

1 We thank the reviewers for their thorough and positive reviews. Overall, we were glad to see the reviewers found the
2 paper to be clearly written, technically sound, and the results to be of interest to the (fair) ML community.

3 We will of course incorporate all the suggested edits by the reviewers as well as more clarifications.

4 **General comments:**

5 Regarding our lower bound: Our lower bound holds both for finite predictor classes and infinite predictor classes with
6 finite graph dimension. We will restate the theorem statement so it would state precisely what is proven. Indeed, it
7 remains to be an open problem whether it can be improved, and in particular exhibit a dependency on graph dimension,
8 which we assume is necessary.

9 As for other complexity measures such as fat-shattering and margin classifiers, it could be an interesting direction for
10 future work. In this paper, we chose to derive the generalization bounds using Graph dimension and VC-dimension.

11 **Reviewer 1:** Thank you for your positive review. The reviewer expressed concern regarding the brevity of the discussion
12 about the potential societal impacts of the results. We will address these important issues in detail in the final version.

13 **Reviewer 2:** Thanks for your careful reading and constructive questions. We address your concerns as follows:

14 1. The upper bounds of [Hebert-Johnson et al.] are implicitly guaranteed only for “interesting categories” with
15 large enough mass.

16 2. *Dependence on ψ :* In the case of finite predictor classes, we applied union-bound directly on all the (large
17 enough) “interesting categories”, that have finite cardinality due to the finite number of hypotheses. In the
18 case of infinite predictor classes, this technique is no longer possible. To derive the concentration bound of
19 $c(h, U, v)$ we bound the deviations of $p_1(h, U, v)$ and $p_2(h, U, v)$ (see line 592, in the proof of Theorem 10 in
20 the supplementary material). To achieve our desired accuracy, we need to bound the deviation by $\psi\epsilon/3$ (see
21 Lemma 19, line 331). This increased accuracy translates to the ψ^{-2} in the sample size.

22 3. *Instance optimality:* Your intuition is correct, there are some problem instances with smaller VCs for all the
23 values of v that determines the “interesting categories”, thus their sample complexity is lower. On the other
24 hand, the bound is tight in the worst case. Namely, for any graph dimension d , we can consider a predictor
25 classes with $VC(\mathcal{H}_v) = VC(\Phi_{\mathcal{H}_v}) = d$ for any prediction value v .

26 4. *Applicability of the proposed techniques:* Our techniques are suitable to derive generalization bounds for
27 complex measures such as F-score. Specifically for F-score, one can derive generalization bounds, assuming
28 that at least one of TP, FP, or FN is “not negligible”.

29 5. *Objectives which contain both a loss component as well as a calibration error component* Our generalization
30 bounds extend to such settings. For example, for zero-one loss, we can guarantee simultaneously for each
31 predictor in the class that all the empirical multicalibration errors of “interesting categories” and the empirical
32 loss are close to their expectations. This does not even require an increase in the sample complexity.

33 **Reviewer 3:** We thank the reviewer for their positive and constructive review. As for multi-category pattern classifica-
34 tion, it can be reduced to binary classification, which we address in this work.

35 **Reviewer 4:** Thank you for diving deeply into the paper and for your detailed comments and corrections. Sure, one can
36 choose α to be the discretization scale, but this is a limitation we wanted to avoid in order to give more flexibility.