

1 We thank the reviewers for the thorough and thoughtful reviews. Below are responses to the main comments.

2 **(R2) Boston experiment:** We agree with R2 that equalizing false positive rates could be helpful in some practical  
 3 applications regarding criminality. We ran an experiment where we added the constraint of equalizing false positive  
 4 rates in addition to equalizing true positive rates simultaneously. This constraint is also known as *equalized odds* (Hardt  
 5 et al. 2016) [22] and has some precedent for use in the policing context.<sup>1</sup> Enforcing equalized odds additionally tests the  
 6 proposed robust approaches on stricter fairness criteria compared to the “easier” criteria of equalizing true positive rates  
 7 alone (another concern pointed out by R2). The model can no longer arbitrarily add positive predictions to increase the  
 8 true positive rates, as this would also increase the false positive rates. Results are reported in Table 1 below.

9 Importantly, the label in the Boston dataset is whether an individual was searched/frisked in the past, and does not  
 10 indicate whether the search was justified. If a predictive model trained with this label were used to determine whether to  
 11 search/frisk someone in the future, then that application itself could harmfully carry forward biases in past search/frisk  
 12 decisions in the training set. We want to emphasize that we do not endorse the application of the predictive model in  
 13 this way. To ensure that we are not implicitly endorsing any problematic downstream applications of this model, we  
 14 have included a paragraph specifically highlighting this label bias issue in the description of the Boston experiment.

15 **(R2) Additional dataset:** We ran an additional experiment with the NYPD Stop, Question, and Frisk dataset<sup>2</sup> (SQF),  
 16 and train a classifier in the same way as prior work<sup>134</sup> to predict whether an individual when stopped does *not* possess  
 17 an illegal weapon. Unlike the Boston dataset, this label is not subject to the same historical decision bias (though it still  
 18 suffers from sampling bias as we have labels only for those who are searched). Table 1 contains results of enforcing  
 19 an *equalized odds* constraint. As with the Boston experiment, the proxy race groups are estimated from the *precinct*  
 20 feature in combination with public US census data, with a similar overall noise level 0.53 (see Appendix F.1.2). We’ll  
 21 be happy to include these results, either in addition to or as a replacement for the Boston experiments.

Table 1: Error rate and true positive rates (TPR) / false positive rates (FPR) constraint violations on test for Boston (top three rows) and SQF (bottom three rows) (mean and std. err. over 10 splits).

Algorithm	Unconstrained	$G$ known	Naïve	DRO	Soft assign.
Error rate (Boston)	$0.278 \pm 0.001$	$0.290 \pm 0.001$	$0.290 \pm 0.001$	$0.315 \pm 0.001$	$0.320 \pm 0.001$
TPR Max $G$ viol.	$0.059 \pm 0.012$	$-0.010 \pm 0.005$	$0.008 \pm 0.005$	$-0.029 \pm 0.003$	$-0.018 \pm 0.002$
FPR Max $G$ viol.	$0.007 \pm 0.002$	$-0.011 \pm 0.001$	$-0.007 \pm 0.002$	$-0.008 \pm 0.008$	$-0.020 \pm 0.001$
Error rate (SQF)	$0.148 \pm 0.004$	$0.181 \pm 0.008$	$0.165 \pm 0.004$	$0.343 \pm 0.015$	$0.382 \pm 0.033$
TPR Max $G$ viol.	$0.004 \pm 0.007$	$-0.014 \pm 0.007$	$0.043 \pm 0.009$	$0.023 \pm 0.013$	$-0.006 \pm 0.012$
FPR Max $G$ viol.	$0.112 \pm 0.039$	$-0.006 \pm 0.028$	$0.087 \pm 0.012$	$-0.008 \pm 0.025$	$-0.020 \pm 0.022$

22 **(R3) Upper bound on fairness violation:** R3 raises the question on scenarios where the upper bound on the fairness  
 23 violation w.r.t the true group using TV distance is tight (Theorem 1). In particular, Theorem 1 is tight for the family of  
 24 functions that satisfy  $|h(\theta, x_1, y_1) - h(\theta, x_2, y_2)| \leq 1$ . This condition holds for any fairness metrics based on rates  
 25 such as demographic parity, where  $h$  is simply some scaled combination of indicator functions. However, for a different  
 26 particular given set of  $h$ , Theorem 1 may not be tight. We have clarified this discussion in the main text.

27 **(R3) Cases where the noisy proxy variables  $\hat{G}$  are high-dimensional compared to the true groups  $G$ :** We agree  
 28 with R3 that an evaluation of the SA approach when  $\hat{G}$  and  $G$  have different dimensionalities would be valuable in a  
 29 future empirical study. Theoretically, this can be handled by both the SA and DRO approaches. In particular, for DRO,  
 30 Lemma 1 can be generalized to  $TV(p_j, \hat{p}_i) \leq P(\hat{G} \neq i | G = j)$ ,  $j \in \mathcal{G}, i \in \hat{\mathcal{G}}$ , and the true group distribution  $p_j$  can be  
 31 bounded in a TV distance ball centered at  $\hat{p}_i$ . It is interesting future work to compare the robust approaches when the  
 32 noisy proxy variables have different dimensionality from the true groups. We’ve added this discussion to the main text.

33 **(R1 & R4) Estimating the bound on TV distance:** Lemma 1 provides a practical way to estimate an upper bound on  
 34 the TV distance between the data distributions under the true groups and noisy groups from an auxiliary dataset. As R1  
 35 points out, the looseness of the bound will lead to over-conservativeness of DRO approach. Furthermore, as R4 points  
 36 out, if estimates from the auxiliary dataset are off, the robust approaches will also not be calibrated correctly. We agree  
 37 that these are both important to note to practitioners. Developing methods to better calibrate the DRO neighborhood,  
 38 and further study of the impact of distribution mismatch between the main dataset and the auxiliary dataset would be  
 39 valuable future work. We have added more thorough discussions on these limitations and future work in the main text.

<sup>1</sup>Nathan Kallus, Angela Zhou. Residual Unfairness in Fair Machine Learning from Prejudiced Data. *ICML*, 2018.

<sup>2</sup><http://www1.nyc.gov/site/nypd/stats/reports-analysis/stopfrisk.page>

<sup>3</sup>Zafar et al. From Parity to Preference-based Notions of Fairness in Classification. *NeurIPS*, 2017.

<sup>4</sup>Goel et al. Precinct or prejudice? understanding racial disparities in new york city’s stop-and-frisk policy *Annals of Applied Statistics*, 2016.