

1 **To all reviewers** We would like to thank the reviewers for their thoughtful comments and useful suggestions, which
2 helped us improve the manuscript. Below we provide point-by-point responses.

3 **To Reviewer 1**

4 [R1-1] (*Extension to multi-class classifiers*) As per your great suggestion, we will present a multi-class classifier
5 counterpart that we developed yet not included in the current draft. The idea is to make a binary hard decision
6 *individually for each element* in the softmax output, and then to estimate the pmf of each hard decision via KDE. For a
7 3-way classification synthetic dataset (3 Gaussian mixture), this extended method offers respectful improvements over
8 [44] - trends of the gains are similar to those in Figs. 2 and 3. In a revision, we will provide details on our extension
9 together with the experimental results while moving experimental details to the appendix as suggested.

10 [R1-2] (*Hyperparameter tuning*) (a) Yes, we exhaustively searched hyperparameters for the baselines. For instance, the
11 learning rate was best-picked among several log-scaled candidates (b) Similarly the narrow MLP was chosen as a result
12 of search - we found the depth and width do not affect too much in performance at least for the considered benchmark
13 datasets. We will clarify these in a revision.

14 [R1-3] (*Motivation of the kernel method and writing-flow*) (a) In fact, the stability is also a key issue that we wanted
15 to highlight, although it is not well motivated in the current introduction. We will rewrite the introduction to better
16 balance the issues while citing relevant papers suggested. (b) Yes, the hard/soft decision issue is a key motivation,
17 and the relevant insight allowed us to use the kernel method powerfully and beneficially. For a smooth logical flow,
18 we will re-balance Secs. 3 & 4 so that the issue is strongly emphasized together with an ablation study (w/ and w/o
19 differentiability) in a synthetic setting (as suggested).

20 [R1-4] (*Other comments*) Thanks for your suggestion of the "majority"/"minority" naming as well as pointing out typos.
21 We will fix them.

22 **To Reviewer 2**

23 [R2-1] (*Lack of theory*) We fully agree that the theory is missing for the main claim re. improved performance. Yes,
24 the analysis was not that simple due to the non-convexity of the problem. We will acknowledge this with a proper
25 discussion in a revision.

26 [R2-2] (*Comprehensive experiments for supporting performance gains?*) As per your great suggestion, we will conduct
27 an ablation study in which one may be able to separate improvement due to our regularizer from that due to the use of a
28 more flexible model. Specifically we will compare ours to a kernel SVM by Zafar et al. and include this result in a
29 revision.

30 [R2-3] (*Complexity comparison w.r.t. Agarwal et al*) As you may guess, the theoretical complexity analysis of our
31 algorithm was not done although we empirically demonstrated that ours exhibits lower complexity relative to Agarwal
32 et al. (requiring multiple rounds of training); see Table 1 in supplementary. Instead we will provide in-depth empirical
33 analysis by plotting the running time as a function of the number of data points, as suggested.

34 [R2-4] (*Problem statement & organization*) As per your suggestion, we will make the problem statement clearer and
35 more formal, as well as move the regularizer part into the "Proposed Approach" section.

36 [R2-5] (*Relation to prior work & additional feedback*) Thanks for pointing out the approaches (Wasserstein Fair
37 Classifier etc) that directly estimate fairness measures via information theory. We will cite them with enough discussion.
38 Also we will include the kernel SVM by Zafar et al. and Agarwal et al. in the synthetic setting.

39 **To Reviewer 3**

40 [R3-1] (*Extension to multi-class settings and performance comparison*) As mentioned in response to [R1-1], we actually
41 developed a generalized kernel method that is applicable to multi-class settings. As per your suggestion, we also made
42 performance comparison on a 3-way classification synthetic dataset, observing similar gains as those exhibited in Figs.
43 2 and 3. We will discuss all of these in a revision.

44 [R3-2] (*Multiple sensitive attributes*) Yes, our approach works well for the complex setting. We now conducted
45 experiments on one such setting (AdultCensus with two sensitive attributes: race, gender), observing similar performance
46 improvements, as those in Figs. 2 and 3. We will include the results in a revision.

47 [R3-3] (*Comparison to [44] in many aspects*) While our approach offers key benefits in training stability and tradeoff
48 performance, we do agree that [44] is more flexible in terms of application domains. For a fair comparison, we will
49 summarize pros-&-cons of our approach relative to [44] in many aspects.

50 **To Reviewer 4**

51 [R4-1] (*Dataset set in Fig. 1 and accuracy of the pmf estimate*) We employed a Gaussian mixture: $0.3 \cdot \mathcal{N}(0.37, 0.0055) +$
52 $0.7 \cdot \mathcal{N}(0.74, 0.0055)$. Here the true probability is around 0.7, and this is very close to the pmf estimates in Fig. 1(Right).
53 We will provide this in a revision.

54 [R4-2] (*Robustness quantification and its relation with accuracy*) One way of quantification is to compute the variance
55 of the pmf estimates over different h 's. We will mention this in a revision while plotting accuracy as a function of h .

56 [R4-3] (*Removal of sensitive attributes?*) Yes, that is one natural trial. However, such removal does not ensure fairness
57 especially when X is correlated with Z . Please see [42] for details.