1  We thank all the reviewers for the time reading our paper! We will fix all the minor issues, and below we only address
2  the main concerns. We restate those questions below.

3  • **R2:** The two-layer part of this submission might be similar to Daniely's?
4  No. Daniely's result trains only the *last layer* of (any depth) network, and the changes of hidden layers are negligible
5  (can be set zero). Daniely proves generalization using conjugate kernel, but it's unclear what *explicit functions* can
6  be learned and what's the *explicit sample complexity*. In contrast, we consider target functions consisting of 2-layer
7  networks, and show how they can be PAC-learned and what's the sample complexity.

8  • **R2:** Why the 3-layer result is *not based on NTK* (neural tangent kernel)?
9  Recall $W = W_0 + W'$ and $V = V_0 + V'$ where $W_0, V_0$ are initializations. In NTK, by ignoring higher-order terms,
10  the network is a linear function in $(W', V')$, so entries of $W'$ will *never* be multiplied with $V'$. In contrast, see
11  Lemma 6.10, we track $DV'DW'$ so $V'$ and $W'$ are multiplied together. This is non-linear so is not NTK. In fact, as
12  we explained in the paper, our concept class (learnable by three-layer networks) is not captured by the NTK of a
13  three layer network. We shall make it more clear in the revision. (We thank **R5** for carefully reading our paper and
14  acknowledge "unlike prior work" we have considered "non-convex interactions between weight matrices.")

15  • **R2:** Does the sample complexity really depend polylog in the network size?
16  Sorry for the confusion. What we mean is sample complexity *grows* polylog in $m$. Indeed, the sample complexity
17  shares some poly terms with $M_0$, but as $m \geq M_0$, sample complexity only *grows* polylog in $m$.

18  • **R2+R4:** Why can't standard SGD (without noise) work?
19  There are many reasons (see footnotes on Page 7). For instance, *all* existing escape-saddle point papers need noise.
20  Since we rely on such existing work, we need noise for theoretical purpose.

21  • **R4:** Why increasing $m$ supports more target functions? What is $R$?
22  Sorry for the confusion and it is actually simpler than you may have thought. For any target function, there is
23  a corresponding $M_0$ such that our theorem applies whenever $m \geq M_0$. Here, $M_0$ depends on the complexity
24  notion introduced on line 122. In other words, if we increase $M_0$, there will be more functions to be supported by
25  this threshold $M_0$. Finally, $R$ only plays some role in our 3-layer theorem, where it allows the complexity to be
26  composed with another complexity function.

27  • **R4:** Modify references to e.g. Table 1 in line 239. Great suggestion and we will do that!

28  • **R5** raises concerns about the significance of the generalization part of this paper.
29  Although the generalization lemmas *on their own* are simple, they are not our main contribution (and constitute
30  5% of this paper). **Instead**, our main contribution is to make the convergence theorems and generalization lemmas
31  *compatible*: SGD can find solutions with small norms so that generalization lemmas apply (almost independent of
32  $m$). This cannot be done by combining any "convergence theorem" and any "generalization lemma". For instance,
33  the prior work Allen-Zhu et al [2] proves convergence, but it is not compatible with our generalization lemmas.

34  • **R5** has concerns about the practical relevance of our generalization lemmas.
35  – **R5**: What will happen if learning rate is independent of the parameter count?
36  Under our (wlog.) choice of initialization $W, V \sim \mathcal{N}(0, 1/m)$ and output layer $\mathcal{N}(0, 1)$, it is ***not*** **a good idea**
37  to use constant learning rate. For instance, in our 3-layer experiment below, any learning rate $lr \in [0.01, 0.1]$
38  works for $m = 50$, but any learning rate $lr \in (0.01, \infty)$ gives NaN error for $m \geq 5000$. In general, learning
39  rate *depends on initialization*: if hidden weights are scaled up and output layer is scaled down, then the learning
40  rate will increase.
41  – **R5:** If learning rate decreases as $m$, does $\|W'\|$ and $\|V'\|$ decrease experimentally?
42  **Yes.** For 2-layer network on MNIST with $lr = 400/m$ and target accuracy 95%, norm $\|W'\|_F / \|W_0\|_F$
43  decreases, see Li-Liang [30, Fig 5 of page 26 of NeurIPS camera ready]. Below in left figure, we show $\|W'\|_{2,\infty}$
44  also decreases. As another example with 3-layer networks, say inputs are random $\|x\|_2 = 1$ and the true label
45  $y = x_1 \cdot x_2 + x_3 \cdot x_4$. Using our standard initialization, then $\|W'\|_F, \|W'\|_{2,4}, \|V'\|_F$ all decrease as $m$ increases
46  (for $lr = 1/m$, test errors drop below 0.003 within 30 epochs), see right figure. We have plugged them into the
47  generalization formula and they do give very meaningful bounds on sample complexity even for large $m$.