

1 We appreciate all the reviewers’ valuable comments. Here is our response to the major questions raised by the reviewers.

2 **Reviewer #1:**

3 **Q:** Regarding the notation $B_{w,*}$, B_w , λ_{\min} , and $q(z)$. Missing $\sqrt{\log(Dm)}$ factor in Theorem 1.

4 **A:** We used $B_{w,*}$ to highlight it’s the radius of the ball that contains \mathbf{W}_* , and B_w to denote an arbitrary radius (later we
5 apply this result with B_w larger than $B_{w,*}$). λ_{\min} denotes the minimum eigenvalue. $q(z)$ in Line 174 is the dummy
6 variable (under the min) for a degree- k polynomial. ϵ in Algorithm 1 determines the choice of D (as in Theorem 2). We
7 appreciate these suggestions and will make a supplementary table for all the symbols in our next version. Our Theorem
8 1 does contain a $\sqrt{\log(Dm)}$ factor (complete version is provided in Appendix Line 464 with $\sqrt{\log(Dm)}$ in the first
9 term). We used $\tilde{O}(\cdot)$ to hide dependency on any log factor in the theorem statement.

10 **Reviewer #2:**

11 **Q:** What is responsible for the difference between the quadratic model and the linear models. How much does the main
12 result say about the advantages of a nonlinear NN over a linearized NN or a kernel model?

13 **A:** Our main results show that the quadratic model achieves $\tilde{O}(d^{\lceil p/2 \rceil})$ sample complexity with neural representations,
14 while the linearized model / kernel suffers from at least $\Omega(d^p)$. The key thing behind is that the generalization
15 performance of the quadratic model depends on (and can benefit from) the conditioning of the covariance of the
16 input (Line 143-147). This enables the sample complexity to be reduced when we feed it with an expressive and
17 isotropic feature map. In contrast, linearized models/neural tangent kernels cannot benefit from feature isotropicity, and
18 generalizes at most as well as a kernel, as stated in the lower bound.

19 **Reviewer #3:**

20 **Q:** How does whitening affect the proof in Section 4.1? Why not also whiten the raw input?

21 **A:** The only effect of whitening is to make the features isotropic, which does not change the expressivity of the features
22 (since it is only a linear transformation) but is beneficial to generalization, as discussed in Line 208-217. We also
23 showed that whitening is not the only option — using unwhitened features \mathbf{g} along with a proper data dependent
24 regularizer on \mathbf{W} gives us exactly the same result as whitening the features (Appendix C.5). On the other hand, existing
25 results on raw representations already assumed \mathbf{x} is exactly or nearly isotropic (Line 244 for NTK-Raw; Line 196-198
26 for Quad-Raw). Those bounds won’t be improved if we further whiten \mathbf{x} to be exactly isotropic.

27 **Q:** What properties of your random representations actually make the difference for the sample complexity here?

28 **A:** The key thing that allowed a fixed neural representation to be helpful is that the nonlinearities (along with the width)
29 give us strong expressive power. Linear combinations of these fixed neurons can already express high-complexity
30 nonlinear functions, and such expressivity can be used by the top trainable model to reduce the complexity of the
31 function it has to learn itself. In comparison, when we use the raw input, linear combinations of the input is only a
32 linear function, thus all the “heavylifting” is on the top trainable model, causing the sample complexity to be higher.
33 Therefore, our theory shows that lower-layer representations can ease the burden of learning in the upper layers, which
34 we suspect is also the case in practice even when the representation function is trainable as well.

35 **Reviewer #4:**

36 **Q:** Are our models overparameterized? How would our results change with smaller m .

37 **A:** Our model is overparameterized (we chose D, m to be large) for the purpose of approximating the ground truth
38 function and making the optimization landscape nice. However, our choice of D, m does not explicitly depend on n
39 (Line 181), making our model not necessarily wide enough to memorize the training data. In this sense we are not as
40 overparametrized as the memorization regime.

41 **Q:** It would be nice for you to discuss how your work relates to other works that construct random features.

42 **A:** Prior work e.g. of Rahimi and Recht considered training only the output layer of the network (a_r in our notation),
43 which is effectively a linear model/kernel method. In contrast our model is non-linear (quadratic) in the trainable
44 parameter \mathbf{W} and has different optimization/generalization behaviors from kernel methods.

45 **Q:** Further motivations for considering the quadratic Taylor model.

46 **A:** As one example apart from the theoretical benefits shown in this paper, empirically the (full) quadratic Taylor model
47 also approximates the training trajectories of standard neural networks better than the linearized model, as shown in Bai
48 et al. 2020 “Taylorized Training”.

49 We appreciate all the above questions and will incorporate these discussions in our next version.