

1 We thank all the reviewers for their comments.

2 **Reviewer #1.** Regarding the relationship between the query lower bound, the tolerance  $\tau$  and the error  $\epsilon$ , note first  
3 that we must have  $\tau \lesssim \epsilon \lesssim \beta$  (where  $\beta$  is a norm lower bound for the concept class), as we discuss in Appendix A.  
4 Moreover, because of technical requirements on  $\tau$  for the DKKZ20 result, we end up picking  $\tau$  as a function of  $\epsilon$  in our  
5 reduction. So as  $\epsilon$  decreases,  $\tau$  decreases as well, and we get a series of incomparable (though still exponential) bounds  
6 due to the tradeoffs between query complexity and tolerance. We will include a clearer discussion of these issues in  
7 Appendix A. We will also address your additional feedback (including expanding on lines 32-34) when revising the  
8 manuscript.

9 **Reviewer #4.** (1) A major goal in deep learning is to find provably efficient algorithms for learning classes of neural  
10 networks. There are several papers a year on this topic. Here we are showing for the first time that if there is noise in  
11 the labels, this goal is impossible even for the simplest networks, ones that consist of a single activation. As we mention  
12 in the paper, our results rule out polynomial-time algorithms for *any* nonpolynomial activation (halfspaces and ReLU  
13 are just examples). The activation can take Boolean or real-valued outputs. Since algorithms for learning polynomial  
14 activations are known, this characterizes the computational complexity of learning single activations.

15 (2) Regarding the context for Theorems 1, 2, 3, we note that these lower bounds hold broadly for all statistical query  
16 algorithms and hold regardless of the structure of the learner’s output hypothesis. As for the tolerance, much of its  
17 significance lies in capturing what in traditional PAC algorithms would be the sample complexity. Specifically, it takes  
18  $\Theta(1/\tau^2)$  samples to simulate an SQ of tolerance  $\tau$ , and this is sometimes known as the *estimation complexity* of an SQ  
19 algorithm.

20 As for a comparison with prior results such as KKMS08 and DKN10 for halfspaces, our SQ lower bound (like all SQ  
21 lower bounds) states that any SQ algorithm must use *either*  $\exp(n)$  queries *or* very small ( $n^{-1/\epsilon}$ ) tolerance (which  
22 corresponds to sample complexity). KKMS08 (with its sample complexity of  $n^{1/\epsilon^2}$ ) falls into the latter category, and  
23 our tolerance bounds show that this is nearly optimal.

24 (3) Regarding your comments on the difficulty of reading the paper without prior background in the SQ literature, we  
25 acknowledge some of your points on where the writing could be improved. It is true that Corollary 2.4 is somewhat  
26 confusing. The “accompanying norm lower bound” refers to a  $\beta$  such that  $\|g\| \geq \beta$  for all  $g \in \mathcal{G}$ , as used in Theorem  
27 2.2. Part (a) of the corollary ((b) and (c) are similar) follows from “instantiating” the generic construction  $\mathcal{G}$  of Theorem  
28 2.3 with  $\phi = \text{ReLU}$  to obtain say  $\mathcal{G}_{\text{ReLU}}$ , noting that functions in  $\mathcal{G}_{\text{ReLU}}$  satisfy a norm lower bound of  $\beta = \Omega(1/k^2)$   
29 (with proofs deferred to Appendix D), and then using  $\mathcal{C} = \mathcal{G}_{\text{ReLU}}$  in Theorem 2.2 to obtain the final lower bound. As  
30 for Theorem 5.1, we do state “suppose that Assumption 3.1 holds for  $\mathcal{H}_{\text{ReLU}}$ ” to try to be as clear as possible. As  
31 for Theorem 1.1 (and 1.2 and 1.3), they are proved in Section 5, and we will note this. We will also try to make the  
32 organization of the lemmas clearer.

33 Regarding additional feedback:

- 34 • The concept class in Def 2.1 can be either real-valued or Boolean.
- 35 • As defined in lines 124-126,  $D_{f^*}$  in this case refers to the unique distribution on  $\mathbb{R}^n \times \{\pm 1\}$  defined by the  
36 (Boolean)  $p$ -concept  $f^*$ .
- 37 • One important reason for using Frank–Wolfe is because it avoids a projection step, as would be required in  
38 say standard projected GD. In our  $L^2$  function space, it is not natural to require the base learner to find such  
39 a projection of the functional gradient onto  $\text{conv}(\mathcal{H})$ . Another important reason is that Frank–Wolfe uses a  
40 *linear* optimization subproblem, and this is important in order to preserve “SQness” during the reduction (see  
41 lines 215-227 and 242-245).

42 We will address all these points when revising the manuscript.