
How Much Over-parameterization Is Sufficient to Learn Deep ReLU Networks?

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A recent line of research on deep learning focuses on the extremely over-
2 parameterized setting, and shows that when the network width is larger than a high
3 degree polynomial of the training sample size n and the inverse of the target error
4 ϵ^{-1} , deep neural networks learned by (stochastic) gradient descent enjoy nice opti-
5 mization and generalization guarantees. Very recently, it is shown that under certain
6 margin assumptions on the training data, a polylogarithmic width condition suffices
7 for two-layer ReLU networks to converge and generalize [15]. However, whether
8 deep neural networks can be learned with such a mild over-parameterization is
9 still an open question. In this work, we answer this question affirmatively and
10 establish sharper learning guarantees for deep ReLU networks trained by (stochas-
11 tic) gradient descent. In specific, under certain assumptions made in previous
12 work, our optimization and generalization guarantees hold with network width
13 polylogarithmic in n and ϵ^{-1} . Our results push the study of over-parameterized
14 deep neural networks towards more practical settings.

15 1 Introduction

16 Deep neural networks have become one of the most important and prevalent machine learning models
17 due to their remarkable power in many real-world applications. However, the success of deep learning
18 has not been well-explained in theory. It remains mysterious why standard optimization algorithms
19 tend to find a globally optimal solution, despite the highly non-convex landscape of the training loss
20 function. Moreover, despite the extremely large amount of parameters, deep neural networks rarely
21 over-fit, and can often generalize well to unseen data and achieve good test accuracy. Understanding
22 these mysterious phenomena on the optimization and generalization of deep neural networks is one
23 of the most fundamental problems in deep learning theory.

24 Recent breakthroughs have shed light on the optimization [12, 2, 23, 24] and generalization Allen-
25 Zhu et al. [1], Arora et al. [4], Cao and Gu [8] of deep neural networks (DNNs) under the over-
26 parameterized setting, where the hidden layer width is extremely large, which is typically a high
27 degree polynomial of the training sample size n and the inverse of the target error ϵ^{-1} . As there
28 still remains a huge gap between such network width requirement and the practice, many attempts
29 have been made to improve the over-parameterization condition. For two-layer ReLU networks, a
30 recent work [15] showed that when the training data are well separated, polylogarithmic width is
31 sufficient to guarantee good optimization and generalization performances. However, their results
32 cannot be extended to deep ReLU networks since their proof technique largely relies on the fact that
33 the network model is 1-homogeneous, which cannot be satisfied by DNNs. Therefore, whether deep
34 neural networks can be learned with such a mild over-parameterization is still an open problem.

35 In this paper, we resolve this open problem by showing that polylogarithmic network width is
 36 sufficient to learn DNNs. In particular, unlike the existing works that require the DNNs to behave
 37 very close to a linear model (up to some small approximation error), we show that a constant linear
 38 approximation error is sufficient to establish nice optimization and generalization guarantees for
 39 DNNs. Thanks to the relaxed requirement on the linear approximation error, a milder condition on
 40 the network width and tighter bounds on the convergence rate and generalization error can be proved.
 41 We summarize our contributions as follows:

- 42 • We establish the global convergence guarantee of GD for training deep ReLU networks based on
 43 the so-called NTRF function class [8], a set of linear functions over random features. Specifically,
 44 we prove that GD can learn deep ReLU networks with width $m = \text{poly}(R)$ to compete with the
 45 best function in NTRF function class, where R is the radius of the NTRF function class, which
 46 can be demonstrated to be $\tilde{O}(1)$ under commonly used data separability assumptions.
- 47 • We also establish the generalization guarantees for both GD and SGD in the same setting. Specifi-
 48 cally, we prove a diminishing statistical error for a wide range of network width $m \in (\tilde{\Omega}(1), \infty)$,
 49 while most of the previous generalization bounds in the NTK regime only works in the setting
 50 where the network width m is much greater than the sample size n . Moreover, we establish $\tilde{O}(\epsilon^{-2})$
 51 $\tilde{O}(\epsilon^{-1})$ sample complexities for GD and SGD respectively, which are tighter than existing bounds
 52 for learning deep ReLU networks [8], and match the best results when reduced to the two-layer
 53 cases [5, 15].

54 For the ease of comparison, we summarize our results along with the most related previous results in
 55 Table 1, in terms of data assumption, the over-parameterization condition and sample complexity.
 56 It can be seen that under data separation assumption (See Sections A.1, A.2), our result improves
 57 existing results for learning deep neural networks by only requiring a $\text{polylog}(n, \epsilon^{-1})$ network width.

Table 1: Comparison of neural network learning results in terms of over-parameterization condition and sample complexity. Here ϵ is the target error rate, n is the sample size, L is the network depth.

	Assumptions	Algorithm	Over-para. Condition	Sample Complexity	Network
Zou et al. [23]	Data nondegeneration	GD	$\tilde{\Omega}(n^{12}L^{16}(n^2 + \epsilon^{-1}))$	-	Deep
This paper	Data nondegeneration	GD	$\tilde{\Omega}(L^{22}n^{12})$	-	Deep
Cao and Gu [9]	Data separation	GD	$\tilde{\Omega}(\epsilon^{-14}) \cdot e^{\Omega(L)}$	$\tilde{O}(\epsilon^{-4}) \cdot e^{O(L)}$	Deep
Ji and Telgarsky [15]	Data separation	GD	$\text{polylog}(n, \epsilon^{-1})$	$\tilde{O}(\epsilon^{-2})$	Shallow
This paper	Data separation	GD	$\text{polylog}(n, \epsilon^{-1}) \cdot \text{poly}(L)$	$\tilde{O}(\epsilon^{-2}) \cdot e^{O(L)}$	Deep
Cao and Gu [8]	Data separation	SGD	$\tilde{\Omega}(\epsilon^{-14}) \cdot \text{poly}(L)$	$\tilde{O}(\epsilon^{-2}) \cdot \text{poly}(L)$	Deep
Ji and Telgarsky [15]	Data separation	SGD	$\text{polylog}(\epsilon^{-1})$	$\tilde{O}(\epsilon^{-1})$	Shallow
This paper	Data separation	SGD	$\text{polylog}(\epsilon^{-1}) \cdot \text{poly}(L)$	$\tilde{O}(\epsilon^{-1}) \cdot \text{poly}(L)$	Deep

58 2 Preliminaries on Learning Neural Networks

59 In this section, we introduce the problem setting in this paper, including definitions of the neural
 60 network and loss functions, and the training algorithms, i.e., GD and SGD with random initialization.

61 **Neural network function.** Given an input $\mathbf{x} \in \mathbb{R}^d$, the output of deep fully-connected ReLU network
 62 is defined as follows,

$$f_{\mathbf{W}}(\mathbf{x}) = m^{1/2} \mathbf{W}_L \sigma(\mathbf{W}_{L-1} \cdots \sigma(\mathbf{W}_1 \mathbf{x}) \cdots),$$

63 where $\mathbf{W}_1 \in \mathbb{R}^{m \times d}$, $\mathbf{W}_2, \dots, \mathbf{W}_{L-1} \in \mathbb{R}^{m \times m}$ and $\mathbf{W}_L \in \mathbb{R}^{1 \times m}$. We denote the collection of all
 64 weight matrices as $\mathbf{W} = \{\mathbf{W}_1, \dots, \mathbf{W}_L\}$.

65 **Loss function.** Given training dataset $\{\mathbf{x}_i, y_i\}_{i=1, \dots, n}$ with input $\mathbf{x}_i \in \mathbb{R}^d$ and output $y_i \in \{-1, +1\}$,
 66 we define the training loss function as

$$L_S(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n L_i(\mathbf{W}),$$

67 where $L_i(\mathbf{W}) = \ell(y_i f_{\mathbf{W}}(\mathbf{x}_i)) = \log(1 + \exp(-y_i f_{\mathbf{W}}(\mathbf{x}_i)))$ is defined as the cross-entropy loss.

68 **3 Main Theory**

69 In this section, we present the optimization and generalization guarantees of GD and SGD for learning
 70 deep ReLU networks. We first make the following assumption on the training data points.

71 **Assumption 3.1.** All training data points satisfy $\|\mathbf{x}_i\|_2 = 1, i = 1, \dots, n$.

72 This assumption has been widely made in many previous works [2, 3, 12, 11, 23] in order to simplify
 73 the theoretical analysis.

74 In the following, we give the definition of Neural Tangent Random Feature (NTRF) [8], which
 75 characterizes the functions learnable by over-parameterized ReLU networks.

76 **Definition 3.2** (Neural Tangent Random Feature, [8]). Let $\mathbf{W}^{(0)}$ be the initialization weights, and
 77 $F_{\mathbf{W}^{(0)}, \mathbf{W}}(\mathbf{x}) = f_{\mathbf{W}^{(0)}}(\mathbf{x}) + \langle \nabla f_{\mathbf{W}^{(0)}}(\mathbf{x}), \mathbf{W} - \mathbf{W}^{(0)} \rangle$ be a function with respect to the input \mathbf{x} .
 78 Then the NTRF function class is defined as follows

$$\mathcal{F}(\mathbf{W}^{(0)}, R) = \{F_{\mathbf{W}^{(0)}, \mathbf{W}}(\cdot) : \mathbf{W} \in \mathcal{B}(\mathbf{W}^{(0)}, R \cdot m^{-1/2})\}.$$

79 The function class $F_{\mathbf{W}^{(0)}, \mathbf{W}}(\mathbf{x})$ consists of linear models over random features defined based on
 80 the network gradients at the initialization. Therefore it captures the key ‘‘almost linear’’ property of
 81 wide neural networks in the NTK regime [17, 8]. In this paper, we use the NTRF function class as a
 82 reference class to measure the difficulty of a learning problem. In what follows, we deliver our main
 83 theoretical results regarding the optimization and generalization guarantees of learning deep ReLU
 84 networks. We study both GD and SGD with random initialization.

85 **3.1 Gradient Descent**

86 The following theorem establishes the optimization guarantee of GD for training deep ReLU networks
 87 for binary classification.

88 **Theorem 3.3.** For $\delta, R > 0$, let $\epsilon_{\text{NTRF}} = \inf_{F \in \mathcal{F}(\mathbf{W}^{(0)}, R)} n^{-1} \sum_{i=1}^n \ell[y_i F(\mathbf{x}_i)]$ be the minimum
 89 training loss achievable by functions in $\mathcal{F}(\mathbf{W}^{(0)}, R)$. Then there exists

$$m^*(\delta, R, L) = \tilde{\mathcal{O}}(\text{poly}(R, L) \cdot \log^{4/3}(n/\delta)),$$

90 such that if $m \geq m^*(\delta, R, L)$, with probability at least $1 - \delta$ over the initialization, GD with step
 91 size $\eta = \Theta(L^{-1}m^{-1})$ can train a neural network to achieve at most $3\epsilon_{\text{NTRF}}$ training loss within
 92 $T = \mathcal{O}(L^2 R^2 \epsilon_{\text{NTRF}}^{-1})$ iterations.

93 Theorem 3.3 shows that the deep ReLU network trained by GD can compete with the best function in
 94 the NTRF function class $\mathcal{F}(\mathbf{W}^{(0)}, R)$ if the network width has a polynomial dependency in R and L
 95 and a logarithmic dependency in n and $1/\delta$. Moreover, if the NTRF function class with $R = \tilde{\mathcal{O}}(1)$
 96 can learn the training data well (i.e., ϵ_{NTRF} is less than a small target error ϵ), a polylogarithmic (in
 97 terms of n and ϵ^{-1}) network width suffices to guarantee the global convergence of GD, which directly
 98 improves over-paramterization condition in the most related work [8]. In Appendix A, we show that
 99 under commonly used data separability assumptions, NTRF function class with $R = \text{polylog}(n, \epsilon^{-1})$
 100 can achieve $\epsilon_{\text{NTRF}} \leq \epsilon$ for arbitrarily small $\epsilon > 0$. Moreover, under a much weaker data assumption
 101 which covers the case of random labels, we also have $\epsilon_{\text{NTRF}} \leq \epsilon$ for $R = \Omega(n^{3/2} \log(n/\epsilon))$, which
 102 implies global convergence of GD when $m = \tilde{\Omega}(n^{12})$. For all cases, our over-parameterization
 103 requirement is better than existing results for DNNs.

104 Compared with the results in [15] which give similar network width requirements for two-layer
 105 networks, our result works for deep networks. Moreover, while Ji and Telgarsky [15] essentially
 106 required all training data to be separable by a function in the NTRF function class with a constant
 107 margin, our result does not require such data separation assumptions, and allows the NTRF function
 108 class to misclassify a small proportion of the training data points¹.

