

Retraction note for “Nonparametric Regression with Shallow Overparameterized Neural Networks Trained by GD with Early Stopping”

Ilja Kuzborskij
DeepMind, London

ILJAK@DEEPMIND.COM

Csaba Szepesvári
DeepMind, Canada and University of Alberta, Edmonton

SZEPI@DEEPMIND.COM

Editors: Mikhail Belkin and Samory Kpotufe

Abstract

Theorems 1 and 2 appearing in the paper contain a shared critical issue and this note serves as a retraction note. The issue is traced to the Theorem 4, which concerns the control of the Lipschitzness of a Gradient Descent (GD)-trained overparameterized shallow neural network in the input. In the following we describe the issue in some detail.

Setting. We first briefly recall the setting the paper. The paper considers the problem of nonparametric regression with random design where inputs $\mathbf{X}, \mathbf{X}_1, \dots, \mathbf{X}_n \sim P_X \in \mathcal{M}_1(\mathbb{S}^{d-1})$ are i.i.d. and targets are generated by the Lipschitz (possibly non-differentiable) and bounded regression function f^* on \mathbb{S}^{d-1} , as $Y_i = f^*(\mathbf{X}_i) + \varepsilon_i$, with σ -subgaussian independent noise $\varepsilon_1, \dots, \varepsilon_n$. The paper considers learning of f^* by the shallow neural network predictor

$$\hat{f}_{\mathbf{W}}(\mathbf{x}) = \sum_{k=1}^m u_k \phi\left(\mathbf{w}_k^\top \mathbf{x}\right), \quad (\mathbf{x} \in \mathbb{S}^{d-1}, \mathbf{W} \in \mathbb{R}^{d \times m})$$

defined with respect to a fixed twice differentiable activation function $\phi : \mathbb{R} \rightarrow \mathbb{R}$, where ϕ', ϕ'' are bounded on \mathbb{R} . The predictor is parameterized by an output layer, a (non-tunable) random weight vector $\mathbf{u} \stackrel{\text{iid}}{\sim} \text{unif}(\{\pm 1/\sqrt{m}\})^m$ (sampled once at the beginning of the training), and a tunable hidden layer weight matrix $\mathbf{W} \in \mathbb{R}^{d \times m}$, where m is the width of the network. In particular, weights of a hidden layer are obtained by the standard GD procedure by approximately minimizing the empirical risk

$$\hat{L}(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n (\hat{f}_{\mathbf{W}}(\mathbf{X}_i) - Y_i)^2.$$

Here, \mathbf{W}_0 is obtained by drawing each entry from $\mathcal{N}(0, \nu_{\text{init}}^2)$ independently from each other, while \mathbf{W}_T is obtained by the recursive update rule $\mathbf{W}_{t+1} = \mathbf{W}_t - \eta \nabla \hat{L}(\mathbf{W}_t)$ for $t = 0, \dots, T-1$ and where $\eta > 0$ is a fixed step size.

Claimed result and proof sketch. The aim of the paper is to control the excess risk defined as

$$\mathcal{E}_T = \mathbb{E} [(f_{\mathbf{W}_T}(\mathbf{X}) - f^*(\mathbf{X}))^2 \mid \mathbf{u}, \mathbf{W}_0] .$$

In particular, for $\sigma = 0^1$ the paper claims a minimax optimal rate with high probability over $(\mathbf{u}, \mathbf{W}_0)$, assuming that the network is overparameterized, $m \geq \text{poly}((n/d)(1 + \nu_{\text{init}}^2))$:

$$\mathcal{E}_T = \mathcal{O}_{\mathbb{P}} \left((\text{Lip}(f^*)^2 + (1 + d\nu_{\text{init}}^2)^2) n^{-\frac{2}{2+d}} \right) \quad \text{as } n \rightarrow \infty ,$$

which is based on the following proof idea. Introduce a nearest-neighbor operator (with ties broken arbitrarily)

$$\pi(\mathbf{x}) = \arg \min_{i \in [n]} \|\mathbf{x} - \mathbf{X}_i\|_2 , \quad (\mathbf{x} \in \mathbb{S}^{d-1})$$

and consider a decomposition of the excess risk (here $\mathbb{E}[\cdot] = \mathbb{E}[\cdot \mid \mathbf{u}, \mathbf{W}_0]$):

