

1 Thank you all for the careful reviews and useful feedback.

2 **Typos:** Thank you for pointing these out: they’ve been fixed.

3 **More general metrics:** The “right” metrics to use depend on the application. Every metric we consider is an expected
4 loss over a subset of the instance space, and while that covers a surprisingly wide variety of real-world problems (see
5 e.g. references (1) and (2), below), we of course agree that an even more general setting would be superior. There are a
6 lot of possibilities, so we think that is a fruitful area for future work. One possible extension, for example, would be to
7 consider smooth functions of such aggregate rate metrics, which would cover cases such as the F-score or G-mean
8 (Appendix A.3, in fact, includes some preliminary experiments on the G-mean metric). Or one could look at metrics
9 over *pairs* of instances, which would cover ROC AUC and variants such as Pinned AUC (see reference (3), below).

10 **Statistical fairness metrics and disorderliness:** Section 3 is meant as a discussion section, and our intention, when
11 we introduce “orderliness”, is to provide an intuitive framework for how one should think about the relative pros and
12 cons of the approaches that we consider. We deliberately do not give a strict mathematical definition, but if we did, we
13 could do so only for classifiers that are based on the idea of subdividing the space, such as thresholding each subdivision
14 at a different threshold (Section 2), or applying a different ensemble element on each subdivision (Section 4). In the
15 subdivision context, a more orderly classifier would have larger subdivisions, and a less orderly classifier would have
16 smaller ones. As we show in Theorem 3, a highly disorderly hashing classifier—i.e. one based on sufficiently small
17 subdivisions—will, with high probability, perform well w.r.t. any m aggregate rate metrics. In other words, while
18 Reviewer #3 is correct that disorderliness is not *necessary* for group fairness, it is *sufficient*, at least in the context of
19 Section 2.3. Our experiments explore this trade-off.

20 Conversely, while a particular orderly deterministic classifier *could* perform well w.r.t. group metrics, it won’t
21 *necessarily* perform well. So, if your only desire is to be confident that your classifier will perform well w.r.t. group
22 metrics, then you should generally favor a hashing classifier that is more disorderly.

23 An open question is how to construct more orderly classifiers that also have guarantees w.r.t group metrics. We will
24 clarify these points in the paper.

25 **Guarantees compared with lower bound:** As we discuss on lines 142-151, what our tightest upper bound (Theorem
26 3) and our lower bound (Theorem 1) have in common is that they both go to zero as the amount of stochasticity on large
27 point masses goes to zero, but the way that the two bounds measure this quantity differs, so there is an opportunity
28 to further close the gap (here and throughout our paper, we view this paper as opening the discussion of these issues,
29 rather than providing the last word).

30 A more precise derivation of the differences between the two bounds can be found In Appendix B.4 (entitled “sanity
31 check”), in which we progressively lower-bound Theorem 3 to verify that it does indeed upper-bound Theorem 1. This
32 appendix (which we point to in a footnote on Page 4) was initially written for our own peace of mind, but it’s useful in
33 that it explicitly lists the steps required to reduce the upper bound of Theorem 3 to the lower bound of Theorem 1.

34 **Significance:** This work was motivated by our realization (and sometimes frustration) that in some parts of industry—
35 and this is admittedly anecdotal—a *stochastic* classifier will often be rejected out of hand by engineers or their
36 managers, regardless of its performance. Part of the contribution of this paper is in studying why practitioners are
37 often uncomfortable with stochastic classifiers, and then digging into those issues to understand whether and how
38 deterministic approximations can address these concerns. We show that there are indeed definite practical downsides
39 to stochasticity, but that the practitioner’s usual default choice of a thresholded deterministic classifier (Section 2.2)
40 does not enjoy the guarantees that theoreticians have worked so hard to prove, and that such a classifier often (as we
41 show in our experiments) performs worse in practice than the original stochastic classifier. We further show that this
42 standard technique for converting a stochastic classifier to a deterministic one (thresholding) does not work as well as
43 hashing (compare Theorem 2 with Theorem 3). Thus this paper provides important first theoretical analyses of what
44 practitioners generally do, and what they could be better-off doing instead, when a stochastic classifier is unacceptable.

45 References

46 [1] Goh, Cotter, Gupta and Friedlander. “Satisfying real-world goals with dataset constraints”. NIPS, 2016.

47 [2] Narasimhan. “Learning with Complex Loss Functions and Constraints”. AISTATS, 2018.

48 [3] Dixon, Li, Sorensen, Thain and Vasserman. “Measuring and Mitigating Unintended Bias in Text Classification”.
49 AIES, 2018.