109 We now characterize the generalization performance of neural networks trained by GD. We denote
 110 $L_{\mathcal{D}}^{0-1}(\mathbf{W}) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[\mathbb{1}\{f_{\mathbf{W}}(\mathbf{x}) \cdot y < 0\}]$ as the expected 0-1 loss (i.e., expected error) of $f_{\mathbf{W}}(\mathbf{x})$.

¹A detailed discussion is given in Section A.2.

111 **Theorem 3.4.** Under the same assumptions as Theorem 3.3, with probability at least $1 - \delta$, the iterate
 112 $\mathbf{W}^{(t)}$ of gradient descent satisfies that

$$L_{\mathcal{D}}^{0-1}(\mathbf{W}^{(t)}) \leq 2L_S(\mathbf{W}^{(t)}) + \tilde{\mathcal{O}}\left(4^L L^2 R \sqrt{\frac{m}{n}} \wedge \left(\frac{L^{3/2} R}{\sqrt{n}} + \frac{L^{11/3} R^{4/3}}{m^{1/6}}\right)\right) + \mathcal{O}\left(\sqrt{\frac{\log(1/\delta)}{n}}\right)$$

113 for all $t = 0, \dots, T$.

114 Theorem 3.4 shows that the test error of the trained neural network can be bounded by its training error
 115 plus statistical error terms. Note that the statistical error terms is in the form of a minimum between
 116 two terms $4^L L^2 R \sqrt{m/n}$ and $L^{3/2} R / \sqrt{n} + L^{11/3} R^{4/3} / m^{1/6}$. Depending on the network width m ,
 117 one of these two terms will be the dominating term and diminishes for large n : (1) if $m = o(n)$,
 118 the statistical error will be $4^L L^2 R \sqrt{m/n}$, and diminishes as n increases; and (2) if $m = \Omega(n)$, the
 119 statistical error is $L^{3/2} R / \sqrt{n} + L^{11/3} R^{4/3} / m^{1/6}$, and again goes to zero as n increases. Moreover,
 120 in this paper we have a specific focus on the setting $m = \tilde{\mathcal{O}}(1)$, under which Theorem 3.4 gives a
 121 statistical error of order $\tilde{\mathcal{O}}(n^{-1/2})$. This distinguishes our result from previous generalization bounds
 122 for deep networks [9, 8], which cannot be applied to the setting $m = \tilde{\mathcal{O}}(1)$.

123 3.2 Stochastic Gradient Descent

124 Here we study the performance of SGD for training deep ReLU networks. The following theorem
 125 establishes a generalization error bound for the output of SGD.

126 **Theorem 3.5.** For $\delta, R > 0$, let $\epsilon_{\text{NTRF}} = \inf_{F \in \mathcal{F}(\mathbf{W}^{(0)}, R)} n^{-1} \sum_{i=1}^n \ell[y_i F(\mathbf{x}_i)]$ be the minimum
 127 training loss achievable by functions in $\mathcal{F}(\mathbf{W}^{(0)}, R)$. Then there exists

$$m^*(\delta, R, L) = \tilde{\mathcal{O}}(\text{poly}(R, L) \cdot \log^{4/3}(n/\delta)),$$

128 such that if $m \geq m^*(\delta, R, L)$, with probability at least $1 - \delta$, SGD with step size $\eta = \Theta(m^{-1} \cdot$
 129 $(LR^2 n^{-1} \epsilon_{\text{NTRF}}^{-1} \wedge L^{-1}))$ achieves

$$\mathbb{E}[L_{\mathcal{D}}^{0-1}(\widehat{\mathbf{W}})] \leq \frac{8L^2 R^2}{n} + \frac{8 \log(2/\delta)}{n} + 24\epsilon_{\text{NTRF}},$$

130 where the expectation is taken over the uniform draw of $\widehat{\mathbf{W}}$ from $\{\mathbf{W}^{(0)}, \dots, \mathbf{W}^{(n-1)}\}$.

131 For any $\epsilon > 0$, Theorem 3.5 gives a $\tilde{\mathcal{O}}(\epsilon^{-1})$ sample complexity for deep ReLU networks trained with
 132 SGD to achieve $O(\epsilon_{\text{NTRF}} + \epsilon)$ test error. Again, under commonly used data separability assumptions,
 133 NTRF function class with $R = \text{polylog}(n, \epsilon^{-1})$ can achieve $\epsilon_{\text{NTRF}} \leq \epsilon$ for arbitrarily small $\epsilon > 0$
 134 (See Appendix A), which implies an $m = \tilde{\Omega}(1)$ over-parameterization condition and an $n = \tilde{\omega}(\epsilon^{-1})$
 135 sample complexity. Our result extends the result for two-layer networks proved in [15] to multi-layer
 136 networks. Theorem 3.5 also provides sharper results compared with Allen-Zhu et al. [1], Cao and Gu
 137 [8] in two aspects: (1) the sample complexity is improved from $n = \tilde{\mathcal{O}}(\epsilon^{-2})$ to $n = \tilde{\mathcal{O}}(\epsilon^{-1})$; and (2)
 138 the overparamterization condition is improved from $m \geq \text{poly}(\epsilon^{-1})$ to $m = \tilde{\Omega}(1)$.

139 4 Conclusion

140 In this paper, we established the global convergence and generalization error bounds of GD and SGD
 141 for training deep ReLU networks for the binary classification problem. We show that a network width
 142 condition that is polylogarithmic in the sample size n and the inverse of target error ϵ^{-1} is sufficient
 143 to guarantee the learning of deep ReLU networks. Our results resolve an open question raised in Ji
 144 and Telgarsky [15].

References

- 145
- 146 [1] ALLEN-ZHU, Z., LI, Y. and LIANG, Y. (2019). Learning and generalization in overparameter-
147 ized neural networks, going beyond two layers. In *Advances in Neural Information Processing*
148 *Systems*.
- 149 [2] ALLEN-ZHU, Z., LI, Y. and SONG, Z. (2019). A convergence theory for deep learning via
150 over-parameterization. In *International Conference on Machine Learning*.
- 151 [3] ALLEN-ZHU, Z., LI, Y. and SONG, Z. (2019). On the convergence rate of training recurrent
152 neural networks. In *Advances in Neural Information Processing Systems*.
- 153 [4] ARORA, S., DU, S., HU, W., LI, Z. and WANG, R. (2019). Fine-grained analysis of opti-
154 mization and generalization for overparameterized two-layer neural networks. In *International*
155 *Conference on Machine Learning*.
- 156 [5] ARORA, S., DU, S. S., HU, W., LI, Z., SALAKHUTDINOV, R. and WANG, R. (2019). On exact
157 computation with an infinitely wide neural net. In *Advances in Neural Information Processing*
158 *Systems*.
- 159 [6] BARTLETT, P. L., FOSTER, D. J. and TELGARSKY, M. J. (2017). Spectrally-normalized
160 margin bounds for neural networks. In *Advances in Neural Information Processing Systems*.
- 161 [7] BARTLETT, P. L. and MENDELSON, S. (2002). Rademacher and Gaussian complexities: Risk
162 bounds and structural results. *Journal of Machine Learning Research* **3** 463–482.
- 163 [8] CAO, Y. and GU, Q. (2019). Generalization bounds of stochastic gradient descent for wide and
164 deep neural networks. In *Advances in Neural Information Processing Systems*.
- 165 [9] CAO, Y. and GU, Q. (2020). Generalization error bounds of gradient descent for learning
166 over-parameterized deep relu networks. In *the Thirty-Fourth AAAI Conference on Artificial*
167 *Intelligence*.
- 168 [10] CESA-BIANCHI, N., CONCONI, A. and GENTILE, C. (2004). On the generalization ability of
169 on-line learning algorithms. *IEEE Transactions on Information Theory* **50** 2050–2057.
- 170 [11] DU, S., LEE, J., LI, H., WANG, L. and ZHAI, X. (2019). Gradient descent finds global minima
171 of deep neural networks. In *International Conference on Machine Learning*.
- 172 [12] DU, S. S., ZHAI, X., POZOS, B. and SINGH, A. (2019). Gradient descent provably optimizes
173 over-parameterized neural networks. In *International Conference on Learning Representations*.
- 174 [13] FREI, S., CAO, Y. and GU, Q. (2019). Algorithm-dependent generalization bounds for
175 overparameterized deep residual networks. In *Advances in Neural Information Processing*
176 *Systems*.
- 177 [14] JI, Z. and TELGARSKY, M. (2018). Risk and parameter convergence of logistic regression.
178 *arXiv preprint arXiv:1803.07300*.
- 179 [15] JI, Z. and TELGARSKY, M. (2020). Polylogarithmic width suffices for gradient descent to
180 achieve arbitrarily small test error with shallow relu networks. In *International Conference on*
181 *Learning Representations*.
- 182 [16] KRIZHEVSKY, A. ET AL. (2009). Learning multiple layers of features from tiny images .
- 183 [17] LEE, J., XIAO, L., SCHOENHOLZ, S. S., BAHRI, Y., SOHL-DICKSTEIN, J. and PENNINGTON,
184 J. (2019). Wide neural networks of any depth evolve as linear models under gradient descent.
185 In *Advances in Neural Information Processing Systems*.
- 186 [18] MOHRI, M., ROSTAMIZADEH, A. and TALWALKAR, A. (2018). *Foundations of machine*
187 *learning*. MIT press.

- 188 [19] NITANDA, A. and SUZUKI, T. (2019). Refined generalization analysis of gradient descent
189 for over-parameterized two-layer neural networks with smooth activations on classification
190 problems. *arXiv preprint arXiv:1905.09870* .
- 191 [20] SHALEV-SHWARTZ, S. and BEN-DAVID, S. (2014). *Understanding machine learning: From*
192 *theory to algorithms*. Cambridge university press.
- 193 [21] SHAMIR, O. (2020). Gradient methods never overfit on separable data. *arXiv preprint*
194 *arXiv:2007.00028* .
- 195 [22] VERSHYNIN, R. (2010). Introduction to the non-asymptotic analysis of random matrices. *arXiv*
196 *preprint arXiv:1011.3027* .
- 197 [23] ZOU, D., CAO, Y., ZHOU, D. and GU, Q. (2019). Gradient descent optimizes over-
198 parameterized deep ReLU networks. *Machine Learning* .
- 199 [24] ZOU, D. and GU, Q. (2019). An improved analysis of training over-parameterized deep neural
200 networks. In *Advances in Neural Information Processing Systems*.

201 **A Discussion on the NTRF Class**

202 Our theoretical results in Section 3 rely on the radius (i.e., R) of the NTRF function class $\mathcal{F}(\mathbf{W}^{(0)}, R)$
 203 and the minimum training loss achievable by functions in $\mathcal{F}(\mathbf{W}^{(0)}, R)$, i.e., ϵ_{NTRF} . Note that a larger
 204 R naturally implies a smaller ϵ_{NTRF} , but also leads to worse conditions on m . In this section, for
 205 any (arbitrarily small) target error rate $\epsilon > 0$, we discuss various data assumptions studied in the
 206 literature under which our results can lead to $\mathcal{O}(\epsilon)$ training/test errors, and specify the network width
 207 requirement.

208 **A.1 Data Separability by Neural Tangent Random Feature**

209 In this subsection, we consider the setting where a large fraction of the training data can be linearly
 210 separated by the neural tangent random features. The assumption is stated as follows.

211 **Assumption A.1.** There exists a collection of matrices $\mathbf{U}^* = \{\mathbf{U}_1^*, \dots, \mathbf{U}_L^*\}$ satisfying
 212 $\sum_{l=1}^L \|\mathbf{U}_l^*\|_F^2 = 1$, such that for at least $(1 - \rho)$ fraction of training data we have

$$y_i \langle \nabla f_{\mathbf{W}^{(0)}}(\mathbf{x}_i), \mathbf{U}^* \rangle \geq m^{1/2} \gamma,$$

213 where γ is an absolute positive constant² and $\rho \in [0, 1)$.

