

1 We thank all reviewers for their time and understanding of the results, and we are glad that all reviewers appreciate the
2 introduction of the first CIs that feature a double-descent behavior. We first address 2 points raised by several referees
3 and then address comments by individual reviewers. All typos/writing issues will be fixed as suggested in the reviews.
4 We believe the rebuttal addresses all actionable concerns, and we invite the referees to revisit their scores in light of this.

5 • **Why does (8) hold?** By (3), the weighted average in (8) is equal to $-\text{trace}[\mathbf{T}_0(\mathbf{e}_j)]/\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2$ in Theorem 1. The
6 correctness of (8) then follows by Theorem 1. The term [additive correction] in (8) includes traces of $\mathbf{T}_1, \mathbf{T}_L$ and the
7 first term in (11). We'll clarify in the camera-ready version. **This addresses a major concern from Reviewers 2, 3, 4.**

8 • **Why vector $\boldsymbol{\xi}$ is the argument of \hat{f} in (3)?** The variable name $\boldsymbol{\xi}$ in (3) was picked to avoid confusion with the
9 observed feature vectors $(\mathbf{x}_i)_{i \in [m]}$, since the function \hat{f} itself in (3) depends implicitly on $(\mathbf{x}_i)_{i \in [m]}$ through $\hat{\boldsymbol{\alpha}}$. It is
10 thus important to avoid \mathbf{x}_i as argument for the definition of the function $\hat{f}: \mathbb{R}^p \rightarrow \mathbb{R}$ in (3); for clarity we avoided \mathbf{x} as
11 well. Applied to \mathbf{x}_i , quantity $\hat{f}(\mathbf{x}_i)$ depends on \mathbf{x}_i directly through the input given to \hat{f} and implicitly through $\hat{\boldsymbol{\alpha}}$ in the
12 definition of \hat{f} . The implicit dependence of \hat{f} on \mathbf{x}_i through $\hat{\boldsymbol{\alpha}}$ is the reason why classical CLT would not apply to
13 $n^{-1/2} \sum_i \hat{f}(\mathbf{x}_i)$. Since \hat{f} is the prediction function after training, $\hat{f}(\boldsymbol{\xi})$ can be thought of as the prediction at a point of
14 interest such as an unlabeled future observation, though this is only a thought experiment and $\boldsymbol{\xi}$ is only used as the
15 argument of \hat{f} to set the notation for partial derivatives. The quantity ξ_L in Theorem 1 is unrelated to the vector $\boldsymbol{\xi}$ in (3);
16 to avoid confusion we will clarify and use ζ_L instead. **This addresses a major concern from Reviewers 1, 2 and 3.**

17 **Reviewer 1:** *"the linear part of the target is the only thing that RF model can learn"* thanks—we'll add this comment
18 with a reference as a motivation for the β_j s being the targets of the CIs. *"(...) the CI is only established for the linear*
19 *part of the target (...)* *In this case what would be the motivation of considering a nonlinear perturbation"* \Rightarrow The
20 nonlinear perturbation adds noise correlated with the features, which is possibly more hostile than the independent
21 additive noise ε . This model of nonlinear perturbation wasn't introduced in the submission but in [Mei and Montanari
22 2019] to study the double-descent curves for the risk. It is not an ad hoc model to obtain CIs. *"I find the "double*
23 *descent" phenomenon in the CI length to be interesting.(...) Can the author comment on the plausible mechanism of this*
24 *observation?"* \Rightarrow A possible mechanism is the conjecture that the length of the CIs is an increasing function of the risk.
25 Proving this appears highly non-trivial. Such phenomena were observed for M-estimators in [Celentano and Montanari
26 2019, Prop. 4.3(iii)]. *"The result in Section 2.4 (based on Mei and Montanari 2019) seems to be under the assumption*
27 *of iid weight matrix W"* \Rightarrow Results in § 2.4 build upon results of Mei and Montanari. Hence, for that section only we
28 require \mathbf{W} with iid entries, but, e.g., Theorems 1 and 2 hold for bounded $\|\mathbf{W}\|_{op}$ with no iid assumption. *"(minor)*
29 *A few related works(...)"* \Rightarrow Thanks—we'll add those. *"(minor) Does the characterization also hold for the ridgeless*
30 *limit ($\lambda = 0$)?"* \Rightarrow Not currently. A lead to study $\lambda = 0$ is to first take $n, p \rightarrow +\infty$ (our result), then take $\lambda \rightarrow 0$ and
31 study whether interchangeability of limits applies. This appears non-trivial. *"(minor) On Figure 2 Left, why is there a*
32 *discrepancy between the predicted and simulated boxplot?"* \Rightarrow The predicted theory is the thick blue line only. The
33 discrepancy is mild and expected for some boxplots when plotting that many boxplots. We'll increase number of runs.

