

1 Thanks for the valuable comments! Responses to each reviewer follow.

2 **R1:**

3 > Restrictive definition of fairness

4 Our results in fact hold for most widely used group-based fairness definitions, including: (i) demographic parity and
5 equality of opportunity (per Lines 120–124), (ii) equalized odds (by adding an additional constraint for $Y = 0$ in Λ_D^{EO}),
6 and (iii) disparate mistreatment, also known as accuracy parity (by using the 0-1 loss and Λ_D^{DP}). Thus, we do not believe
7 our results are narrow in scope.

8 > Non-binary sensitive features

9 Our theorem can be generalized to non-binary sensitive features, provided one makes stronger assumptions on the noise
10 (e.g., that it is symmetric). Further developing this is certainly of interest, but studying a binary feature is a common
11 starting point in both the fairness and label noise literature. As ours is the first study of the issue of noisy sensitive
12 features (to our knowledge), we focused on getting a clean and practical result in this important case.

13 > MC noise model is restrictive

14 MC learning is an active, widely-used noise model (e.g., [1, 2, 3, 4, 5]). By ensuring that our technique works for this
15 noise model, we have covered the two important and pervasive special cases of CCN [3, 1, 5] and PU learning [2, 4].

16 > Constraint $\alpha + \beta < 1$ in MC model is restrictive

17 The constraint imposes no loss of generality: when $\alpha + \beta > 1$, we can simply flip the two labels and apply our theorem.
18 When $\alpha + \beta = 1$, all information about the sensitive attribute is lost. This pathological case is equivalent to not
19 measuring the sensitive attribute at all.

20 > Unclear/undefined notation

21 The order of equations in (4) was accidentally swapped (our apologies). For all the other points raised: (i) Λ_D is defined
22 in the immediately following para (Lines 111-115) and examples are given in Equations 2 and 3; (ii) \bar{L} is defined on
23 Line 118, and $\bar{\ell}$ immediately before on Line 117; (iii) regarding redefinition of L_D , in both cases is the same quantity;
24 the RHS in the second usage explicates the involvement of the sensitive attribute.

25 > (c) “In equation (5) . . . given that α and β don’t sum up to 1, what happens to the remaining data points?”

26 There seems to be a slight misunderstanding: the mixture weights need only sum to 1 *within* and not *across* each
27 corrupted class-conditional. One must take $\Pr(Y_{\text{corr}} = 1)$ into account when reasoning about the samples. As a concrete
28 example, consider a CCN setup where $\Pr(Y = 1) = \frac{1}{2}$, and +ve and –ves have a 0% and 50% chance respectively of
29 having their label flipped. One may verify that in this case, $D_{1,\text{corr}} = \frac{2}{3} \cdot D_1 + \frac{1}{3} \cdot D_0$ and $D_{0,\text{corr}} = D_0$. Clearly, here
30 $\alpha + \beta = \frac{1}{3} \neq 1$. The “missing” weight on D_1 is compensated by there being a greater fraction of corrupted +ves, as one
31 can verify $\Pr(Y_{\text{corr}} = 1) = \frac{3}{4}$. Indeed, the fraction of true +ves remains at $\frac{2}{3} \cdot \frac{3}{4} = \frac{1}{2}$, and so no sample goes missing.

32 **R2:**

33 > Post-processing of Jagielski et al. can be applied to any method, and can handle demographic parity

34 It is correct that the post-processing method of Jagielski et al. can be applied after any fairness-unaware learner. By
35 contrast, our method can be applied to any *in-process* fairness-preserving learner (e.g. [Donini et al., 2018], [Agarwal
36 et al., 2018], [Zafar et al., 2017b]) In-processing algorithms generally result in better tradeoffs than post-processing
37 (e.g., [Agarwal et al., 2018]). We agree regarding demographic parity, and we will note this in our revision.

38 **R3:**

39 > Suggestions on expanding differential privacy discussion

40 We appreciate the insightful suggestions, and will expand our discussion accordingly. In particular, we do not foresee
41 any difficulties in combining our approach with any privacy-preserving technique for standard (non-sensitive) features.

42 **References**

- 43 [1] N. Charoenphakdee, J. Lee, and M. Sugiyama. On symmetric losses for learning from corrupted labels. In *ICML*, 2019.
44 [2] S. Jain, M. White, and P. Radivojac. Recovering true classifier performance in positive-unlabeled learning. In *AAAI*, 2017.
45 [3] J. Katz-Samuels, G. Blanchard, and C. Scott. Decontamination of mutual contamination models. *JMLR*, 20(41):1–57, 2019.
46 [4] R. Kiryo, G. Niu, M. du Plessis, and M. Sugiyama. Positive-unlabeled learning with non-negative risk estimator. In *NIPS*. 2017.
47 [5] B. van Rooyen and R. C. Williamson. A theory of learning with corrupted labels. *JMLR*, 18(228):1–50, 2018.