214 The following corollary provides an upper bound of ϵ_{NTRF} under Assumption A.1 for some R .

Proposition A.2. Under Assumption A.1, for any $\epsilon, \delta > 0$, if $R \geq C[\log^{1/2}(n/\delta) + \log(1/\epsilon)]/\gamma$
 for some absolute constant C , then with probability at least $1 - \delta$,

$$\epsilon_{\text{NTRF}} := \inf_{F \in \mathcal{F}(\mathbf{W}^{(0)}, R)} n^{-1} \sum_{i=1}^n \ell(y_i F(\mathbf{x}_i)) \leq \epsilon + \rho \cdot \mathcal{O}(R).$$

215 Proposition A.2 covers the setting where the NTRF function class is allowed to misclassify training
 216 data, while most of existing work typically assumes that all training data can be perfectly separated
 217 with constant margin (i.e., $\rho = 0$) [15, 21]. Our results show that for sufficiently small misclassifica-
 218 tion ratio $\rho = \mathcal{O}(\epsilon)$, we have $\epsilon_{\text{NTRF}} = \tilde{\mathcal{O}}(\epsilon)$ by choosing the radius parameter R logarithmic in n ,
 219 δ^{-1} , and ϵ^{-1} . Substituting this result into Theorems 3.3, 3.4 and 3.5, it can be shown that a neural
 220 network with width $m = \text{poly}(L, \log(n/\delta), \log(1/\epsilon))$ suffices to guarantee good optimization and
 221 generalization performances for both GD and SGD.

222 **A.2 Data Separability by Shallow Neural Tangent Model**

223 In this subsection, we study the data separation assumption made in Ji and Telgarsky [15] and show
 224 that our results cover this particular setting. We first restate the assumption as follows.

225 **Assumption A.3.** There exists $\bar{\mathbf{u}}(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ and $\gamma \geq 0$ such that $\|\bar{\mathbf{u}}(\mathbf{z})\|_2 \leq 1$ for all $\mathbf{z} \in \mathbb{R}^d$, and

$$y_i \int_{\mathbb{R}^d} \sigma'(\langle \mathbf{z}, \mathbf{x}_i \rangle) \cdot \langle \bar{\mathbf{u}}(\mathbf{z}), \mathbf{x}_i \rangle d\mu_N(\mathbf{z}) \geq \gamma$$

226 for all $i \in [n]$, where $\mu_N(\cdot)$ denotes the standard normal distribution.

227 Assumption A.3 is related to the linear separability of the gradients of the first layer parameters at
 228 random initialization, where the randomness is replaced with an integral by taking the infinite width
 229 limit. Note that similar assumptions have also been studied in [9, 19, 13]. The assumption made in
 230 [9, 13] uses gradients with respect to the second layer weights instead of the first layer ones. In the
 231 following, we mainly focus on Assumption A.3, while our result can also be generalized to cover the
 232 setting in [9, 13].

²The factor $m^{1/2}$ is introduced here since $\|\nabla_{\mathbf{W}^{(0)}} f(\mathbf{x}_i)\|_F$ is typically of order $\mathcal{O}(m^{1/2})$.

233 In order to make a fair comparison, we reduce our results for multilayer networks to the two-layer
 234 setting. In this case, the neural network function takes form

$$f_{\mathbf{W}}(\mathbf{x}) = m^{1/2} \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}).$$

235 Then we provide the following proposition, which states that Assumption A.3 implies a certain choice
 236 of $R = \tilde{O}(1)$ such that the minimum training loss achieved by the function in the NTRF function
 237 class $\mathcal{F}(\mathbf{W}^{(0)}, R)$ satisfies $\epsilon_{\text{NTRF}} = O(\epsilon)$, where ϵ is the target error.

238 **Proposition A.4.** Suppose the training data satisfies Assumption A.3. For any $\epsilon, \delta > 0$, let $R =$
 239 $C[\log(n/\delta) + \log(1/\epsilon)]/\gamma$ for some large enough absolute constant C . If the neural network width
 240 satisfies $m = \Omega(\log(n/\delta)/\gamma^2)$, then with probability at least $1 - \delta$, there exist $F_{\mathbf{W}^{(0)}, \overline{\mathbf{W}}}(\mathbf{x}_i) \in$
 241 $\mathcal{F}(\mathbf{W}^{(0)}, R)$ such that $\ell(y_i \cdot F_{\mathbf{W}^{(0)}, \overline{\mathbf{W}}}(\mathbf{x}_i)) \leq \epsilon, \forall i \in [n]$.

242 Proposition A.4 shows that under Assumption A.3, there exists $F_{\mathbf{W}^{(0)}, \overline{\mathbf{W}}}(\cdot) \in \mathcal{F}(\mathbf{W}^{(0)}, R)$ with
 243 $R = \tilde{O}(1)$ such that the cross-entropy loss of $F_{\mathbf{W}^{(0)}, \overline{\mathbf{W}}}(\cdot)$ at each training data point is bounded by ϵ .
 244 This implies that $\epsilon_{\text{NTRF}} \leq \epsilon$. Moreover, by applying Theorem 3.3 with $L = 2$, the condition on the
 245 neural network width becomes $m = \Omega(\text{poly}(\log(n/\delta), \log(1/\epsilon)))$ ³, which matches the condition
 246 proved in Ji and Telgarsky [15].

247 A.3 Class-dependent Data Nondegeneration

248 In previous subsections, we have shown that under certain data separation conditions ϵ_{NTRF} can be
 249 sufficiently small while the corresponding NTRF function class has R of order $\tilde{O}(1)$. Thus neural
 250 networks with polylogarithmic width enjoy nice optimization and generalization guarantees. In this
 251 part, we consider the following much milder data separability assumption made in Zou et al. [23].

252 **Assumption A.5.** For all $i \neq i'$ if $y_i \neq y_{i'}$, then $\|\mathbf{x}_i - \mathbf{x}_{i'}\|_2 \geq \phi$ for some absolute constant ϕ .

253 In contrast to the conventional data nondegeneration assumption (i.e., no duplicate data points)
 254 made in Allen-Zhu et al. [2], Du et al. [12, 11], Zou and Gu [24]⁴, Assumption A.5 only requires
 255 that the data points from different classes are nondegenerate, thus we call it class-dependent data
 256 nondegeneration.

257 We have the following proposition which shows that Assumption A.5 also implies the existence of a
 258 good function that achieves ϵ training error, in the NTRF function class with a certain choice of R .

259 **Proposition A.6.** Under Assumption A.5, if

$$R = \Omega(n^{3/2} \phi^{-1/2} \log(n \delta^{-1} \epsilon^{-1})), \quad m = \tilde{\Omega}(L^{22} n^{12} \phi^{-4}),$$

260 we have $\epsilon_{\text{NTRF}} \leq \epsilon$ with probability at least $1 - \delta$.

261 Proposition A.6 suggests that under Assumption A.5, in order to guarantee $\epsilon_{\text{NTRF}} \leq \epsilon$, the size of
 262 NTRF function class needs to be $\Omega(n^{3/2})$. Plugging this into Theorems 3.4 and 3.5 leads to vacuous
 263 bounds on the test error. This makes sense since Assumption A.5 basically covers the ‘‘random label’’
 264 setting, which is impossible to be learned with small generalization error. Moreover, we would like
 265 to point out our theoretical analysis leads to a sharper over-parameterization condition than that
 266 proved in Zou et al. [23], i.e., $m = \tilde{\Omega}(n^{14} L^{16} \phi^{-4} + n^{12} L^{16} \phi^{-4} \epsilon^{-1})$, if the network depth satisfies
 267 $L \leq \tilde{O}(n^{1/3} \vee \epsilon^{-1/6})$.

³Similar to Ji and Telgarsky [15], the margin parameter is considered as a constant and thus does not appear in the condition on m .

⁴Specifically, Allen-Zhu et al. [2], Zou and Gu [24] require that any two data points (rather than data points from different classes) are separated by a positive distance. Zou and Gu [24] shows that this assumption is equivalent to those made in Du et al. [12, 11], which require that the composite kernel matrix is strictly positive definite.

268 **B Proof of Main Theorems**

269 In this section we provide the full proof of Theorems 3.3, 3.4 and 3.5.

270 **B.1 Proof of Theorem 3.3**

271 Here we introduce a key technical lemma used in the proof of Theorem 3.3.

272 Our proof is based on the key observation that near initialization, the neural network function can be
 273 approximated by its first-order Taylor expansion. In the following, we first give the definition of the
 274 linear approximation error in a τ -neighborhood around initialization.

$$\epsilon_{\text{app}}(\tau) := \sup_{i=1, \dots, n} \sup_{\mathbf{W}', \mathbf{W} \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)} |f_{\mathbf{W}'}(\mathbf{x}_i) - f_{\mathbf{W}}(\mathbf{x}_i) - \langle \nabla f_{\mathbf{W}}(\mathbf{x}_i), \mathbf{W}' - \mathbf{W} \rangle|.$$

275 If all the iterates of GD stay inside a neighborhood around initialization with small linear approxima-
 276 tion error, then we may expect that the training of neural networks should be similar to the training of
 277 the corresponding linear model, where standard optimization techniques can be applied. Motivated
 278 by this, we also give the following definition on the gradient upper bound of neural networks around
 279 initialization, which is related to the Lipschitz constant of the optimization objective function.

$$M(\tau) := \sup_{i=1, \dots, n} \sup_{l=1, \dots, L} \sup_{\mathbf{W} \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)} \|\nabla_{\mathbf{w}_l} f_{\mathbf{W}}(\mathbf{x}_i)\|_F.$$

280 By definition, we can choose $\mathbf{W}^* \in \mathcal{B}(\mathbf{W}^{(0)}, Rm^{-1/2})$ such that $n^{-1} \sum_{i=1}^n \ell(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)) =$
 281 ϵ_{NTRF} . Then we have the following lemma.

282 **Lemma B.1.** Set $\eta = \mathcal{O}(L^{-1}M(\tau)^{-2})$. Suppose that $\mathbf{W}^* \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)$ and $\mathbf{W}^{(t)} \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)$
 283 for all $0 \leq t \leq t' - 1$. Then it holds that

$$\frac{1}{t'} \sum_{t=0}^{t'-1} L_{\mathcal{S}}(\mathbf{W}^{(t)}) \leq \frac{\|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F^2 - \|\mathbf{W}^{(t')} - \mathbf{W}^*\|_F^2 + 2t'\eta\epsilon_{\text{NTRF}}}{t'\eta(\frac{3}{2} - 4\epsilon_{\text{app}}(\tau))}.$$

284 Lemma B.1 plays a central role in our proof. In specific, if $\mathbf{W}^{(t)} \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)$ for all $t \leq t'$, then
 285 Lemma B.1 implies that the average training loss is in the same order of ϵ_{NTRF} as long as the linear
 286 approximation error $\epsilon_{\text{app}}(\tau)$ is bounded by a positive constant. This is in contrast to the proof in
 287 Cao and Gu [8], where $\epsilon_{\text{app}}(\tau)$ appears as an additive term in the upper bound of the training loss,
 288 thus requiring $\epsilon_{\text{app}}(\tau) = \mathcal{O}(\epsilon_{\text{NTRF}})$ to achieve the same error bound as in Lemma B.1. Since we can
 289 show that $\epsilon_{\text{app}} = \tilde{\mathcal{O}}(m^{-1/6})$ (See Section B.1), this suggests that $m = \tilde{\Omega}(1)$ is sufficient to make the
 290 average training loss in the same order of ϵ_{NTRF} .

291 Compared with the recent results for two-layer networks by [15], Lemma B.1 is proved with different
 292 techniques. In specific, the proof by [15] relies on the 1-homogeneous property of the ReLU activation
 293 function, which limits their analysis to two-layer networks with fixed second layer weights. In
 294 comparison, our proof does not rely on homogeneity, and is purely based on the linear approximation
 295 property of neural networks and some specific properties of the loss function. Therefore, our proof
 296 technique can handle deep networks, and is potentially applicable to non-ReLU activation functions
 297 and other network architectures (e.g. Convolutional neural networks and Residual networks).

298 We provide the following lemma which is useful in the subsequent proof.

299 **Lemma B.2** (Lemmas 4.1 and B.3 in Cao and Gu [8]). There exists an absolute constant κ such that,
 300 with probability at least $1 - \mathcal{O}(nL^2) \exp[-\Omega(m\tau^{2/3}L)]$, for any $\tau \leq \kappa L^{-6}[\log(m)]^{-3/2}$, it holds
 301 that

$$\epsilon_{\text{app}}(\tau) \leq \tilde{\mathcal{O}}(\tau^{4/3}L^3m^{1/2}), \quad M(\tau) \leq \tilde{\mathcal{O}}(\sqrt{m}).$$

302 Now we provide the detailed proof which consists of two steps: (i) showing that all T iterates stay
 303 close to initialization, and (ii) bounding the empirical loss achieved by gradient descent. Both of
 304 these steps are proved based on Lemma B.1.

