

1 We thank the reviewers for their helpful feedback. Below we address specific comments.

2 **Reviewer #1:** The reviewer’s point about the clarity of the algorithm in Section 5 is well taken. Please note that some  
3 omitted details were relegated to Appendix E. We’d be happy to flesh out the algorithm’s description in the body of the  
4 paper, as the reviewer suggests.

5 **Reviewer #2:** The reviewer’s main concern has to do with the experimental methodology. We understand the comment  
6 that our methodology (generate instances in a way that there is an envy-free classifier with zero loss) is “slightly  
7 backward.” We wish to clarify, though, that the reasoning behind this experimental design is that, if we generated  
8 instances as the reviewer suggests, we wouldn’t be able to identify the envy-free classifier that minimizes loss, and,  
9 consequently, we wouldn’t be able to measure how far our algorithm is from the optimum. Since there are no existing  
10 algorithmic benchmarks, we would argue that it’s sensible to evaluate our algorithm by generating instances in a way  
11 that the optimal solution — i.e., the minimum loss achievable by an envy-free classifier — is known upfront, as we say  
12 in lines 281–283.

13 Let us now answer the reviewer’s specific questions, using the given numbering; in our revision we will clarify all  
14 issues the reviewer listed.

- 15 (i) This is referring to expected utility. Please see lines 131–136. The word “overall” is confusing and will be  
16 deleted.
- 17 (ii) Exactly, envy is computed for each pair and then averaged over pairs. In Figure 3, “negative envy is replaced  
18 with 0, to avoid obfuscating positive envy.” Please see lines 310–311.
- 19 (iii) Absolutely, these fractions correspond to  $\alpha$  and  $\beta$  from Definition 1. The purpose here is not to make a technical  
20 connection to these parameters, though, but rather to make Figure 4 more concrete by giving examples of two  
21 points on the solid orange and magenta lines.

22 **Reviewer #3:** The reviewer notes that the “paper offers an interesting, original notion of envy-free fairness,” that “the  
23 paper is well-written,” and that “it offers nice technical results.” On the negative side, the reviewer raises two issues,  
24 regarding the notion of fairness and the existence of utility functions. Since both points are rather terse, we weren’t  
25 quite sure what specifically the reviewer is concerned about — we apologize if we misunderstood.

26 *On the notion of envy-freeness:* The reviewer writes that envy-freeness is “well-studied in other disciplines such as  
27 sociology, psychology and economics,” so it seems that the importance of envy-freeness as a notion of fairness isn’t in  
28 question. Rather, if we understand correctly, the reviewer is questioning whether the insights from these other disciplines  
29 carry over to the machine learning domain in practice. While human-subject experiments are needed to definitively  
30 answer this question,<sup>1</sup> we note that there is a significant body of empirical work on envy-freeness in computational fair  
31 division [1], HCI [3], and behavioral economics [2]. The main insight is that people perceive situations where they  
32 are envious — i.e., those where they have a higher utility for someone else’s outcome than for their own — as unfair;  
33 there is no reason why the same conclusion wouldn’t hold in the classification setting, as it shares many of the same  
34 characteristics. We’d be happy to elaborate on this point in our revision.

35 *On the existence of utilities:* The reviewer writes that the paper “makes a very strong assumption of the existence of a  
36 particular form of utility function  $u(x, h(x))$ .” We actually view this as a very mild assumption. All we’re assuming is  
37 that each individual has a utility for each outcome, and the utility for a distribution over outcomes is the expected utility.  
38 Such utility functions, known as von Neumann-Morgenstern utilities, are the basis for much of the work in economics,  
39 decision theory, algorithmic game theory, and related disciplines.

## 40 References

- 41 [1] Y. Gal, M. Mash, A. D. Procaccia, and Y. Zick. Which is the fairest (rent division) of them all? *Journal of the*  
42 *ACM*, 64(6): article 39, 2017.
- 43 [2] D. K. Herreiner and C. D. Puppe. Envy freeness in experimental fair division problems. *Theory and decision*,  
44 67(1):65–100, 2009.
- 45 [3] M. K. Lee and S. Baykal. Algorithmic mediation in group decisions: Fairness perceptions of algorithmically  
46 mediated vs. discussion-based social division. In *Proceedings of the ACM Conference on Computer Supported*  
47 *Cooperative Work and Social Computing (CSCW)*, pages 1035–1048, 2017.

---

<sup>1</sup>There are preciously few such studies even with respect to the well established notions of fairness in machine learning.