

1 Thank you for your reviews! Your comments will be useful in revising our paper. We appreciate the positive
2 comments, e.g., that our work addresses an “under-studied”/“timely” problem “of extreme importance,” the method
3 is “nice”/“simple”/“clean”/“well validated,” the paper is “(exceptionally) clear”/“well written,” results are “solid”/
4 “significant”/“clear,” with “significant performance gains.”

5 We briefly summarize the different *major* concerns of the 3 reviewers before we respond to them separately. We note
6 that **there were no shared weaknesses pointed out by the 3 reviewers**. The main concern of @R1 (overall score=4)
7 is about the novelty and contribution of the paper, especially in the context of a specific recent line of work by Foulds et
8 al. The reviewer also viewed focusing on assessment of fairness instead of learning fair models as a weakness. @R2
9 (overall score=7) characterizes our work as being a useful approach to quantifying and reducing uncertainty in fairness
10 metrics, with broad applicability and significant performance gains, and “does not think the weaknesses of our work
11 is consequential enough to prevent publication”. @R4 (overall score=4) identified some potential weaknesses (prior
12 sensitivity, desirability of more theoretical results) but did not identify major issues with the paper.

13 **@R1: “in light of Foulds et al 2019, the contribution is rather minor”** Thank you for pointing out this paper. We
14 will certainly cite and discuss it in our revised paper since it shares a common starting point with our work of using
15 Bayesian methods to assess fairness. *However, we disagree that our contribution is minor in light of this paper. The two*
16 *papers complement each other: they address different problems and take different technical approaches.* We believe our
17 work is substantially different from [1] in objective, methods and results and we ask Reviewer 1 to reconsider our paper
18 in this light. **Objective:** our work specifically focuses on the problem of *leveraging unlabeled data to generate better*
19 *estimates of fairness metrics* given limited labeled data; in contrast, [1] focuses on assessing *intersectional fairness*
20 when the amount of labeled data is extremely small and unlabeled data is unavailable. We also make a point in our work
21 of emphasizing that estimates of fairness metrics can suffer from high variance even in the presence of relatively large
22 amounts of labeled data (see Fig 1). **Methods:** We use Bayesian methods to calibrate scores for unlabeled datapoints
23 for improved estimation of group fairness metrics, whereas [1] uses Bayesian methods to provide parametric smoothing
24 among groups for improved estimation of intersectional fairness metrics based on labeled data. **Results:** We evaluate
25 our methods across a broad range of different datasets and different models (mechanisms) and demonstrate substantial
26 improvements in accuracy of estimates of multiple group fairness metrics, whereas [1]’s results are focused mainly
27 about the accuracy ratio between intersectional groups on two (semi-synthetic) datasets.

28 **@R1: “focuses only on estimation and not on how to obtain fair policies”**. We agree that obtaining fair policies
29 is an important problem, but we also believe that for a given model, trained fairly or not, *independent and accurate*
30 *assessment of its fairness is important and under-studied* (also the motivation in [1]), particularly when users only have
31 access to a blackbox predictor. We will make this point clearer in revision. **@R1: “seems limited to fairness settings**
32 **where Bayesian calibration is applicable”**: We emphasize that our approach is applicable to any classification setting
33 and we **do not need** any special setup to apply our method. The only requirement is that there is both labeled and
34 unlabeled data available from the deployment environment. @R1: notation in 1.105-107: We agree this notation is
35 confusing and will remove this equation (its not needed). **@R1: “discuss how you would extend this idea to other**
36 **fairness metrics.”** Good point, we agree. We can directly extend our approach to handle metrics such as calibration and
37 balance as well as ratio-based metrics and we will this discussion of such extensions to the paper. All fairness metrics
38 which are defined as deterministic functions of model score S , label Y and sensitive attribute A (for concreteness we
39 demonstrated with 3 popular fairness metrics in the paper) can be approximated on unlabeled data with our proposed
40 method.

41 **@R2: “the possibility that the CIs might be overconfident if there are many labeled examples.”** Good point. On
42 page 13 of Appendix, we empirically validated that our method provides reasonably well-calibrated CIs, but we believe
43 there is room for further improvement in this area. **@R2: “the challenge of balancing the bias-variance tradeoff for**
44 **this method.”** We agree that this is an interesting direction for future work. We provide some theoretical considerations
45 in lines 165-177 and also acknowledge this issue in a brief discussion in lines 242-250: but there is certainly room for
46 more work on this front. **@R2:Eskimo -> Inuit:** Thank you for spotting this! We will update it in the paper.

47 **@R4: “proper sensitivity analysis”**: We agree that systematic sensitivity analysis is lacking and we will add it to the
48 Supplement. The values for the priors were selected based on consideration on knowledge of ranges of miscalibration
49 one typically sees with trained classifiers—and we also found that the same priors worked well across a large range of
50 datasets and models without any need for tuning (see also lines 153-156). **@R4: “non-hierarchical variance:** Thanks
51 for suggesting the ablation study, we will add it to the Supplement. **@R4: Equation 1:** The phrase “For example”
52 was loosely worded and we will remove it in the revision to avoid potential confusion. **@R4: “more theoretical**
53 **results to establish the relevance”**: We agree that more theoretical results are important going forward—please see
54 our response to R2 on this point. **@R4: “posterior computations contain no novelties or contributions”**: We use
55 standard MCMC for posterior computation and it works well in our experiments in terms of both accuracy and runtime.