305 *Proof of Theorem 3.3.* Recall that \mathbf{W}^* is chosen such that

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)) = \epsilon_{\text{NTRF}}$$

306 and $\mathbf{W}^* \in \mathcal{B}(\mathbf{W}^{(0)}, Rm^{-1/2})$. Note that to apply Lemma B.1, we need the region $\mathcal{B}(\mathbf{W}^{(0)}, \tau)$ to in-
 307 clude both \mathbf{W}^* and $\{\mathbf{W}^{(t)}\}_{t=0, \dots, t'}$. This motivates us to set $\tau = \tilde{\mathcal{O}}(L^{1/2}m^{-1/2}R)$, which is slightly
 308 larger than $m^{-1/2}R$. With this choice of τ , by Lemma B.2 we have $\epsilon_{\text{app}}(\tau) = \tilde{\mathcal{O}}(\tau^{4/3}m^{1/2}L^3) =$
 309 $\tilde{\mathcal{O}}(R^{4/3}L^{11/3}m^{-1/6})$. Therefore, we can set

$$m = \tilde{\Omega}(R^8 L^{22}) \quad (\text{B.1})$$

310 to ensure that $\epsilon_{\text{app}}(\tau) \leq 1/8$, where $\tilde{\Omega}(\cdot)$ hides polylogarithmic dependencies on network depth L ,
 311 NTRF function class size R , and failure probability parameter δ . Then by Lemma B.1, we have with
 312 probability at least $1 - \delta$, we have

$$\|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F^2 - \|\mathbf{W}^{(t')} - \mathbf{W}^*\|_F^2 \geq \eta \sum_{t=0}^{t'-1} L_S(\mathbf{W}^{(t)}) - 2t'\eta\epsilon_{\text{NTRF}} \quad (\text{B.2})$$

313 as long as $\mathbf{W}^{(0)}, \dots, \mathbf{W}^{(t'-1)} \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)$. In the following proof we choose $\eta = \Theta(L^{-1}m^{-1})$
 314 and $T = \lceil LR^2 m^{-1} \eta^{-1} \epsilon_{\text{NTRF}}^{-1} \rceil$.

315 We prove the theorem by two steps: 1) we show that all iterates $\{\mathbf{W}^{(0)}, \dots, \mathbf{W}^{(T)}\}$ will stay inside
 316 the region $\mathcal{B}(\mathbf{W}^{(0)}, \tau)$; and 2) we show that GD can find a neural network with at most $3\epsilon_{\text{NTRF}}$
 317 training loss within T iterations.

318 **All iterates stay inside $\mathcal{B}(\mathbf{W}^{(0)}, \tau)$.** We prove this part by induction. Specifically, given $t' \leq T$, we
 319 assume the hypothesis $\mathbf{W}^{(t)} \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)$ holds for all $t < t'$ and prove that $\mathbf{W}^{(t')} \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)$.
 320 First, it is clear that $\mathbf{W}^{(0)} \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)$. Then by (B.2) and the fact that $L_S(\mathbf{W}) \geq 0$, we have

$$\|\mathbf{W}^{(t')} - \mathbf{W}^*\|_F^2 \leq \|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F^2 + 2\eta t' \epsilon_{\text{NTRF}}$$

321 Note that $T = \lceil LR^2 m^{-1} \eta^{-1} \epsilon_{\text{NTRF}}^{-1} \rceil$ and $\mathbf{W}^* \in \mathcal{B}(\mathbf{W}^{(0)}, R \cdot m^{-1/2})$, we have

$$\sum_{l=1}^L \|\mathbf{W}_l^{(t')} - \mathbf{W}_l^*\|_F^2 = \|\mathbf{W}^{(t')} - \mathbf{W}^*\|_F^2 \leq CLR^2 m^{-1},$$

322 where $C \geq 4$ is an absolute constant. Therefore, by triangle inequality, we further have the following
 323 for all $l \in [L]$,

$$\begin{aligned} \|\mathbf{W}_l^{(t')} - \mathbf{W}_l^{(0)}\|_F &\leq \|\mathbf{W}_l^{(t')} - \mathbf{W}_l^*\|_F + \|\mathbf{W}_l^{(0)} - \mathbf{W}_l^*\|_F \\ &\leq \sqrt{CLR} m^{-1/2} + R m^{-1/2} \\ &\leq 2\sqrt{CLR} m^{-1/2}. \end{aligned} \quad (\text{B.3})$$

324 Therefore, it is clear that $\|\mathbf{W}_l^{(t')} - \mathbf{W}_l^{(0)}\|_F \leq 2\sqrt{CLR} m^{-1/2} \leq \tau$ based on our choice of τ
 325 previously. This completes the proof of the first part.

326 **Convergence of gradient descent.** (B.2) implies

$$\|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F^2 - \|\mathbf{W}^{(T)} - \mathbf{W}^*\|_F^2 \geq \eta \left(\sum_{t=0}^{T-1} L_S(\mathbf{W}^{(t)}) - 2T\epsilon_{\text{NTRF}} \right).$$

327 Dividing by ηT on the both sides, we get

$$\frac{1}{T} \sum_{t=0}^{T-1} L_S(\mathbf{W}^{(t)}) \leq \frac{\|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F^2}{\eta T} + 2\epsilon_{\text{NTRF}} \leq \frac{LR^2 m^{-1}}{\eta T} + 2\epsilon_{\text{NTRF}} \leq 3\epsilon_{\text{NTRF}},$$

328 where the second inequality is by the fact that $\mathbf{W}^* \in \mathcal{B}(\mathbf{W}^{(0)}, R \cdot m^{-1/2})$ and the last in-
 329 equality is by our choices of T and η which ensure that $T\eta \geq LR^2 m^{-1} \epsilon_{\text{NTRF}}^{-1}$. Notice that
 330 $T = \lceil LR^2 m^{-1} \eta^{-1} \epsilon_{\text{NTRF}}^{-1} \rceil = \mathcal{O}(L^2 R^2 \epsilon_{\text{NTRF}}^{-1})$. This completes the proof of the second part, and
 331 we are able to complete the proof. \square

332 **B.2 Proof of Theorem 3.4**

333 Following Cao and Gu [9], we first introduce the definition of surrogate loss of the network, which is
 334 defined by the derivative of the loss function.

335 **Definition B.3.** We define the empirical surrogate error $\mathcal{E}_S(\mathbf{W})$ and population surrogate error
 336 $\mathcal{E}_D(\mathbf{W})$ as follows:

$$\mathcal{E}_S(\mathbf{W}) := -\frac{1}{n} \sum_{i=1}^n \ell'[y_i \cdot f_{\mathbf{W}}(\mathbf{x}_i)], \quad \mathcal{E}_D(\mathbf{W}) := \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \{-\ell'[y \cdot f_{\mathbf{W}}(\mathbf{x})]\}.$$

337 The following lemma gives uniform-convergence type of results for $\mathcal{E}_S(\mathbf{W})$ utilizing the fact that
 338 $-\ell'(\cdot)$ is bounded and Lipschitz continuous.

339 **Lemma B.4.** For any $\tilde{R}, \delta > 0$, suppose that $m = \tilde{\Omega}(L^{12}\tilde{R}^2) \cdot [\log(1/\delta)]^{3/2}$. Then with probability
 340 at least $1 - \delta$, it holds that

$$|\mathcal{E}_D(\mathbf{W}) - \mathcal{E}_S(\mathbf{W})| \leq \tilde{\mathcal{O}}\left(\min\left\{4^L L^{3/2} \tilde{R} \sqrt{\frac{m}{n}}, \frac{L\tilde{R}}{\sqrt{n}} + \frac{L^3 \tilde{R}^{4/3}}{m^{1/6}}\right\}\right) + \mathcal{O}\left(\sqrt{\frac{\log(1/\delta)}{n}}\right)$$

341 for all $\mathbf{W} \in \mathcal{B}(\mathbf{W}^{(0)}, \tilde{R} \cdot m^{-1/2})$

342 We are now ready to prove Theorem 3.4, which combines the trajectory distance analysis in the proof
 343 of Theorem 3.3 with Lemma B.4.

344 *Proof of Theorem 3.4.* With exactly the same proof as Theorem 3.3, by (B.3) and induction we have
 345 $\mathbf{W}^{(0)}, \mathbf{W}^{(1)}, \dots, \mathbf{W}^{(T)} \in \mathcal{B}(\mathbf{W}^{(0)}, \tilde{R}m^{-1/2})$ with $\tilde{R} = \mathcal{O}(\sqrt{LR})$. Therefore by Lemma B.4, we
 346 have

$$|\mathcal{E}_D(\mathbf{W}^{(t)}) - \mathcal{E}_S(\mathbf{W}^{(t)})| \leq \tilde{\mathcal{O}}\left(\min\left\{4^L L^2 R \sqrt{\frac{m}{n}}, \frac{L^{3/2} R}{\sqrt{n}} + \frac{L^{11/3} R^{4/3}}{m^{1/6}}\right\}\right) + \mathcal{O}\left(\sqrt{\frac{\log(1/\delta)}{n}}\right)$$

347 for all $t = 0, 1, \dots, T$. Note that we have $\mathbb{1}\{z < 0\} \leq -2\ell'(z)$. Therefore,

$$\begin{aligned} \mathbb{E}L_{\mathcal{D}}^{0-1}(\mathbf{W}^{(t)}) &\leq 2\mathcal{E}_D(\mathbf{W}^{(t)}) \\ &\leq 2L_S(\mathbf{W}^{(t)}) + \tilde{\mathcal{O}}\left(\min\left\{4^L L^2 R \sqrt{\frac{m}{n}}, \frac{L^{3/2} R}{\sqrt{n}} + \frac{L^{11/3} R^{4/3}}{m^{1/6}}\right\}\right) + \mathcal{O}\left(\sqrt{\frac{\log(1/\delta)}{n}}\right) \end{aligned}$$

348 $t = 0, 1, \dots, T$. This finishes the proof. \square

349 **B.3 Proof of Theorem 3.5**

350 In this section we provide the full proof of Theorem 3.5. We first give the following result, which is
 351 the counterpart of Lemma B.1 for SGD. Again we pick $\mathbf{W}^* \in \mathcal{B}(\mathbf{W}^{(0)}, Rm^{-1/2})$ such that the loss
 352 of the corresponding NTRF model $F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x})$ achieves ϵ_{NTRF} .

353 **Lemma B.5.** Set $\eta = \mathcal{O}(L^{-1}M(\tau)^{-2})$. Suppose that $\mathbf{W}^* \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)$ and $\mathbf{W}^{(n')} \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)$
 354 for all $0 \leq n' \leq n - 1$. Then it holds that

$$\|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F^2 - \|\mathbf{W}^{(n')} - \mathbf{W}^*\|_F^2 \geq \left(\frac{3}{2} - 4\epsilon_{\text{app}}(\tau)\right)\eta \sum_{i=1}^{n'} L_i(\mathbf{W}^{(i-1)}) - 2n\eta\epsilon_{\text{NTRF}}.$$

355 We introduce a surrogate loss $\mathcal{E}_i(\mathbf{W}) = -\ell'[y_i \cdot f_{\mathbf{W}}(\mathbf{x}_i)]$ and its population version $\mathcal{E}_D(\mathbf{W}) =$
 356 $\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[-\ell'[y \cdot f_{\mathbf{W}}(\mathbf{x})]]$, which have been used in [14, 8, 15]. Our proof is based on the application
 357 of Lemma B.5 and an online-to-batch conversion argument [10, 8, 15]. We introduce a surrogate
 358 loss $\mathcal{E}_i(\mathbf{W}) = -\ell'[y_i \cdot f_{\mathbf{W}}(\mathbf{x}_i)]$ and its population version $\mathcal{E}_D(\mathbf{W}) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[-\ell'(y \cdot f_{\mathbf{W}}(\mathbf{x}))]$,
 359 which have been used in [14, 8, 19, 15].

360 *Proof of Theorem 3.5.* Recall that \mathbf{W}^* is chosen such that

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)) = \epsilon_{\text{NTRF}}$$

361 and $\mathbf{W}^* \in \mathcal{B}(\mathbf{W}^{(0)}, Rm^{-1/2})$. To apply Lemma B.5, we need the region $\mathcal{B}(\mathbf{W}^{(0)}, \tau)$ to include
 362 both \mathbf{W}^* and $\{\mathbf{W}^{(t)}\}_{t=0, \dots, t'}$. This motivates us to set $\tau = \tilde{\mathcal{O}}(L^{1/2}m^{-1/2}R)$, which is slightly
 363 larger than $m^{-1/2}R$. With this choice of τ , by Lemma B.2 we have $\epsilon_{\text{app}}(\tau) = \tilde{\mathcal{O}}(\tau^{4/3}m^{1/2}L^3) =$
 364 $\tilde{\mathcal{O}}(R^{4/3}L^{11/3}m^{-1/6})$. Therefore, we can set

$$m = \tilde{\Omega}(R^8L^{22})$$

365 to ensure that $\epsilon_{\text{app}}(\tau) \leq 1/8$, where $\tilde{\Omega}(\cdot)$ hides polylogarithmic dependencies on network depth L ,
 366 NTRF function class size R , and failure probability parameter δ .