$$\underbrace{\mathbb{E} (f^*(\mathbf{X}) - f^*(\mathbf{X}_{\pi(\mathbf{X})}))^2}_{(i)} + \underbrace{\mathbb{E} (f^*(\mathbf{X}_{\pi(\mathbf{X})}) - \hat{f}_{\mathbf{W}_T}(\mathbf{X}_{\pi(\mathbf{X})}))^2}_{(ii)} + \underbrace{\mathbb{E} (\hat{f}_{\mathbf{W}_T}(\mathbf{X}_{\pi(\mathbf{X})}) - \hat{f}_{\mathbf{W}_T}(\mathbf{X}))^2}_{(iii)} .$$

Here term (i) is controlled thanks to the Lipschitzness of f^* and the standard partitioning analysis of the one nearest-neighbor rule (Lemma 7 in the paper) gives

$$\mathbb{E} \|\mathbf{X} - \mathbf{X}_{\pi(\mathbf{X})}\|_2^2 \lesssim n^{-\frac{2}{2+d}} .$$

Term (ii) is the expected empirical risk $\mathbb{E} \hat{L}(\mathbf{W}_T)$, which is controlled by the convergence of optimization error of GD to zero for shallow overparameterized neural networks (Du et al., 2018; Oymak and Soltanolkotabi, 2020).

Issue with term (iii). The issue occurs in the control of term (iii) which captures Lipschitzness of the predictor trained for T steps. Before discussing the issue we need to introduce some technical notions. In particular, for the analysis we require the Neural Tangent Random Feature (NTRF) predictor, which is obtained by linearization of $\mathbf{W} \mapsto \hat{f}_{\mathbf{W}}(\mathbf{x})$ around \mathbf{W}_0 , namely

$$\hat{f}_{\mathbf{W}}^{\text{rf}}(\mathbf{x}) = \hat{f}_{\mathbf{W}_0}(\mathbf{x}) + \Psi(\mathbf{x})^\top \text{vec}(\mathbf{W} - \mathbf{W}_0) ,$$

where $\Psi(\cdot)$ is called the NTRF feature map and defined as

$$\Psi(\mathbf{x}) = [u_1 \phi'(\mathbf{W}_{0,1}^\top \mathbf{x}) \mathbf{x}^\top, \dots, u_m \phi'(\mathbf{W}_{0,m}^\top \mathbf{x}) \mathbf{x}^\top]^\top .$$

In addition, we introduce NTRF iterates $(\mathbf{W}_t^{\text{rf}})_{t=0}^T$ obtained by running GD as discussed in the beginning of the note, however with the empirical risk replaced by the least-squares objective $\mathbf{W} \mapsto ((\hat{f}_{\mathbf{W}}^{\text{rf}}(\mathbf{X}_1) - Y_1)^2 + \dots + (\hat{f}_{\mathbf{W}}^{\text{rf}}(\mathbf{X}_n) - Y_n)^2)/n$.

The attractive property of the NTRF iterates is that the distance $\|\mathbf{W}_T - \mathbf{W}_T^{\text{rf}}\|_F$ can be made small by overparameterization of the neural network thanks to the so-called ‘coupling’ argument, see for instance (Bartlett et al., 2021, Theorem 5.1) and Appendix B in the paper we discuss. Then,

1. In this note we only consider a noise-free setting, but the same issue applies for the case $\sigma > 0$.

assuming that in addition to the boundedness of derivatives, activation function ϕ also satisfies $\phi'(x) \leq B_{\phi''}|x|$ with $\sup_{x \in \mathbb{R}} |\phi''(x)| \leq B_{\phi''}$ for all $x \in \mathbb{R}$, we attempt to control the Lipschitzness:

$$\begin{aligned} \text{Lip}(\hat{f}_{\mathbf{W}_T}) &= \sup_{\mathbf{x} \in \mathbb{S}^{d-1}} \left\| \sum_{k=1}^m u_k \phi'(\mathbf{W}_{T,k}^\top \mathbf{x}) \mathbf{W}_{T,k} \right\|_2 \\ &\leq \frac{1}{\sqrt{m}} \sup_{\mathbf{x} \in \mathbb{S}^{d-1}} \sum_{k=1}^m |\phi'(\mathbf{W}_{T,k}^\top \mathbf{x})| \|\mathbf{W}_{T,k}\|_2 \\ &\leq \frac{B_{\phi''}}{\sqrt{m}} \sum_{k=1}^m \|\mathbf{W}_{T,k}\|_2^2 \\ &\leq \frac{2B_{\phi''}}{\sqrt{m}} \underbrace{\|\mathbf{W}_T - \mathbf{W}_T^{\text{rf}}\|_F^2}_{(a)} + \frac{2B_{\phi''}}{\sqrt{m}} \underbrace{\|\mathbf{W}_T^{\text{rf}}\|_F^2}_{(b)}. \end{aligned}$$

