

1 We thank the reviewers for taking the time to read and comment on our submission despite these challenging times.

2 **Reviewer #1.** You wrote: “This work studies an important problem . . . The paper is mostly clearly written and most  
3 of the proofs in the paper seem correct.” Thanks for acknowledging the strengths of our submission.

- 4 • *Regarding major technical mistake.* There is no mistake here. We are simply using the fact that the zero-one loss is  
5 bounded from above by the logistic loss, and taking expectations on both sides. This is what gives the inequality after  
6 line 259. What it states is that the expected zero-one loss, i.e., the probability of misclassification, is bounded by  
7 expected logistic loss; there is no *training error* involved in this expression. A detailed derivation is in the Appendix  
8 (look for Proof of Theorem 4.3 on lines 559-561).
- 9 • *No generalization analysis.* Theorems 4.3 and 4.6 are our generalization error bounds. Instead of analyzing the  
10 Rademacher complexity, we use a martingale concentration inequality to bound the generalization error. Please see  
11 the Proof of Theorem 4.6. in the appendix for more details.
- 12 • Our generalization guarantees do not *worsen* with the gradient updates, they *improve*. In both Theorem 4.3 as well as  
13 Theorem 4.6, we bound the generalization error by  $\tilde{O}(1/T)$ , which implies that the generalization error decreases  
14 with  $T$ , i.e., with the number of gradient updates in Algorithm 1.
- 15 • The constant strategy  $W_t = W_1$  would incur a constant error at each iteration, and thus not decay with  $T$ . In sharp  
16 contrast, the error of the iterates of Algorithm 1 decays as  $\tilde{O}(1/T)$ . The advantage of Algorithm 1 over the constant  
17 strategy shows up when we apply Lemma 5.4 to the right hand side of the inequality given by Lemma 5.3.
- 18 • Thanks for your comment on Nagarajan and Kolter – we will make the *uniform convergence* discussion explicit in the  
19 final version of our paper to better convey the following quote from their paper: “. . . uniform convergence provably  
20 cannot “explain generalization” – even if we take into account the implicit bias of GD to the fullest extent possible.”
- 21 • Regarding your minor comments, 1) we will extend the language in Remark 4.5 to emphasize on computational  
22 learning theoretic advantages of dropout over SGD; 2) we will make it more clear and precise; 3-4) thanks for  
23 catching the typos,  $p$  should be  $q$ .

24 We hope that our response has clarified all of your technical concerns, in which case, please reconsider your score.

25 **Reviewer #2.** Thanks for acknowledging the novelty of our results.

- 26 • The definition of logistic loss on line 123 is not wrong; feel free to refer to any ML book. You can also refer to the  
27 prior work, e.g., [Cao and Gu, 2019], [Ji and Telgarsky, 2019], etc. It is typical to define a loss function as  $\ell(y, f(x))$   
28 and overload the notation to write  $\ell(z) = \ell(yf(x))$ ; e.g., for logistic loss this would be  $\log(1 + \exp(-yf(x)))$ .  
29 Please see the definition of hinge loss for another example.
- 30 • As we state in Remark 4.4, with high probability, dropout iterates will not violate the max-norm constraints, i.e.,  
31 these constraints are not active in a typical run of dropout whatsoever. We only need the max-norm constraints to  
32 give generalization error bounds in expectation.
- 33 • Yes,  $k \in \{1, \dots, T\}$  represents an iterate. Linearization means we approximate a function with its best linear  
34 approximation, given in terms of the gradient of the function at that point. This is what is happening on lines 167-169.  
35 Neural tangent kernel (NTK) analysis is based on using  $\nabla g_t$  as features; see the preliminaries on lines 141 – 151.
- Proof of basic fact on line 174:

$$\langle \nabla g_t(W_t), W_t \rangle = \sum_{r=1}^m \langle \frac{1}{\sqrt{m}} a_r b_{r,t} \mathbb{I}\{w_{r,t}^\top x_t \geq 0\} x_t, w_{r,t} \rangle = \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r b_{r,t} \sigma(w_{r,t}^\top x_t) = \frac{1}{\sqrt{m}} a^\top B_t \sigma(W_t x_t) = g_t(W_t)$$

- 36 • The weights returned by the algorithm are scaled by  $q$  to account for the dropout. Hence, the risk takes the argument  
37  $qW_t$ . We make the following remark on lines 139 – 140: “at test time, the weights are multiplied by  $q$  so as to make  
38 sure that the output at test time is on par with the expected output at training time.”

39 We hope that our response has clarified all of your technical concerns, in which case, please reconsider your score.

40 **Reviewer #3.** Thanks for your encouraging comments! We define the *neural tangent feature*  $\phi_x : z \mapsto x \mathbb{I}\{z^\top x > 0\}$   
41 on line 150 after Eq. 1. Assumption 1 is simply a margin assumption in the RKHS of the NTK. It is a mild and  
42 reasonable assumption that the data becomes linearly separable after mapping it into a high-dimensional non-linear  
43 feature space; this idea has been the cornerstone of kernel methods using the RBF kernel, for example, and for learning  
44 with neural networks. Further, this assumption yields bounds under mild over-parametrization compared to other works.

45 We hope that our response has clarified all of your technical concerns, in which case, please reconsider your score.

46 **Reviewer #4** Thanks for your positive review!