

1 We thank all the reviewers for their helpful comments and suggestions.

2 **Reviewer 1** *RE: needing fair unlabeled data.* Just to clarify, when you say we need fair unlabeled data, do you mean  
3 unlabeled data ( $D$ ) without selection bias? In that case, point taken, but we would argue that in some cases label bias is a  
4 worse problem, and that to some extent, a good data scientist can help mitigate some of the known biases in the data  $D'$   
5 to make it look more like what  $D$ . Moreover, if you could identify a subset of the data that is less biased, then maybe  
6 you could use this subset to improve the more biased subset through semi-supervision.

7 *RE: synthetic data.* Synthetic data is a powerful tool that is frequently employed in ML. For some scientific questions,  
8 like ours, it is the only empirical way to study them. If NeurIPS requires experiments on real world data then we would  
9 have to find another venue. Unfortunately, real-world data does not have the observable quantities we need to measure  
10 the true accuracy which we require to investigate the relationship between fairness and accuracy. We could follow R3's  
11 suggestion, but in the end, this still requires synthetic modifications to the data. We had initially considered a similar  
12 idea, but opted to discuss the perspective learning theory provides instead.

13 **Reviewer 2** First, thank you very much for the thorough and thoughtful review, especially about the related work. It  
14 is very helpful and constructive.

15 *RE: related work* We actually agree with your characterization of the literature and see it as complementary to ours.  
16 Perhaps our language was sometimes too strong in our attempt to highlight what we perceive as a problem. E.g., the use  
17 of “most” instead of “many” may be severe. At the very least, “many” seems appropriate, and as you say, many fairness  
18 papers — including some highly cited papers by highly respected authors — fail to clearly state the assumptions. Take  
19 for example, Zemel et al *Learning Fair Representations*, one of the most highly cited papers in the area. The authors  
20 defer to two previous papers for accuracy reporting (Kamishima et al 2011 and Kamiran&Calders 2009), both of  
21 which discuss accuracy and fairness, but neither of which acknowledge the problem of label bias in that discussion.  
22 Conclusions are then drawn, discussed, disseminated (and repeated) without the assumptions needed to interpret them.<sup>1</sup>

23 It's true that some papers do mention label bias, and we directly quote the relevant passage from such a paper in the  
24 related work. Though, our perception is that more papers (than not) fail to mention label bias (and this might be where  
25 our characterizations of the literature differ). We believe that label bias is omitted frequently enough that it affects  
26 people's thinking on matters of fairness. For example, sometimes when presenting the semi-supervision idea to friends  
27 in the field for the first time, it is initially dismissed because of the alleged tradeoff. Once assumptions are accurately  
28 communicated and agreed upon, everything becomes copacetic again and the idea is accepted as realistic. To speculate:  
29 part of the problem is that in machine learning, we're so used to our gold standard labels being the indisputable truth  
30 (modulo Cohen's kappa), that it's easy to overlook label bias in fair ML, for which that's no longer true.

31 Thanks for pointing out our error about the Fish et al paper, which indeed does use simulated data to get access  
32 to the unbiased labels. They are exploring different questions than us, and so they don't end up investigating the  
33 accuracy fairness tradeoff, which we see as one of the main contributions of our paper and a significant distinguishing  
34 characteristic of our work.

35 *RE: contribution* We agree that adding yet another term to an objective function is not usually particularly novel, but  
36 in this case it first requires getting over the intellectual hump formed by the fairness accuracy trade-off. Maybe this  
37 is obvious, but the idea was initially met with criticism because of this purported obstacle. The synthetic bias is a  
38 mechanism that we use to investigate the fairness-accuracy tradeoff and we see the investigation itself as a key part of  
39 the contribution.

40 *RE: other comments.* (a) Agree that fair ERM might be a better term (b) Data set size was chosen to accommodate  
41 dimensionality while being small for computational expedience (we repeat each experiment 10x so there is actually  
42 10 times more data in total), but in retrospect we could've made it larger (c) it's non-trivial (but possible via more  
43 sophisticated rejection sampling methods) to adjust for perfect class balance while turning the various other experimental  
44 knobs (like protected/unprotected ratio) since there are some dependencies between the experimental parameters;  
45 fortunately, the imbalance was not severe. (d) We will look at the language and try to simplify it for a broader audience.

46 **Reviewer 3** This review very accurately identifies our specific contributions, and their relative importance to each  
47 other. We initially planned on doing something similar to your suggestion (3) in which we synthetically modified the  
48 COMPAS data, but decided to use the space for the learning theory example instead. Maybe a more compact way of  
49 conveying the same data would free up some room as you suggested in (2). With regard to (1), we cited that work as an  
50 example of in-processing and used the same objective function as our semi-supervised method, but on the testing set  
51 instead. Perhaps we should've cited a classification paper for in-processing too.

---

<sup>1</sup>To be clear, we have great respect for these works and these researchers, and are simply using them to illustrate the point.