While (a) is controlled by the aforementioned ‘coupling’ argument, it turns out that $\|\mathbf{W}_T^{\text{rf}}\|_F^2$ cannot be made small enough. The problem comes in assuming that $\mathbf{W}_0^{\text{rf}} = \mathbf{0}$ while it must be $\mathbf{W}_0^{\text{rf}} = \mathbf{W}_0$ for the coupling argument of (Bartlett et al., 2021, Theorem 5.1) to hold. Then, having $\mathbf{W}_0^{\text{rf}} = \mathbf{W}_0$, it is clear that $\|\mathbf{W}_T^{\text{rf}}\|_F^2 \lesssim \|\mathbf{W}_0\|_F^2$ while $\|\mathbf{W}_0\|_F^2 = \mathcal{O}(m)$ which renders the above bound on the Lipschitz constant of order $\mathcal{O}(\sqrt{m})$.

Prospects on fixing the issue. One potential resolution to the above could be to follow a more careful analysis and replace $\|\mathbf{W}_T^{\text{rf}}\|_F^2$ with $\|\mathbf{W}_T^{\text{rf}} - \mathbf{W}_0\|_F^2$, where the later quantity is definitely no larger since GD takes the shortest ℓ^2 path from the initialization for least-squares problems of the considered kind (Oymak and Soltanolkotabi, 2020).² Such an alternative analysis of term (iii) would be based upon a decomposition (for any $\mathbf{x}, \tilde{\mathbf{x}} \in \mathbb{S}^{d-1}$),

$$\hat{f}_{\mathbf{W}_T}(\mathbf{x}) - \hat{f}_{\mathbf{W}_T}(\tilde{\mathbf{x}}) = \underbrace{\hat{f}_{\mathbf{W}_T}(\mathbf{x}) - \hat{f}_{\mathbf{W}_T^{\text{rf}}}(\mathbf{x})}_{(iii.a)} + \underbrace{\hat{f}_{\mathbf{W}_T^{\text{rf}}}(\mathbf{x}) - \hat{f}_{\mathbf{W}_T^{\text{rf}}}(\tilde{\mathbf{x}})}_{(iii.b)} + \underbrace{\hat{f}_{\mathbf{W}_T^{\text{rf}}}(\tilde{\mathbf{x}}) - \hat{f}_{\mathbf{W}_T}(\tilde{\mathbf{x}})}_{(iii.c)}.$$

Here terms (iii.a) and (iii.c) could be controlled thanks to the ‘coupling’ argument as mentioned above. It is not hard to see that

$$\begin{aligned} (iii.b) &= (\hat{f}_{\mathbf{W}_0}(\mathbf{x}) - \hat{f}_{\mathbf{W}_0}(\tilde{\mathbf{x}})) + \text{vec}(\mathbf{W} - \mathbf{W}_0)^\top (\Psi(\mathbf{x}) - \Psi(\tilde{\mathbf{x}})) \\ &\leq c \|\mathbf{x} - \tilde{\mathbf{x}}\|_2 + \|\mathbf{W}_T^{\text{rf}} - \mathbf{W}_0\|_F \|\Psi(\mathbf{x}) - \Psi(\tilde{\mathbf{x}})\|_2 \\ &\leq c \|\mathbf{x} - \tilde{\mathbf{x}}\|_2 + c' \|\mathbf{W}_T^{\text{rf}} - \mathbf{W}_0\|_F \|\mathbf{x} - \tilde{\mathbf{x}}\|_2 \quad (\Psi(\cdot) \text{ is Lipschitz}) \end{aligned}$$

