

1 *General response:*

2 We wish to thank the reviewers for their valuable feedback. We are pleased that they generally appreciated the novelty
3 of the analysis presented in this paper, which sheds light on many recent observations in the double-descent literature.
4 We also appreciate extensive and thoughtful comments regarding clarity and pointers to relevant literature that was
5 missing. We will implement all of them which we think will greatly improve the quality of the text, we also include
6 more specific discussions on the main points below (R1 R2 R3 R4).

7 We would like to start by emphasizing that the purpose of our paper is to analyze and distinguish two types of overfitting
8 which are both attracting significant interest and are often conflated with one another. From this point of view, our paper
9 is less about the existence of multiple peaks but more about (1) differentiating the sources and properties of the two
10 types of peaks and (2) investigating how they interact with the nonlinearity of the activation function (which implicitly
11 regularizes the linear peak: as suggested by R2 we will emphasize and add more details on this).

12 We believe these contributions will help understand more complex cases, especially in light of recent works which also
13 study the presence of multiple peaks. On the one hand, data distributions whose covariances have block structure can
14 give rise to multiple *linear peaks*¹. On the other hand, random feature regression with the NTK features of a two-layer
15 network displays two *nonlinear* peaks due to the block structure of the covariance of the NTK features². The work by
16 Adlam and Pennington, which appeared after submission time and was mentioned by R1, provides an answer to R2's
17 question about the possible existence of several nonlinear peaks.

18 R2 and R3 raise valid concerns about how our picture generalizes to more complex data distributions and learning
19 algorithms. Despite its simplicity, the setup we consider is rich enough to capture the two kinds of overfitting as desired.
20 As far as the data distribution is concerned, we considered the unstructured iid case in the main text because it is the
21 historical model to describe the two kinds of overfitting. Appendix C briefly discusses the case of MNIST, but the
22 phenomenology of the linear peak becomes significantly more complex in structured datasets as illustrated by the work
23 of Chen et. al. As for the impact of the learning algorithm, an unclear sentence pointed out by R4 might be underselling
24 how realistic our teacher-student framework is: by '2-layer networks', we mean 3-layer networks with 2 *hidden* layers.
25 Furthermore, we have verified that adding extra layers has little impact on the results, a point which we will further
26 stress.

27 *Specific responses:*

28 **@R1 Parameter-wise vs. sample-wise:** We are not sure to understand the first weakness stated by the reviewer. Double
29 descent was initially investigated parameter-wise but several works have subsequently studied it sample-wise. We
30 certainly agree with the reviewer that changing the width and the dataset size is not symmetric. We vividly illustrate this
31 fact by presenting a full visualization of the 2D phase-space, which highlights the role played by the input dimension.
32 **'Aren't both peaks due to small eigenvalues in $Z^T Z$?'**: The linear peak is indeed related to small eigenvalues (orange
33 curve of fig.5), but not due to strictly vanishing eigenvalues like the nonlinear peak (purple curve of fig.5). In fact,
34 the smallest eigenvalue decreases monotonously sample-wise (for $N \leq P$) reaching zero in correspondence of the
35 non-linear peak. An explanation for this counter-intuitive behavior is provided in appendix B. We will clarify this point.

36 **@R2 Bias-variance decomposition:** We understand the concerns of the reviewer about the relevance of the section
37 presenting a BV decomposition. The aim of this section was to highlight the fact that the linear peak is only caused by
38 noise variance, which is why it vanishes in absence of noise in contrast to the nonlinear peak. However we acknowledge
39 that the density of the figures may drown the message conveyed, and will consider simplifying them or moving some
40 material to the appendix.

41 **@R3 The case of MNIST:** Our contribution indeed leverages pre-existing theory, but we believe it brings along new
42 insights and will do our best to better highlight the key takeaways. Regarding the MNIST experiment, as we will explain
43 better, the absence of the linear peak is due to the very small intrinsic dimension of such a simple and structured dataset.
44 A more thorough investigation of the impact of the structure of the dataset is left for future work.

45 **@R4 Role of the nonlinearity:** We agree that the wording 'purely nonlinear' ($r = 0$) should be clarified. In the
46 high-dimensional regime we focus on, the linear part of an activation function is measured by $\zeta = \mathbb{E}_{z \sim \mathcal{N}(0,1)} [z\sigma(z)]$,
47 which vanishes for even functions such as $\sigma(x) = |x|$ but also $\sigma(x) = x^2$. The non-smoothness of generalization
48 behavior around the point of pure linearity ($r = 1$) is also an interesting question which we will expand in SM. As for
49 the divergence of the bias in Mei and Montanari's work³, it is due to the fact that the bias still contains the diverging
50 initialization variance. The latter was disentangled to yield a well-behaved bias⁴.

¹L. Chen, Y. Min, M. Belkin, and A. Karbasi. Multiple descent: Design your own generalization curve.

²B. Adlam, J. Pennington. The neural tangent kernel in high dimensions: Triple descent and a multi-scale theory of generalization.

³S. Mei and A. Montanari. The generalization error of random features regression: Precise asymptotics and double descent curve.

⁴S. d'Ascoli, M. Refinetti, G. Biroli, and F. Krzakala. Double trouble in double descent: Bias and variance(s) in the lazy regime.