367 Then by Lemma B.5, we have with probability at least $1 - \delta$,

$$\|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F^2 - \|\mathbf{W}^{(n')} - \mathbf{W}^*\|_F^2 \geq \eta \sum_{i=1}^{n'} L_i(\mathbf{W}^{(i-1)}) - 2n\eta\epsilon_{\text{NTRF}} \quad (\text{B.4})$$

368 as long as $\mathbf{W}^{(0)}, \dots, \mathbf{W}^{(n'-1)} \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)$.

369 We then prove Theorem 3.5 in two steps: 1) all iterates stay inside $\mathcal{B}(\mathbf{W}^{(0)}, \tau)$; and 2) convergence
 370 of online SGD.

371 **All iterates stay inside $\mathcal{B}(\mathbf{W}^{(0)}, \tau)$.** Similar to the proof of Theorem 3.3, we prove this part by
 372 induction. Assuming $\mathbf{W}^{(i)}$ satisfies $\mathbf{W}^{(i)} \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)$ for all $i \leq n' - 1$, by (B.4), we have

$$\begin{aligned} \|\mathbf{W}^{(n')} - \mathbf{W}^*\|_F^2 &\leq \|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F^2 + 2n\eta\epsilon_{\text{NTRF}} \\ &\leq LR^2 \cdot m^{-1} + 2n\eta\epsilon_{\text{NTRF}}, \end{aligned}$$

373 where the last inequality is by $\mathbf{W}^* \in \mathcal{B}(\mathbf{W}^{(0)}, Rm^{-1/2})$. Then by triangle inequality, we further get

$$\begin{aligned} \|\mathbf{W}_l^{(n')} - \mathbf{W}_l^{(0)}\|_F &\leq \|\mathbf{W}_l^{(n')} - \mathbf{W}_l^*\|_F + \|\mathbf{W}_l^* - \mathbf{W}_l^{(0)}\|_F \\ &\leq \|\mathbf{W}^{(n')} - \mathbf{W}^*\|_F + \|\mathbf{W}_l^* - \mathbf{W}_l^{(0)}\|_F \\ &\leq \mathcal{O}(\sqrt{LR}m^{-1/2} + \sqrt{n\eta\epsilon_{\text{NTRF}}}). \end{aligned}$$

374 Then by our choices of $\eta = \Theta(m^{-1} \cdot (LR^2n^{-1}\epsilon_{\text{NTRF}}^{-1} \wedge L^{-1}))$, we have $\|\mathbf{W}^{(n')} - \mathbf{W}^{(0)}\|_F \leq$
 375 $2\sqrt{LR}m^{-1/2} \leq \tau$. This completes the proof of the first part.

376 **Convergence of online SGD.** By (B.4), we have

$$\|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F^2 - \|\mathbf{W}^{(n)} - \mathbf{W}^*\|_F^2 \geq \eta \left(\sum_{i=1}^n L_i(\mathbf{W}^{(i-1)}) - 2n\epsilon_{\text{NTRF}} \right).$$

377 Dividing by ηn on the both sides and rearranging terms, we get

$$\frac{1}{n} \sum_{i=1}^n L_i(\mathbf{W}^{(i-1)}) \leq \frac{\|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F^2 - \|\mathbf{W}^{(n)} - \mathbf{W}^*\|_F^2}{\eta n} + 2\epsilon_{\text{NTRF}} \leq \frac{L^2R^2}{n} + 3\epsilon_{\text{NTRF}},$$

378 where the second inequality follows from facts that $\mathbf{W}^* \in \mathcal{B}(\mathbf{W}^{(0)}, R \cdot m^{-1/2})$ and $\eta = \Theta(m^{-1} \cdot$
 379 $(LR^2n^{-1}\epsilon_{\text{NTRF}}^{-1} \wedge L^{-1}))$. By Lemma 4.3 in [15] and the fact that $\mathcal{E}_i(\mathbf{W}^{(i-1)}) \leq L_i(\mathbf{W}^{(i-1)})$, we
 380 have

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n L_D^{0-1}(\mathbf{W}^{(i-1)}) &\leq \frac{2}{n} \sum_{i=1}^n \mathcal{E}_D(\mathbf{W}^{(i-1)}) \\ &\leq \frac{8}{n} \sum_{i=1}^n \mathcal{E}_i(\mathbf{W}^{(i-1)}) + \frac{8 \log(1/\delta)}{n} \\ &\leq \frac{8L^2R^2}{n} + \frac{8 \log(1/\delta)}{n} + 24\epsilon_{\text{NTRF}}. \end{aligned}$$

381 This completes the proof of the second part. \square

382 **C Proof of Results in Section A**

383 **C.1 Proof of Proposition A.2**

384 We first provide the following lemma which gives an upper bound of the neural network output at the
385 initialization.

386 **Lemma C.1** (Lemma 4.4 in Cao and Gu [8]). Under Assumptions 3.1, if $m \geq \bar{C}L \log(nL/\delta)$ with
387 some absolute constant \bar{C} , with probability at least $1 - \delta$, we have

$$|f_{\mathbf{W}^{(0)}}(\mathbf{x}_i)| \leq C\sqrt{\log(n/\delta)}$$

388 for some absolute constant C .

389 *Proof of Proposition A.2.* Under assumption A.1, we can find a collection of matrices $\mathbf{U}^* =$
390 $\{\mathbf{U}_1^*, \dots, \mathbf{U}_L^*\}$ with $\sum_{i=1}^L \|\mathbf{U}_i^*\|_F^2 = 1$ such that $y_i \langle \nabla f_{\mathbf{W}^{(0)}}(\mathbf{x}_i), \mathbf{U}^* \rangle \geq m^{1/2}\gamma$ for at least $1 - \sigma$
391 fraction of training data. By Lemma C.1, for all $i \in [n]$ we have $|f_{\mathbf{W}^{(0)}}(\mathbf{x}_i)| \leq C\sqrt{\log(n/\delta)}$ for
392 some absolute constant C . Then for any positive constant λ , we have for at least $1 - \sigma$ portion of
393 data,

$$y_i (f_{\mathbf{W}^{(0)}}(\mathbf{x}_i) + \langle \nabla f_{\mathbf{W}^{(0)}}, \lambda \mathbf{U}^* \rangle) \geq m^{1/2}\lambda\gamma - C\sqrt{\log(n/\delta)}.$$

394 For this fraction of data, we can set

$$\lambda = \frac{C' [\log^{1/2}(n/\delta) + \log(1/\epsilon)]}{m^{1/2}\gamma},$$

395 where C' is an absolute constant, and get

$$m^{1/2}\lambda\gamma - C\sqrt{\log(n/\delta)} \geq \log(1/\epsilon).$$

396 Now we let $\mathbf{W}^* = \mathbf{W}^{(0)} + \lambda \mathbf{U}^*$. By the choice of R in Proposition A.2, we have $\mathbf{W}^* \in$
397 $\mathcal{B}(\mathbf{W}^{(0)}, R \cdot m^{-1/2})$. The above inequality implies that for this at least $1 - \sigma$ fraction of data,
398 we have $\ell(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)) \leq \epsilon$. For the rest data, we have

$$y_i (f_{\mathbf{W}^{(0)}}(\mathbf{x}_i) + \langle \nabla f_{\mathbf{W}^{(0)}}, \lambda \mathbf{U}^* \rangle) \geq -C\sqrt{\log(n/\delta)} - \lambda \|\nabla f_{\mathbf{W}^{(0)}}\|_2^2 \geq -C_1 R$$

399 for some absolute positive constant C_1 , where the last inequality follows from fact that $\|\nabla f_{\mathbf{W}^{(0)}}\|_2 =$
400 $\tilde{\mathcal{O}}(m^{1/2})$ (see Lemma B.2 for detail). Then note that we use cross-entropy loss, it follows that for this
401 fraction of training data, we have $\ell(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)) \leq C_2 R$ for some constant C_2 . Combining the
402 results of these two fractions of training data, we can conclude

$$\epsilon_{\text{NTRF}} \leq n^{-1} \sum_{i=1}^n \ell(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)) \leq (1 - \sigma)\epsilon + \rho \cdot \mathcal{O}(R)$$

403 This completes the proof.

404

□

405 **C.2 Proof of Proposition A.4**

406 *Proof of Proposition A.4.* We are going to prove that Assumption A.3 implies the existence of a good
407 function in the NTRF function class.

408 By Definition 3.2 and the definition of cross-entropy loss, our goal is to prove that there exists
409 a collection of matrices $\bar{\mathbf{W}} = \{\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2\}$ satisfying $\max\{\|\bar{\mathbf{W}}_1 - \mathbf{W}_1^{(0)}\|_F, \|\bar{\mathbf{W}}_2 - \mathbf{W}_2^{(0)}\|_2\} \leq$
410 $R \cdot m^{-1/2}$ such that

$$y_i \cdot [f_{\mathbf{W}^{(0)}}(\mathbf{x}_i) + \langle \nabla_{\mathbf{w}_1} f_{\mathbf{W}^{(0)}}, \bar{\mathbf{W}}_1 - \mathbf{W}_1^{(0)} \rangle + \langle \nabla_{\mathbf{w}_2} f_{\mathbf{W}^{(0)}}, \bar{\mathbf{W}}_2 - \mathbf{W}_2^{(0)} \rangle] \geq \log(2/\epsilon).$$

411 We first consider $\nabla_{\mathbf{w}_1} f_{\mathbf{W}^{(0)}}(\mathbf{x}_i)$, which has the form

$$(\nabla_{\mathbf{w}_1} f_{\mathbf{W}^{(0)}}(\mathbf{x}_i))_j = m^{1/2} \cdot w_{2,j}^{(0)} \cdot \sigma'(\langle \mathbf{w}_{1,j}^{(0)}, \mathbf{x}_i \rangle) \cdot \mathbf{x}_i.$$

412 Note that $w_{2,j}^{(0)}$ and $\mathbf{w}_{1,j}^{(0)}$ are independently generated from $\mathcal{N}(0, 1/m)$ and $\mathcal{N}(0, 2\mathbf{I}/m)$ respectively,
413 thus we have $\mathbb{P}(|w_{2,j}^{(0)}| \geq 0.47m^{-1/2}) \geq 1/2$. By Hoeffding's inequality, we know that with
414 probability at least $1 - \exp(-m/8)$, there are at least $m/4$ nodes, whose union is denoted by \mathcal{S} ,
415 satisfying $|w_{2,j}^{(0)}| \geq 0.47m^{-1/2}$. Then we only focus on the nodes in the set \mathcal{S} . Note that $\mathbf{W}_1^{(0)}$ and
416 $\mathbf{W}_2^{(0)}$ are independently generated. Then by Assumption A.3 and Hoeffding's inequality, there exists
417 a function $\bar{\mathbf{u}}(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that with probability at least $1 - \delta'$,

$$\frac{1}{|\mathcal{S}|} \sum_{j \in \mathcal{S}} y_i \cdot \langle \bar{\mathbf{u}}(\mathbf{w}_{1,j}^{(0)}), \mathbf{x}_i \rangle \cdot \sigma'(\langle \mathbf{w}_{1,j}^{(0)}, \mathbf{x}_i \rangle) \geq \gamma - \sqrt{\frac{2 \log(1/\delta')}{|\mathcal{S}|}}.$$

418 Define $\mathbf{v}_j = \bar{\mathbf{u}}(\mathbf{w}_{1,j}^{(0)})/w_{2,j}$ if $|w_{2,j}| \geq 0.47m^{-1/2}$ and $\mathbf{v}_j = \mathbf{0}$ otherwise. Then we have

$$\begin{aligned} \sum_{j=1}^m y_i \cdot w_{2,j}^{(0)} \cdot \langle \mathbf{v}_j, \mathbf{x}_i \rangle \cdot \sigma'(\langle \mathbf{w}_{1,j}^{(0)}, \mathbf{x}_i \rangle) &= \sum_{j \in \mathcal{S}} y_i \cdot \langle \bar{\mathbf{u}}(\mathbf{w}_{1,j}^{(0)}), \mathbf{x}_i \rangle \cdot \sigma'(\langle \mathbf{w}_{1,j}^{(0)}, \mathbf{x}_i \rangle) \\ &\geq |\mathcal{S}| \gamma - \sqrt{2|\mathcal{S}| \log(1/\delta')}. \end{aligned}$$