with high probability over $(\mathbf{u}, \mathbf{W}_0)$, where using independence of $(\mathbf{u}, \mathbf{W}_0)$, $\hat{f}_{\mathbf{W}_0}(\cdot)$ is Lipschitz with a universal constant c . Now, it remains to see whether $\|\mathbf{W}_T^{\text{rf}} - \mathbf{W}_0\|_F$ can be controlled. To this end let $(\lambda_i)_{i=1}^n$ be the eigenvalues of the covariance matrix $\Psi(\mathbf{X}_1)\Psi(\mathbf{X}_1)^\top + \dots + \Psi(\mathbf{X}_n)\Psi(\mathbf{X}_n)^\top$ and let $\mathbf{Y}^* = [f^*(\mathbf{X}_1), \dots, f^*(\mathbf{X}_n)]^\top$. The standard analysis of GD dynamics (Yao et al., 2007) leads to the identity

$$\|\mathbf{W}_T^{\text{rf}} - \mathbf{W}_0\|_F^2 = \sum_{i=1}^n (1 - (1 - \frac{2\eta}{n} \lambda_i)^T)^2 \frac{(\mathbf{u}_i^\top \mathbf{Y}^*)^2}{\lambda_i}.$$

2. The quantity $\|\mathbf{W}_T^{\text{rf}} - \mathbf{W}_0\|_F^2$ was considered before in the analysis of GD-trained shallow neural networks as a complexity measure of a GD solution, for instance (Arora et al., 2019).

Now, since we focus on a noise-free setting ($\sigma = 0$), we can consider the limiting solution with $T \rightarrow \infty$ (no need to control overfitting). To this end,

$$\|\mathbf{W}_T^{\text{rf}} - \mathbf{W}_0\|_F^2 \leq \mathbf{Y}^{*\top} \hat{\mathbf{K}}^{-1} \mathbf{Y}^* \quad (*)$$

where $\hat{\mathbf{K}} \in \mathbb{R}^{n \times n}$ is a Gram matrix with entries $(\hat{\mathbf{K}})_{i,j} = (\Psi(\mathbf{X}_i)^\top \Psi(\mathbf{X}_j))$. Thus, to have (iii.b) of the right order, we need to have $(*) = \tilde{\mathcal{O}}(1)$ as $n \rightarrow \infty$, however, it is unclear how to analyze such a quantity. Arora et al. (2019) used some Reproducing kernel Hilbert space (RKHS) arguments and considered a closely related quantity $\mathbf{Y}^{*\top} \mathbf{K}^{-1} \mathbf{Y}^*$, where \mathbf{K} is a kernel matrix of a reproducing kernel function $\kappa(\mathbf{x}, \tilde{\mathbf{x}}) = \mathbb{E}[\Psi(\mathbf{x})^\top \Psi(\tilde{\mathbf{x}})]$.³ While, they have established some interesting constant-order upper bounds on $\mathbf{Y}^{*\top} \mathbf{K}^{-1} \mathbf{Y}^*$ by considering specific cases of f^* , all of the considered cases were differentiable and not necessarily Lipschitz, which violates the setting of our paper.

Acknowledgements

We thank Di Wu for drawing our attention to the issue.

References

S. Arora, S. Du, W. Hu, Z. Li, and R. Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *International Conference on Machine Learning (ICML)*, 2019.

P. L. Bartlett, A. Montanari, and A. Rakhlin. Deep learning: a statistical viewpoint. *Acta Numerica*, 2021. URL <https://arxiv.org/abs/2103.09177>. To appear.

S. S. Du, X. Zhai, B. Poczos, and A. Singh. Gradient descent provably optimizes over-parameterized neural networks. In *International Conference on Learning Representations (ICLR)*, 2018.

S. Oymak and M. Soltanolkotabi. Towards moderate overparameterization: global convergence guarantees for training shallow neural networks. *IEEE Journal on Selected Areas in Information Theory*, 2020.

J. A. Tropp. An introduction to matrix concentration inequalities. *Foundations and Trends® in Machine Learning*, 8(1-2):1–230, 2015.

Y. Yao, L. Rosasco, and A. Caponnetto. On early stopping in gradient descent learning. *Constructive Approximation*, 26(2):289–315, 2007.

3. Some standard matrix concentration inequalities (Tropp, 2015) allow to say that $(*)$ is close to $\mathbf{Y}^{*\top} \mathbf{K}^{-1} \mathbf{Y}^*$ whenever $m \geq \text{poly}(n/\lambda_{\min}(\mathbf{K}))$.