] Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In *International Conference on Machine Learning*, pages 1675–1685. PMLR, 2019.
- [DN18] John Duchi and Hongseok Namkoong. Learning models with uniform performance via distributionally robust optimization. *arXiv preprint arXiv:1810.08750*, 2018.

[HPS16] Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29, pages 3315–3323. Curran Associates, Inc., 2016.

[HSNL18] Tatsunori Hashimoto, Megha Srivastava, Hongseok Namkoong, and Percy Liang. Fairness without demographics in repeated loss minimization. In Jennifer Dy and Andreas Krause, editors, *International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1929–1938, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR.

[JDST20] Ziwei Ji, Miroslav Dudík, Robert E. Schapire, and Matus Telgarsky. Gradient descent follows the regularization path for general losses. In Jacob Abernethy and Shivani Agarwal, editors, *Proceedings of Thirty Third Conference on Learning Theory*, volume 125 of *Proceedings of Machine Learning Research*, pages 2109–2136. PMLR, 09–12 Jul 2020.

[JGH18] Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.

[KLRS17] Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness. In *Advances in neural information processing systems*, pages 4066–4076, 2017.

[LHC<sup>+</sup>21] Evan Z Liu, Behzad Haghighi, Annie S Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa, Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness without training group information. In *International Conference on Machine Learning*, pages 6781–6792. PMLR, 2021.

[LXS<sup>+</sup>19] Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. *Advances in neural information processing systems*, 32:8572–8583, 2019.

[OSHL19] Yonatan Oren, Shiori Sagawa, Tatsunori Hashimoto, and Percy Liang. Distributionally robust language modeling. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4227–4237, Hong Kong, China, November 2019. Association for Computational Linguistics.

[Raw01] John Rawls. *Justice as fairness: A restatement*. Harvard University Press, 2001.

[SKHL20] Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. In *International Conference on Learning Representations*, 2020.

[SRKL20] Shiori Sagawa, Aditi Raghunathan, Pang Wei Koh, and Percy Liang. An investigation of why overparameterization exacerbates spurious correlations. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 8346–8356. PMLR, 13–18 Jul 2020.

[Ver10] Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.

[ZDKR21] Runtian Zhai, Chen Dan, Zico Kolter, and Pradeep Ravikumar. Doro: Distributional and outlier robust optimization. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 12345–12355. PMLR, 18–24 Jul 2021.

- 646 [ZVGRG17] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, and Krishna P Gum-  
647 madi. Fairness beyond disparate treatment & disparate impact: Learning classification  
648 without disparate mistreatment. In *Proceedings of the 26th international conference on*  
649 *world wide web*, pages 1171–1180, 2017.
- 650 [ZWS<sup>+</sup>13] Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. Learning fair  
651 representations. In *International Conference on Machine Learning*, pages 325–333,  
652 2013.