419 Set $\delta = 2n\delta'$ and apply union bound, we have with probability at least $1 - \delta/2$,

$$\sum_{j=1}^m y_i \cdot w_{2,j}^{(0)} \cdot \langle \mathbf{v}_j, \mathbf{x}_i \rangle \cdot \sigma'(\langle \mathbf{w}_{1,j}^{(0)}, \mathbf{x}_i \rangle) \geq |\mathcal{S}| \gamma - \sqrt{2|\mathcal{S}| \log(2n/\delta)}.$$

420 Therefore, note that with probability at least $1 - \exp(-m/8)$, we have $|\mathcal{S}| \geq m/4$. Moreover, in
421 Assumption A.3, by $y_i \in \{\pm 1\}$ and $|\sigma'(\cdot)|, \|\bar{\mathbf{u}}(\cdot)\|_2, \|\mathbf{x}_i\|_2 \leq 1$ for $i = 1, \dots, n$, we see that $\gamma \leq 1$.
422 Then if $m \geq 32 \log(n/\delta)/\gamma^2$, with probability at least $1 - \delta/2 - \exp(-4 \log(n/\delta)/\gamma^2) \geq 1 - \delta$,

$$\sum_{j=1}^m y_i \cdot w_{2,j}^{(0)} \cdot \langle \mathbf{v}_j, \mathbf{x}_i \rangle \cdot \sigma'(\langle \mathbf{w}_{1,j}^{(0)}, \mathbf{x}_i \rangle) \geq |\mathcal{S}| \gamma / 2.$$

423 Let $\mathbf{U} = (\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m)^\top / \sqrt{m|\mathcal{S}|}$, we have

$$y_i \langle \nabla_{\mathbf{w}_1} f_{\mathbf{W}^{(0)}}(\mathbf{x}_i), \mathbf{U} \rangle = \frac{1}{\sqrt{|\mathcal{S}|}} \sum_{j=1}^m y_i \cdot w_{2,j}^{(0)} \cdot \langle \mathbf{v}_j, \mathbf{x}_i \rangle \cdot \sigma'(\langle \mathbf{w}_{1,j}^{(0)}, \mathbf{x}_i \rangle) \geq \frac{\sqrt{|\mathcal{S}|} \gamma}{2} \geq \frac{m^{1/2} \gamma}{4},$$

424 where the last inequality is by the fact that $|\mathcal{S}| \geq m/4$. Besides, note that by concentration and
425 Gaussian tail bound, we have $|f_{\mathbf{W}^{(0)}}(\mathbf{x}_i)| \leq C \log(n/\delta)$ for some absolute constant C . Therefore,
426 let $\bar{\mathbf{W}}_1 = \mathbf{W}_1^{(0)} + 4(\log(2/\epsilon) + C \log(n/\delta))m^{-1/2}\mathbf{U}/\gamma$ and $\bar{\mathbf{W}}_2 = \mathbf{W}_2^{(0)}$, we have

$$y_i \cdot [f_{\mathbf{W}^{(0)}}(\mathbf{x}_i) + \langle \nabla_{\mathbf{w}_1} f_{\mathbf{W}^{(0)}}, \bar{\mathbf{W}}_1 - \mathbf{W}_1^{(0)} \rangle + \langle \nabla_{\mathbf{w}_2} f_{\mathbf{W}^{(0)}}, \bar{\mathbf{W}}_2 - \mathbf{W}_2^{(0)} \rangle] \geq \log(2/\epsilon). \quad (\text{C.1})$$

427 Note that $\|\bar{\mathbf{u}}(\cdot)\|_2 \leq 1$, we have $\|\mathbf{U}\|_F \leq 1/0.47 \leq 2.2$. Therefore, we further have $\|\bar{\mathbf{W}}_1 -$
428 $\mathbf{W}_1^{(0)}\|_F \leq 8.8\gamma^{-1}(\log(2/\epsilon) + C \log(n/\delta)) \cdot m^{-1/2}$. This implies that $\bar{\mathbf{W}} \in \mathcal{B}(\mathbf{W}^{(0)}, R)$ with
429 $R = \mathcal{O}(\log(n/(\delta\epsilon))/\gamma)$. Applying the inequality $\ell(\log(2/\epsilon)) \leq \epsilon$ on (C.1) gives

$$\ell(y_i \cdot F_{\mathbf{W}^{(0)}, \bar{\mathbf{W}}}(\mathbf{x}_i)) \leq \epsilon$$

430 for all $i = 1, \dots, n$. This completes the proof. \square

431 C.3 Proof of Proposition A.6

432 Based on our theoretical analysis, the major goal is to show that there exist certain choices of R
433 and m such that the best NTRF model in the function class $\mathcal{F}(\mathbf{W}^{(0)}, R)$ can achieve ϵ training
434 error. In this proof, we will prove a stronger results by showing that given the quantities of R
435 and m specified in Proposition A.6, there exists a NTRF model with parameter \mathbf{W}^* that satisfies
436 $n^{-1} \sum_{i=1}^n \ell(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)) \leq \epsilon$.

437 In order to do so, we consider training the NTRF model via a different surrogate loss function.
438 Specifically, we consider squared hinge loss $\tilde{\ell}(x) = (\max\{\lambda - x, 0\})^2$, where λ denotes the target

margin. In the later proof, we choose $\lambda = \log(1/\epsilon) + 1$ such that the condition $\tilde{\ell}(x) \leq 1$ can guarantee that $x \geq \log(\epsilon)$. Moreover, we consider using gradient flow, i.e., gradient descent with infinitesimal step size, to train the NTRF model. Therefore, in the remaining part of the proof, we consider optimizing the NTRF parameter \mathbf{W} with the loss function

$$\tilde{L}_S(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \tilde{\ell}(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}}(\mathbf{x}_i)).$$

Moreover, for simplicity, we only consider optimizing parameter in the last hidden layer (i.e., \mathbf{W}_{L-1}). Then the gradient flow can be formulated as

$$\frac{d\mathbf{W}_{L-1}(t)}{dt} = -\nabla_{\mathbf{W}_{L-1}} \tilde{L}_S(\mathbf{W}(t)), \quad \frac{d\mathbf{W}_l(t)}{dt} = \mathbf{0} \quad \text{for any } l \neq L-1.$$

Note that the NTRF model is a linear model, thus by Definition 3.2, we have

$$\begin{aligned} \nabla_{\mathbf{W}_{L-1}} \tilde{L}_S(\mathbf{W}(t)) &= y_i \tilde{\ell}'(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}(t)}(\mathbf{x}_i)) \cdot \nabla_{\mathbf{W}_{L-1}} F_{\mathbf{W}^{(0)}, \mathbf{W}(t)}(\mathbf{x}_i) \\ &= y_i \tilde{\ell}'(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}(t)}(\mathbf{x}_i)) \cdot \nabla_{\mathbf{W}_{L-1}^{(0)}} f_{\mathbf{W}^{(0)}}(\mathbf{x}_i). \end{aligned} \quad (\text{C.2})$$

Then it is clear that $\nabla_{\mathbf{W}_{L-1}} \tilde{L}_S(\mathbf{W}(t))$ has fixed direction throughout the optimization.

In order to prove the convergence of gradient flow and characterize the quantity of R , We first provide the following lemma which gives an upper bound of the NTRF model output at the initialization.

Then we provide the following lemma which characterizes a lower bound of the Frobenius norm of the partial gradient $\nabla_{\mathbf{W}_{L-1}} \tilde{L}_S(\mathbf{W})$.

Lemma C.2 (Lemma B.5 in Zou et al. [23]). Under Assumptions 3.1 and A.5, if $m = \tilde{\Omega}(n^2 \phi^{-1})$, then for all $t \geq 0$, with probability at least $1 - \exp(-O(m\phi/n))$, there exist a positive constant C such that

$$\|\nabla_{\mathbf{W}_{L-1}} \tilde{L}_S(\mathbf{W}(t))\|_F^2 \geq \frac{Cm\phi}{n^5} \left[\sum_{i=1}^n \tilde{\ell}'(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}(t)}(\mathbf{x}_i)) \right]^2.$$

We slightly modified the original version of this lemma since we use different models (we consider NTRF model while Zou et al. [23] considers neural network model). However, by (C.2), it is clear that the gradient $\nabla \tilde{L}_S(\mathbf{W})$ can be regarded as a type of the gradient for neural network model at the initialization (i.e., $\nabla_{\mathbf{W}_{L-1}} L_S(\mathbf{W}^{(0)})$) is valid. Now we are ready to present the proof.

Proof of Proposition A.6. Recall that we only consider training the last hidden weights, i.e., \mathbf{W}_{L-1} , via gradient flow with squared hinge loss, and our goal is to prove that gradient flow is able to find a NTRF model within the function class $\mathcal{F}(\mathbf{W}^{(0)}, R)$ around the initialization, i.e., achieving $n^{-1} \sum_{i=1}^n \ell(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)) \leq \epsilon$. Let $\mathbf{W}(t)$ be the weights at time t , gradient flow implies that

$$\frac{d\tilde{L}_S(\mathbf{W}(t))}{dt} = -\|\nabla_{\mathbf{W}_{L-1}} \tilde{L}_S(\mathbf{W}(t))\|_F^2 \leq -\frac{Cm\phi}{n^5} \left(\sum_{i=1}^n \tilde{\ell}'(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}(t)}(\mathbf{x}_i)) \right)^2 = \frac{4Cm\phi \tilde{L}_S(\mathbf{W}(t))}{n^3},$$

where the first equality is due to the fact that we only train the last hidden layer, the first inequality is by Lemma C.2 and the second equality follows from the fact that $\tilde{\ell}'(\cdot) = -2\sqrt{\tilde{\ell}(\cdot)}$. Solving the above inequality gives

$$\tilde{L}_S(\mathbf{W}(t)) \leq \tilde{L}_S(\mathbf{W}(0)) \cdot \exp\left(-\frac{4Cm\phi t}{n^3}\right). \quad (\text{C.3})$$

Then, set $T = \mathcal{O}(n^3 m^{-1} \phi^{-1} \cdot \log(\tilde{L}_S(\mathbf{W}(0))/\epsilon'))$ and $\epsilon' = 1/n$, we have $\tilde{L}_S(\mathbf{W}(t)) \leq \epsilon'$. Then it follows that $\tilde{\ell}(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}(t)}(\mathbf{x}_i)) \leq 1$, which implies that $y_i F_{\mathbf{W}^{(0)}, \mathbf{W}(t)}(\mathbf{x}_i) \geq \log(\epsilon)$ and thus $n^{-1} \sum_{i=1}^n \ell(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)) \leq \epsilon$. Therefore, $\mathbf{W}(T)$ is exactly the NTRF model we are looking for.

469 The next step is to characterize the distance between $\mathbf{W}(T)$ and $\mathbf{W}(0)$ in order to characterize the
 470 quantity of R . Note that $\|\nabla_{\mathbf{W}_{L-1}} \tilde{L}_S(\mathbf{W}(t))\|_F^2 \geq 4Cm\phi \tilde{L}_S(\mathbf{W}(t))/n^3$, we have

$$\frac{d\sqrt{\tilde{L}_S(\mathbf{W}(t))}}{dt} = -\frac{\|\nabla_{\mathbf{W}_{L-1}} \tilde{L}_S(\mathbf{W}(t))\|_F^2}{2\sqrt{\tilde{L}_S(\mathbf{W}(t))}} \leq -\|\nabla_{\mathbf{W}_{L-1}} \tilde{L}_S(\mathbf{W}(t))\|_F \cdot \frac{C^{1/2}m^{1/2}\phi^{1/2}}{n^{3/2}}.$$

471 Taking integral on both sides and rearranging terms, we have

$$\int_{t=0}^T \|\nabla_{\mathbf{W}_{L-1}} \tilde{L}_S(\mathbf{W}(t))\|_F dt \leq \frac{n^{3/2}}{C^{1/2}m^{1/2}\phi^{1/2}} \cdot \left(\sqrt{\tilde{L}_S(\mathbf{W}(0))} - \sqrt{\tilde{L}_S(\mathbf{W}(T))} \right).$$

472 Note that the L.H.S. of the above inequality is an upper bound of $\|\mathbf{W}(t) - \mathbf{W}(0)\|_F$, we have for any
 473 $t \geq 0$,

$$\|\mathbf{W}(t) - \mathbf{W}(0)\|_F \leq \frac{n^{3/2}}{C^{1/2}m^{1/2}\phi^{1/2}} \cdot \sqrt{\tilde{L}_S(\mathbf{W}(0))} = \mathcal{O}\left(\frac{n^{3/2} \log(n/(\delta\epsilon))}{m^{1/2}\phi^{1/2}}\right),$$

