

1 We would like to thank the reviewers for their valuable comments and encouraging feedback. Below, we address the
2 concerns raised in the order they appeared.

3 **Reviewer 1.**

4 **Q. Could the authors elaborate a little on r_j (line 245)? [...] specific because of the construction in appendix E?**

5 As seen in Figure 2, the construction in Appendix E partitions the dataset into several groups of size N_1, N_2, \dots, N_m
6 and uses the construction in Theorem 3.1 as building blocks to fit each group of data. The additional requirement
7 of r_j nodes corresponds to circle and diamond nodes in Figure 2, because we need to propagate input information
8 using a hidden node up to layer l_m (the last “building block” that fits N_m points), and propagate output information
9 using d_y hidden nodes starting from layer $l_1 + 2$ (after the first group of N_1 points are fitted). This is why we need
10 $r_j = d_y \mathbf{1}\{j > 1\} + \mathbf{1}\{j < m\}$ extra nodes in Proposition 3.4. The reason why we need only one hidden node for input
11 information is because we down-project x_i ’s onto a line. This requirement of r_j is indeed specific to the construction.

12 **Reviewer 2.**

13 We appreciate the reviewer’s positive comments about our submission.

14 **Reviewer 3.**

15 **Q. Omega notation throughout the paper is erroneous?**

16 We agree that our choice of this notation might cause confusion. The reason why we used $\Omega(\sqrt{N})$ for sufficiency
17 statements, e.g., in line 2, was that anything of greater order than \sqrt{N} (e.g., N) could also memorize N data points.
18 However, we agree that it is erroneous to say that $\Omega(\sqrt{N})$ is **necessary**, because this may sound as if N nodes are also
19 necessary for memorization. We will correct the Ω notation to Θ throughout the paper. Thank you for the correction.

20 **Q. Comparison to [Soudry and Carmon, 2016]?**

21 Thank you for pointing out a relevant paper. We believe that their results are not directly comparable to ours because
22 there are a few important differences in the problem settings. The biggest difference is that [SC16] consider a setting
23 where there is a multiplicative “dropout noise” at each hidden node and each data point. At i -th node of l -th layer, the
24 slope of the activation function for the n -th data point is either $\epsilon_{i,l}^{(n)} \cdot 1$ (if input is positive) or $\epsilon_{i,l}^{(n)} \cdot s$ (if negative, $s \neq 0$),
25 where $\epsilon_{i,l}^{(n)}$ is the multiplicative random (e.g., Gaussian) dropout noise. Their theorem statements hold for “almost every
26 realization” of these dropout noise factors, so heavily depend on their particular model. In contrast, our setting is free
27 of these noise terms, and hence corresponds to a **specific** realization of such $\epsilon_{i,l}^{(n)}$ ’s. The discussion after Theorem 5
28 (the multiple hidden layer result) of [SC16] suggests that their proof crucially depends on the “except measure zero”
29 argument on some of the noise terms $\epsilon_{\cdot, L-1}^{(\cdot)}$, hence making our results not directly comparable to theirs.

30 **Q. Theorem 5.1: The statement was a bit difficult to parse. [...] in the theorem and their magnitude.**

31 The exact values of positive constants can be found in Appendix G, and they are dependent on a number of terms, such
32 as the number of data points N , batch size B , the radius ρ_s of a ball in which the slopes of activation don’t change,
33 the Taylor expansions of loss $\ell(f_{\theta^*}(x_i); y_i)$ and network output $f_{\theta^*}(x_i)$ around the memorizing global minimum θ^* ,
34 maximum and minimum strictly positive eigenvalues of $H = \sum_{i=1}^N \ell''(f_{\theta^*}(x_i); y_i) \nu_i \nu_i^T$, where $\nu_i = \nabla_{\theta} f_{\theta^*}(x_i)$. We
35 will make sure to add detailed explanation and an improved theorem statement in the next revision.

36 **Q. If I understand correctly, [...] the initialization must be significantly close in the first place?**

37 It is indeed true that in the worst case the point found at t^* can be farther away from θ^* , and also that our theorem
38 requires the initialization to be close to the global minimum. However, an initialization that is ϵ -close (in Euclidean
39 distance) to the global minimum has empirical risk $O(\epsilon^2)$ (shown by Lemma G.1). Thus, in terms of risk value, the
40 initialization is not necessarily as close to the global minimum compared to the point found at t^* , which achieves $O(\epsilon^4)$
41 risk. We’d like to highlight that one can start off at a $O(\epsilon^2)$ -risk point and quickly find a $O(\epsilon^4)$ -risk point.

42 **Q. Comparison to [Zhong et al., 2017]?**

43 Thanks for bringing up this point. We would like to emphasize that the settings are quite different, so one cannot
44 make direct comparisons. [ZSJ+17] consider 1-hidden-layer networks with ± 1 -valued weights at the output layer,
45 with an implicit assumption that network width is smaller than input dimension. Input x_i is Gaussian and output y_i is
46 generated by a “teacher” network. In comparison, we consider arbitrary datasets and networks, under a mild assumption
47 (especially for overparametrized networks) that the network can memorize the data. We are happy to add more detailed
48 comparisons to our next revision.

49 **Q. Note that in their analysis, Zhong et al. show that the strong convexity [...] this technique?**

50 It seems difficult to show strong convexity in our case. For example, for ReLU networks, if we scale one layer by α and
51 the next layer by α^{-1} , we get exactly the same network. This means that for ReLU, the global minimum always has a
52 direction in which the risk value does not change, hence strong convexity cannot hold at global minima. The key to
53 [ZSJ+17]’s result is that they fix the output layer parameters to ± 1 , and only consider hidden layer parameters.