1  We thank the reviewers for their feedback.

2  **Reviewer 1**

3  Thank you for the encouraging words. We have tried as much as possible to provide insight on what xAUC discrepancies
4  mean and to provide diagnostic tools to inspect them in order to inform practice. We agree further work is needed to
5  push this in specific applications. While this is beyond the scope of this present paper, we have employed xAUC/xROC
6  metrics to assess fairness in more recent work and will continue to push this direction in application work.

7  **Reviewer 2**

8  We appreciate the kind words and your expectation for xAUC to be widely adopted. And, yes, thank you: we will
9  mention the connection to separation (separation for a score implies but is strictly stronger than $\Delta$ xAUC $= 0$; note
10 also that $\Delta$ xAUC is a metric for assessing level of unfairness rather than just a criterion). Re code dependencies: we
11 will update the checklist; we use Python 2.7 with numpy/pandas/sklearn/matplotlib/scipy.

12 **Reviewer 3**

13 "fairness literature is awash with fairness metrics [...] make a compelling case for why we should prefer their metric to
14 the existing alternatives."

15 - Firstly, there are no metrics specifically for **disparate impacts of continuous risk probability scores**, as
16   acknowledged by R1. We distinguish between conditional probability estimation with observed binary
17   outcomes, and fair regression which studies parity constraints on loss with respect to observed regression
18   outcomes (Agarwal et al, 2019; Zink and Rose 2019).

19 - We specifically featured the COMPAS example, a commonly studied dataset, to make a compelling case for
20   this metric: attention to xAUC, its probabilistic interpretation, and its decompositions, illuminates causes of
21   unfairness. We highlighted that previous arguments (e.g., Dieterich et al. [18]) which compared within-group
22   AUCs did not provide a meaningful notion of accuracy equity in decision-support settings, so that **relative to
23   such naive metrics that exist, xAUC should be preferred**. We further featured other datasets to demonstrate
24   the breadth of applicability.

25 Re: base rates: In Proposition 1, we provide a decomposition that illustrates the connection between base rates, xAUC
26 disparities, and marginal AUC. This shows that base rates are just one ingredient and disparities can exist even with
27 equal base rates.

28 Re: what action should be taken after xAUC is measured:

29 - The legitimacy of direct post-training adjustment is highly contested both in the fairness literature and in
30   practice and is highly context-dependent, and this is in no way specific to the xAUC. Part of this ambiguity
31   relates to the overloaded use of "predictive risk scores" in the fairness literature, which our work in part seeks
32   to clarify. We highlight that different interpretations may be more or less directly applicable based on the
33   problem setting. If the setting does not admit adjustment (e.g. risk assessment in healthcare to inform patient
34   decisions), we highlight work that discusses alternative approaches for directly improving performance, such
35   as choice of covariates. When these avenues are available, it is unclear if adjustment is advisable.

36 - What action should be taken is highly context-dependent, and we illustrate how different contexts lead to
37   different recommendations. Recognizing that the issue of direct adjustment (we believe) remains contested
38   in the fairness literature and for practitioners, due to legal and Pareto-violation concerns, we (or anyone)
39   wouldn't be able to provide an honest and realistic algorithmic solution to address xAUC disparities. See e.g.
40   Hellman 2019, "Measuring Algorithmic Unfairness", which seeks to establish general principles for preferring
41   calibration vs. impact-based or base-rate sensitive metrics.

42 - Nonetheless, we relate the idea of post-processing adjustment of scores to post-processing classifiers by
43   adjusting thresholds. Furthermore, the results we include in the supplement do provide the tools for $\Delta$ xAUC-
44   minimizing adjustment of scores for those who would choose to do so.