474 where the second inequality is by Lemma C.1 and our choice of $\lambda = \log(1/\epsilon) + 1$. This implies that
 475 there exists a point \mathbf{W}^* within the class $\mathcal{F}(\mathbf{W}^{(0)}, R)$ with

$$R = \mathcal{O}\left(\frac{n^{3/2} \log(n/(\delta\epsilon))}{\phi^{1/2}}\right)$$

476 such that

$$\epsilon_{\text{NTRF}} := n^{-1} \sum_{i=1}^n \ell(y_i f_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)) \leq \epsilon.$$

477 Then by Theorem 3.3, and, more specifically, (B.1), we can compute the minimal required neural
 478 network width as follows,

$$m = \tilde{\Omega}(R^8 L^{22}) = \tilde{\Omega}\left(\frac{L^{22} n^{12}}{\phi^4}\right).$$

479 This completes the proof. \square

480 D Proof of Technical Lemmas

481 Here we provide the proof of Lemmas B.1, B.4 and B.5.

482 D.1 Proof of Lemma B.1

483 The detailed proof of Lemma B.1 is given as follows.

484 *Proof of Lemma B.1.* Based on the update rule of gradient descent, i.e., $\mathbf{W}^{(t+1)} = \mathbf{W}^{(t)} -$
 485 $\eta \nabla_{\mathbf{W}} L_S(\mathbf{W}^{(t)})$, we have the following calculation.

$$\begin{aligned} & \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_F^2 - \|\mathbf{W}^{(t+1)} - \mathbf{W}^*\|_F^2 \\ &= \underbrace{\frac{2\eta}{n} \sum_{i=1}^n \langle \mathbf{W}^{(t)} - \mathbf{W}^*, \nabla_{\mathbf{W}} L_i(\mathbf{W}^{(t)}) \rangle}_{I_1} - \underbrace{\eta^2 \sum_{l=1}^L \|\nabla_{\mathbf{W}_l} L_S(\mathbf{W}^{(t)})\|_F^2}_{I_2}, \end{aligned} \quad (\text{D.1})$$

486 where the equation follows from the fact that $L_S(\mathbf{W}^{(t)}) = n^{-1} \sum_{i=1}^n L_i(\mathbf{W}^{(t)})$. In what follows,
 487 we first bound the term I_1 on the R.H.S. of (D.1) by approximating the neural network functions with
 488 linear models. By assumption, for $t = 0, \dots, t' - 1$, $\mathbf{W}^{(t)}, \mathbf{W}^* \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)$. Therefore by the
 489 definition of $\epsilon_{\text{app}}(\tau)$,

$$y_i \cdot \langle \nabla f_{\mathbf{W}^{(t)}}(\mathbf{x}_i), \mathbf{W}^{(t)} - \mathbf{W}^* \rangle \leq y_i \cdot (f_{\mathbf{W}^{(t)}}(\mathbf{x}_i) - f_{\mathbf{W}^*}(\mathbf{x}_i)) + \epsilon_{\text{app}}(\tau) \quad (\text{D.2})$$

490 Moreover, we also have

$$\begin{aligned} 0 &\leq y_i \cdot (f_{\mathbf{W}^*}(\mathbf{x}_i) - f_{\mathbf{W}^{(0)}}(\mathbf{x}_i) - \langle \nabla f_{\mathbf{W}^{(0)}}(\mathbf{x}_i), \mathbf{W}^* - \mathbf{W}^{(0)} \rangle) + \epsilon_{\text{app}}(\tau) \\ &= y_i \cdot (f_{\mathbf{W}^*}(\mathbf{x}_i) - F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)) + \epsilon_{\text{app}}(\tau), \end{aligned} \quad (\text{D.3})$$

491 where the equation follows by the definition of $F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x})$. Adding (D.3) to (D.2) and canceling
492 the terms $y_i \cdot f_{\mathbf{W}^*}(\mathbf{x}_i)$, we obtain that

$$y_i \cdot \langle \nabla f_{\mathbf{W}^{(t)}}(\mathbf{x}_i), \mathbf{W}^{(t)} - \mathbf{W}^* \rangle \leq y_i \cdot (f_{\mathbf{W}^{(t)}}(\mathbf{x}_i) - F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)) + 2\epsilon_{\text{app}}(\tau). \quad (\text{D.4})$$

493 We can now give a lower bound on first term on the R.H.S. of (D.1). For $i = 1, \dots, n$, applying the
494 chain rule on the loss function gradients and utilizing (D.4), we have

$$\begin{aligned} \langle \mathbf{W}^{(t)} - \mathbf{W}^*, \nabla_{\mathbf{W}} L_i(\mathbf{W}^{(t)}) \rangle &= \ell'(y_i f_{\mathbf{W}^{(t)}}(\mathbf{x}_i)) \cdot y_i \cdot \langle \mathbf{W}^{(t)} - \mathbf{W}^*, \nabla_{\mathbf{W}} f_{\mathbf{W}^{(t)}}(\mathbf{x}_i) \rangle \\ &\geq \ell'(y_i f_{\mathbf{W}^{(t)}}(\mathbf{x}_i)) \cdot (y_i f_{\mathbf{W}^{(t)}}(\mathbf{x}_i) - y_i f_{\mathbf{W}^*}(\mathbf{x}_i) + 2\epsilon_{\text{app}}(\tau)) \\ &\geq (1 - 2\epsilon_{\text{app}}(\tau)) \ell(y_i f_{\mathbf{W}^{(t)}}(\mathbf{x}_i)) - \ell(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)), \end{aligned} \quad (\text{D.5})$$

495 where the first inequality is by the fact that $\ell'(y_i f_{\mathbf{W}^{(t)}}(\mathbf{x}_i)) < 0$, the second inequality is by convexity
496 of $\ell(\cdot)$ and the fact that $-\ell'(y_i f_{\mathbf{W}^{(t)}}(\mathbf{x}_i)) \leq \ell(y_i f_{\mathbf{W}^{(t)}}(\mathbf{x}_i))$.

497 We now proceed to bound the term I_2 on the R.H.S. of (D.1). Note that we have $\ell'(\cdot) < 0$, and
498 therefore the Frobenius norm of the gradient $\nabla_{\mathbf{W}_i} L_S(\mathbf{W}^{(t)})$ can be upper bounded as follows,

$$\begin{aligned} \|\nabla_{\mathbf{W}_i} L_S(\mathbf{W}^{(t)})\|_F &= \left\| \frac{1}{n} \sum_{i=1}^n \ell'(y_i f_{\mathbf{W}^{(t)}}(\mathbf{x}_i)) \nabla_{\mathbf{W}_i} f_{\mathbf{W}^{(t)}}(\mathbf{x}_i) \right\|_F \\ &\leq \frac{1}{n} \sum_{i=1}^n -\ell'(y_i f_{\mathbf{W}^{(t)}}(\mathbf{x}_i)) \cdot \|\nabla_{\mathbf{W}_i} f_{\mathbf{W}^{(t)}}(\mathbf{x}_i)\|_F, \end{aligned}$$

499 where the inequality follows by triangle inequality. We now utilize the fact that cross-entropy loss
500 satisfies the inequalities $-\ell'(\cdot) \leq \ell(\cdot)$ and $-\ell'(\cdot) \leq 1$. Therefore by definition of $M(\tau)$, we have

$$\begin{aligned} \sum_{i=1}^L \|\nabla_{\mathbf{W}_i} L_S(\mathbf{W}^{(t)})\|_F^2 &\leq \mathcal{O}(LM(\tau)^2) \cdot \left(\frac{1}{n} \sum_{i=1}^n -\ell'(y_i f_{\mathbf{W}^{(t)}}(\mathbf{x}_i)) \right)^2 \\ &\leq \mathcal{O}(LM(\tau)^2) \cdot L_S(\mathbf{W}^{(t)}). \end{aligned} \quad (\text{D.6})$$

501 Then we can plug (D.5) and (D.6) into (D.1) and obtain

$$\begin{aligned} &\|\mathbf{W}^{(t)} - \mathbf{W}^*\|_F^2 - \|\mathbf{W}^{(t+1)} - \mathbf{W}^*\|_F^2 \\ &\geq \frac{2\eta}{n} \sum_{i=1}^n \left[(1 - 2\epsilon_{\text{app}}(\tau)) \ell(y_i f_{\mathbf{W}^{(t)}}(\mathbf{x}_i)) - \ell(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)) \right] - \mathcal{O}(\eta^2 LM(\tau)^2) \cdot L_S(\mathbf{W}^{(t)}) \\ &\geq \left[\frac{3}{2} - 4\epsilon_{\text{app}}(\tau) \right] \eta L_S(\mathbf{W}^{(t)}) - \frac{2\eta}{n} \sum_{i=1}^n \ell(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)), \end{aligned}$$

502 where the last inequality is by $\eta = \mathcal{O}(L^{-1}M(\tau)^{-2})$ and merging the third term on the second line
503 into the first term. Taking telescope sum from $t = 0$ to $t = t' - 1$ and plugging in the definition
504 $\frac{1}{n} \sum_{i=1}^n \ell(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)) = \epsilon_{\text{NTRF}}$ completes the proof. \square

505 D.2 Proof of Lemma B.4

506 *Proof of Lemma B.4.* We first denote $\mathcal{W} = \mathcal{B}(\mathbf{W}^{(0)}, \tilde{R} \cdot m^{-1/2})$, and define the corresponding
507 neural network function class and surrogate loss function class as $\mathcal{F} = \{f_{\mathbf{W}}(\mathbf{x}) : \mathbf{W} \in \mathcal{W}\}$ and
508 $\mathcal{G} = \{-\ell[y \cdot f_{\mathbf{W}}(\mathbf{x})] : \mathbf{W} \in \mathcal{W}\}$ respectively.

509 By standard uniform convergence results in terms of empirical Rademacher complexity [7, 18, 20],
510 with probability at least $1 - \delta$ we have

$$\sup_{\mathbf{W} \in \mathcal{W}} |\mathcal{E}_S(\mathbf{W}) - \mathcal{E}_{\mathcal{D}}(\mathbf{W})| = \sup_{\mathbf{W} \in \mathcal{W}} \left| -\frac{1}{n} \sum_{i=1}^n \ell'[y_i \cdot f_{\mathbf{W}}(\mathbf{x}_i)] + \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \ell'[y \cdot f_{\mathbf{W}}(\mathbf{x})] \right|$$

$$\leq 2\hat{\mathfrak{R}}_n(\mathcal{G}) + C_1 \sqrt{\frac{\log(1/\delta)}{n}},$$

511 where C_1 is an absolute constant, and

$$\hat{\mathfrak{R}}_n(\mathcal{G}) = \mathbb{E}_{\xi_i \sim \text{Unif}(\{\pm 1\})} \left\{ \sup_{\mathbf{W} \in \mathcal{W}} \frac{1}{n} \sum_{i=1}^n \xi_i \ell' [y_i \cdot f_{\mathbf{W}}(\mathbf{x}_i)] \right\}$$

512 is the empirical Rademacher complexity of the function class \mathcal{G} . We now provide two bounds on
 513 $\hat{\mathfrak{R}}_n(\mathcal{G})$, whose combination gives the final result of Lemma B.4. First, by Corollary 5.35 in [22], with
 514 probability at least $1 - L \cdot \exp(-\Omega(m))$, $\|\mathbf{W}_l^{(0)}\|_2 \leq 3$ for all $l \in [L]$. Therefore for all $\mathbf{W} \in \mathcal{W}$, we
 515 have $\|\mathbf{W}_l\|_2 \leq 4$. Moreover, standard concentration inequalities on the norm of the first row of $\mathbf{W}_l^{(0)}$
 516 also implies that $\|\mathbf{W}_l\|_2 \geq 0.5$ for all $\mathbf{W} \in \mathcal{W}$ and $l \in [L]$. Therefore, an adaptation of the bound in
 517 [6]⁵ gives

$$\begin{aligned} \hat{\mathfrak{R}}_n(\mathcal{F}) &\leq \tilde{\mathcal{O}} \left(\sup_{\mathbf{W} \in \mathcal{W}} \left\{ \frac{m^{1/2}}{\sqrt{n}} \cdot \prod_{l=1}^L \|\mathbf{W}_l\|_2 \cdot \left[\sum_{l=1}^L \frac{\|\mathbf{W}_l^\top - \mathbf{W}_l^{(0)\top}\|_{2,1}^{2/3}}{\|\mathbf{W}_l\|_2^{2/3}} \right]^{3/2} \right\} \right) \\ &\leq \tilde{\mathcal{O}} \left(\sup_{\mathbf{W} \in \mathcal{W}} \left\{ \frac{4^L m^{1/2}}{\sqrt{n}} \cdot \left[\sum_{l=1}^L (\sqrt{m} \cdot \|\mathbf{W}_l^\top - \mathbf{W}_l^{(0)\top}\|_F)^{2/3} \right]^{3/2} \right\} \right) \\ &\leq \tilde{\mathcal{O}} \left(4^L L^{3/2} \tilde{R} \cdot \sqrt{\frac{m}{n}} \right). \end{aligned} \quad (\text{D.7})$$

