

1 We thank the reviewers for their efforts. Below we address the main comments.

2 **Reviewer #1**

3 The reviewer asked how the techniques in this paper differ from the prior work that established the hardness of neural
4 networks under the RSAT assumption. First, we show that learning neural networks is hard already for networks with
5 a special structure that we call sign-CNN. It requires a reduction from the RSAT problem that is different from the
6 reduction used in prior work. We do use in this reduction a Lemma from [17], so we do not repeat the work that has
7 been done there. Second, we show that the properties of sign-CNNs allow us to transform an unknown sign-CNN to a
8 random network, by a multiplication with a random matrix that has a special structure. This method is new, and does
9 not appear in prior work. Finally, in order to bound the support of the input distribution, we need to analyze the singular
10 values of our special-structure random matrices.

11 The reviewer also commented on the presentation. We did make much effort to present the results and proofs in a
12 simple and clear way. In the "proof ideas" and "proof structure" sections we give a high-level overview of the proofs. In
13 the appendix, we first establish some needed lemmas, and then each theorem is proved in a separate subsection. Since
14 all theorems require common lemmas and constructions, skipping directly from the "results" section to the subsection
15 in the appendix where the theorem is proved is not possible. Also, for the same reason, giving a sketch of each proof in
16 the "results" section is not possible. Instead, the "proof ideas" section is essentially a sketch of the proofs.

17 **Reviewer #2**

18 The reviewer asked about the worst-case nature of the input distribution. Prior hardness results for learning neural
19 networks assume that both the input distribution and the weights are worst-case. We show that the problem is hard
20 already for simple networks with natural weights, but the worst-case nature of the input distribution remains. Thus,
21 we show that even very strong assumptions on the network are not sufficient for efficient learning, and therefore that
22 assumptions on the input distribution are necessary. We believe that positive results on the learnability of neural
23 networks would require a combination of assumptions on the weights and on the input distribution.

24 **Reviewer #4**

25 The reviewer raised three concerns:

- 26 1. "The assumption that the adversary can choose a worst-case input distribution is too strong and impractical":
27 The standard PAC-learning framework requires an algorithm that learns successfully for every input distribution
28 and every hypothesis in the class. Hence, the assumption that the adversary can choose a worst-case input
29 distribution is the standard assumption in PAC learning. While prior hardness results for learning neural
30 networks assume worst-case input distribution and worst-case weights, we show that the problem is hard
31 already for networks with natural weights (but the worst-case nature of the input distribution remains). The
32 practical importance of our result, is that it suggests that in order to establish efficient algorithms for learning
33 neural networks, or to show polynomial-time guarantees for existing algorithms, assumptions on the network's
34 architecture and weights are not sufficient, and assumptions on the input distribution are necessary. Hence,
35 our result may help focussing the efforts on directions where positive results are possible. We believe that
36 a combination of assumptions on the weights and on the input distribution is necessary in order to obtain
37 positive results.
- 38 2. "The weights on the second layer are all 1. This assumption breaks the "naturalness" of the network. Is this
39 essential or can it be generalized?":
40 The choice to focus on networks where the weights in the second layer are all 1 is not essential. Our results
41 hold also where the weights in the second layer have other fixed values, and also in the case where they are
42 random. We may add a remark on this point.
- 43 3. "The result only holds if the hidden dimension $m = O(\log^2(n))$, which makes the result less interesting. Is
44 there a way to prove similar results with more hidden nodes?":
45 The results do extend to more hidden nodes. As we mention in page 4 (lines 168-170), the results hold (with
46 minor changes) for any $m = \omega(\log(n))$, including networks with many hidden neurons. We chose to focus on
47 $m = O(\log^2(n))$ since we wanted to show hardness of learning already where the number of hidden neurons
48 is relatively small.

49 **Reviewer #5**

50 We thank the reviewer for his/her comments and will fix accordingly.