34 **Reviewer 2:** *"Neural Networks, when used for classification, are known to suffer from poor calibration. This paper*
35 *does not touch upon any of these issues prevalent in 2-layer neural networks."* \Rightarrow We understand the referee's concerns
36 but it does not appear reasonable to solve, in a single 8-pages submission, both CIs for regression, classification and
37 solve the mentioned calibration issues. *"The quantity nL^2 remains elusive (...)"* \Rightarrow In our results $L \asymp n^{-1/2}$, as for the
38 length of CIs in classical statistics. nL^2 is of constant order after multiplication by n . That's why we study this quantity
39 in figures. *"weighing in (8) seems to be a matter of analytic convenience (...)"* \Rightarrow This is a contribution of the paper to
40 pinpoint this specific surprising weighting that leads to a pivotal quantity in Theorem 1. This allows the construction of
41 CIs. *"It appears that the training loss appears everywhere. Is this consistent with (...)"* \Rightarrow Training loss, $\text{trace}[\mathbf{H}]$ and
42 population loss are linked through subtle nonlinear relationships studied [Mei and Montanari 2019] and our § 2.4. *"(...)*
43 *overview into the proof techniques needs to be provided"* \Rightarrow A proof outline will be added (using the 9th page allowed).

44 **Reviewer 3:** *"(...) τ is chosen to be $\tau = (d/p)\lambda$, is this scaling (by d/p) technically important?(...)"* \Rightarrow No. Since d/p
45 has a finite limit, parametrization of the tuning parameter can be performed through either τ or λ without changing the
46 conclusions. *"Eq (2) and above, notations for transpose should be aligned (...)* *needs clarification. (...)"* \Rightarrow Thanks!
47 We'll clarify the writing and fix the typos as suggested. *"(...) The connection between Eq (8) and Theorem 1 needs to*
48 *be clarified (...)"* \Rightarrow cf. line 5-7 above. *"Figure 3:(...)or choosing $\max(x, 0.01x)$ could help?"* \Rightarrow Thx! We'll use it.

49 **Reviewer 4:** *"The main issue is that this CLT holds under a linear model(...)* *The coefficients within a linear model are*
50 *not a natural quantity to be concerned with for the neural network(...)* *if we believed that a linear model held, we would*
51 *fit a linear model, not a neural net"* \Rightarrow Thanks. Developing similar CIs where $\mathbb{E}[y|\mathbf{x}]$ is a nonlinear function of \mathbf{x} is an
52 interesting future direction. One reason for starting with a linear model is that the targets for the CIs are canonically
53 defined (the unknown coefficients β_j), while for nonlinear models, it is unclear what to consider as canonical population
54 targets (functionals of $\mathbb{E}[y|\mathbf{x}]$) for the CIs. Studying linear models has been fruitful to understand double descent curves
55 in numerous recent works (cf. related work section or e.g., [Mei and Montanari]), and this submission is the first work
56 to provide CIs within this line of research. *"(...)it is confusing to see that the results address the linear models and*
57 *their betas.(...) Perhaps the authors can clear up this confusion, by answering how (8) exactly follows."* \Rightarrow cf. line 5-7
58 above. The relationship between Theorem 1 and the linear coefficients β_j s are addressed on line 187-188 around (14).