518 We now derive the second bound on $\hat{\mathfrak{R}}_n(\mathcal{G})$, which is inspired by the proof provided in [9]. Since
 519 $y \in \{+1, 1\}$, $|\ell'(z)| \leq 1$ and $\ell'(z)$ is 1-Lipschitz continuous, by standard empirical Rademacher
 520 complexity bounds [7, 18, 20], we have

$$\hat{\mathfrak{R}}_n(\mathcal{G}) \leq \hat{\mathfrak{R}}_n(\mathcal{F}) = \mathbb{E}_{\xi_i \sim \text{Unif}(\{\pm 1\})} \left[\sup_{\mathbf{W} \in \mathcal{W}} \frac{1}{n} \sum_{i=1}^n \xi_i f_{\mathbf{W}}(\mathbf{x}_i) \right],$$

521 where $\hat{\mathfrak{R}}_n(\mathcal{F})$ is the empirical Rademacher complexity of the function class \mathcal{F} . We have

$$\hat{\mathfrak{R}}_n[\mathcal{F}] \leq \underbrace{\mathbb{E}_{\xi} \left\{ \sup_{\mathbf{W} \in \mathcal{W}} \frac{1}{n} \sum_{i=1}^n \xi_i [f_{\mathbf{W}}(\mathbf{x}_i) - F_{\mathbf{W}^{(0)}, \mathbf{W}}(\mathbf{x}_i)] \right\}}_{I_1} + \underbrace{\mathbb{E}_{\xi} \left\{ \sup_{\mathbf{W} \in \mathcal{W}} \frac{1}{n} \sum_{i=1}^n \xi_i F_{\mathbf{W}^{(0)}, \mathbf{W}}(\mathbf{x}_i) \right\}}_{I_2}, \quad (\text{D.8})$$

522 where $F_{\mathbf{W}^{(0)}, \mathbf{W}}(\mathbf{x}) = f_{\mathbf{W}^{(0)}}(\mathbf{x}) + \langle \nabla_{\mathbf{W}} f_{\mathbf{W}^{(0)}}(\mathbf{x}), \mathbf{W} - \mathbf{W}^{(0)} \rangle$. For I_1 , by Lemma 4.1 in [8], with
 523 probability at least $1 - \delta/2$ we have

$$I_1 \leq \max_{i \in [n]} |f_{\mathbf{W}}(\mathbf{x}_i) - F_{\mathbf{W}^{(0)}, \mathbf{W}}(\mathbf{x}_i)| \leq \mathcal{O}(L^3 \tilde{R}^{4/3} m^{-1/6} \sqrt{\log(m)}),$$

524 For I_2 , note that $\mathbb{E}_{\xi} \left[\sup_{\mathbf{W} \in \mathcal{W}} \sum_{i=1}^n \xi_i f_{\mathbf{W}^{(0)}}(\mathbf{x}_i) \right] = 0$. By Cauchy-Schwarz inequality we have

$$I_2 = \frac{1}{n} \sum_{l=1}^L \mathbb{E}_{\xi} \left\{ \sup_{\|\tilde{\mathbf{W}}_l\|_F \leq \tilde{R} m^{-1/2}} \text{Tr} \left[\tilde{\mathbf{W}}_l^\top \sum_{i=1}^n \xi_i \nabla_{\mathbf{W}_l} f_{\mathbf{W}^{(0)}}(\mathbf{x}_i) \right] \right\} \leq \frac{\tilde{R} m^{-1/2}}{n} \sum_{l=1}^L \mathbb{E}_{\xi} \left[\left\| \sum_{i=1}^n \xi_i \nabla_{\mathbf{W}_l} f_{\mathbf{W}^{(0)}}(\mathbf{x}_i) \right\|_F \right].$$

525 Therefore

$$I_2 \leq \frac{\tilde{R} m^{-1/2}}{n} \sum_{l=1}^L \sqrt{\mathbb{E}_{\xi} \left[\left\| \sum_{i=1}^n \xi_i \nabla_{\mathbf{W}_l} f_{\mathbf{W}^{(0)}}(\mathbf{x}_i) \right\|_F^2 \right]} = \frac{\tilde{R} m^{-1/2}}{n} \sum_{l=1}^L \sqrt{\sum_{i=1}^n \|\nabla_{\mathbf{W}_l} f_{\mathbf{W}^{(0)}}(\mathbf{x}_i)\|_F^2} \leq \mathcal{O} \left(\frac{L \cdot \tilde{R}}{\sqrt{n}} \right),$$

⁵Bartlett et al. [6] only proved the Rademacher complexity bound for the composition of the ramp loss and the neural network function. In our setting essentially the ramp loss is replaced with the $-\ell'(\cdot)$ function, which is bounded and 1-Lipschitz continuous. The proof in our setting is therefore exactly the same as the proof given in [6], and we can apply Theorem 3.3 and Lemma A.5 in [6] to obtain the desired bound we present here.

526 where we apply Jensen’s inequality to obtain the first inequality, and the last inequality follows by
 527 Lemma B.3 in [8]. Combining the bounds of I_1 and I_2 gives

$$\hat{\mathfrak{R}}_n[\mathcal{F}] \leq \tilde{\mathcal{O}}\left(\frac{L\tilde{R}}{\sqrt{n}} + \frac{L^3\tilde{R}^{4/3}}{m^{1/6}}\right).$$

528 Further combining this bound with (D.7) and recaling δ completes the proof. \square

529 D.3 Proof of Lemma B.5

530 *Proof of Lemma B.5.* Different from the proof of Lemma B.1, online SGD only queries one data to
 531 update the model parameters in each iteration, i.e., $\mathbf{W}^{i+1} = \mathbf{W}^i - \eta \nabla L_{i+1}(\mathbf{W}^{(i)})$. By this update
 532 rule, we have

$$\begin{aligned} & \|\mathbf{W}^{(i)} - \mathbf{W}^*\|_F^2 - \|\mathbf{W}^{(i+1)} - \mathbf{W}^*\|_F^2 \\ &= 2\eta \langle \mathbf{W}^{(i)} - \mathbf{W}^*, \nabla_{\mathbf{W}} L_{i+1}(\mathbf{W}^{(i)}) \rangle - \eta^2 \sum_{l=1}^L \|\nabla_{\mathbf{W}_l} L_{i+1}(\mathbf{W}^{(i)})\|_F^2. \end{aligned} \quad (\text{D.9})$$

533 With exactly the same proof as (D.5) in the proof of Lemma B.1, we have

$$\langle \mathbf{W}^{(t)} - \mathbf{W}^*, \nabla_{\mathbf{W}} L_i(\mathbf{W}^{(t)}) \rangle \geq (1 - 2\epsilon_{\text{app}}(\tau)) \ell(y_i f_{\mathbf{W}^{(t)}}(\mathbf{x}_i)) - \ell(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)), \quad (\text{D.10})$$

534 for all $i = 0, \dots, n' - 1$. By the fact that $-\ell'(\cdot) \leq \ell(\cdot)$ and $-\ell'(\cdot) \leq 1$, we have

$$\begin{aligned} \sum_{l=1}^L \|\nabla_{\mathbf{W}_l} L_{i+1}(\mathbf{W}^{(i)})\|_F^2 &\leq \sum_{l=1}^L \ell(y_{i+1} f_{\mathbf{W}^{(i)}}(\mathbf{x}_{i+1})) \cdot \|\nabla_{\mathbf{W}_l} f_{\mathbf{W}^{(i)}}(\mathbf{x}_{i+1})\|_F^2 \\ &\leq \mathcal{O}(LM(\tau)^2) \cdot L_{i+1}(\mathbf{W}^{(i)}). \end{aligned} \quad (\text{D.11})$$

535 Then plugging (D.10) and (D.11) into (D.9) gives

$$\begin{aligned} & \|\mathbf{W}^{(i)} - \mathbf{W}^*\|_F^2 - \|\mathbf{W}^{(i+1)} - \mathbf{W}^*\|_F^2 \\ &\geq (2 - 4\epsilon_{\text{app}}(\tau)) \eta L_{i+1}(\mathbf{W}^{(i)}) - 2\eta \ell(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)) - \mathcal{O}(\eta^2 LM(\tau)^2) L_{i+1}(\mathbf{W}^{(i)}) \\ &\geq \left(\frac{3}{2} - 4\epsilon_{\text{app}}(\tau)\right) \eta L_{i+1}(\mathbf{W}^{(i)}) - 2\eta \ell(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)), \end{aligned}$$

536 where the last inequality is by $\eta = \mathcal{O}(L^{-1}M(\tau)^{-2})$ and merging the third term on the second line
 537 into the first term. Taking telescope sum over $i = 0, \dots, n' - 1$, we obtain

$$\begin{aligned} & \|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F^2 - \|\mathbf{W}^{(n')} - \mathbf{W}^*\|_F^2 \\ &\geq \left(\frac{3}{2} - 4\epsilon_{\text{app}}(\tau)\right) \eta \sum_{i=1}^{n'} L_i(\mathbf{W}^{(i-1)}) - 2\eta \sum_{i=1}^{n'} \ell(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)). \\ &\geq \left(\frac{3}{2} - 4\epsilon_{\text{app}}(\tau)\right) \eta \sum_{i=1}^{n'} L_i(\mathbf{W}^{(i-1)}) - 2\eta \sum_{i=1}^n \ell(y_i F_{\mathbf{W}^{(0)}, \mathbf{W}^*}(\mathbf{x}_i)). \\ &\geq \left(\frac{3}{2} - 4\epsilon_{\text{app}}(\tau)\right) \eta \sum_{i=1}^{n'} L_i(\mathbf{W}^{(i-1)}) - 2n\eta \epsilon_{\text{NTRF}}. \end{aligned}$$

538 This finishes the proof. \square

539 E Experiments

540 In this section, we conduct some simple experiments to validate our theory. Since our paper mainly
 541 focuses on binary classification, we use a subset of the original CIFAR10 dataset [16], which only has
 542 two classes of images. We train a 5-layer fully-connected ReLU network on this binary classification

543 dataset with different sample sizes ($n \in \{100, 200, 500, 1000, 2000, 5000, 10000\}$), and plot the
 544 minimal neural network width that is required to achieve zero training error in Figure 1 (solid line).
 545 We also plot $\mathcal{O}(n)$, $\mathcal{O}(\log^3(n))$, $\mathcal{O}(\log^2(n))$ and $\mathcal{O}(\log(n))$ in dashed line for reference. It is evident
 546 that the required network width to achieve zero training error is polylogarithmic on the sample size n ,
 547 which is consistent with our theory.

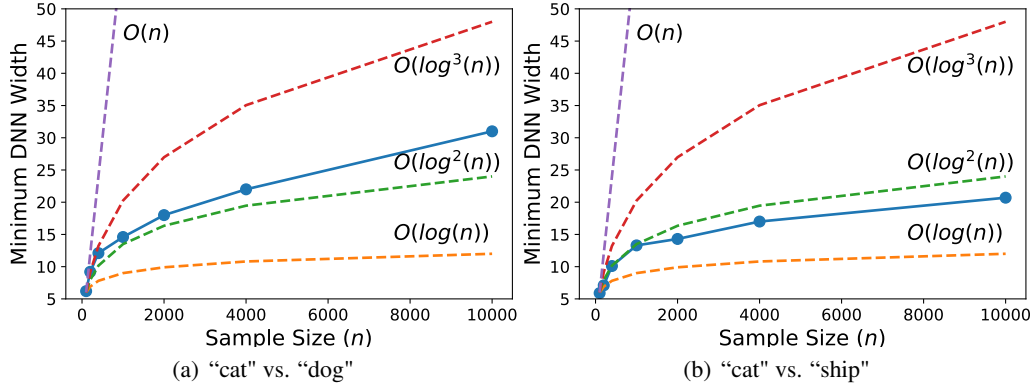


Figure 1: Minimum network width that is required to achieve zero training error with respect to the training sample size (blue solid line). The hidden constants in all $\mathcal{O}(\cdot)$ notations are adjusted to ensure their plots (dashed lines) start from the same point.