

1 **Author Response: Classification Under Misspecification: Halfspaces, Generalized Linear Models, and Evolv-**
2 **ability**

3 We would like to thank all the reviewers for taking the time to understand the contributions of our paper, and for their
4 helpful comments and/or suggestions. We do not have much to add, but would just like to emphasize a few points in
5 response to the reviews.

6 We think our contribution relative to the breakthrough work of Diakonikolas et al. is not just that our algorithm is proper
7 or that the insights behind it lead to algorithms for more general concept classes, but that by avoiding partitioning the
8 domain into a polynomial number of regions, it actually becomes practical and something that we can run on real data.

9 We think that the experimental results are striking, but still only a proof-of-concept in the sense that we added the noise
10 to the data ourselves. A truly compelling demonstration would be, like the first reviewer said, if we could find some real
11 data where our algorithm works well and is demonstrably more fair. This is a direction we are actively pursuing, but we
12 feel that it is a substantial project and would likely be a separate paper if it is successful. Nevertheless some works
13 in fairness work with a graphical model whose causal structure produces confounding effects that lead off-the-shelf
14 algorithms to produce unfair decision rules, see e.g. Kusner et al. [2017]. These types of models naturally lead to
15 situations where noisy observations of some latent quality score might be more variable for some demographics than
16 for others.

17 Also, we agree that it is hard to do justice to all the technical ingredients in just 8 pages. We attempted to give a more
18 detailed outline for our proper learner, and just some of the key ideas for GLMs. It is an interesting suggestion that we
19 could have split it into two papers. However the results actually build on each other, e.g. our algorithm for GLMs in
20 the $\zeta = 0$ case depends on some knowledge distillation primitives which in turn use our proper learning algorithm for
21 halfspaces.

22 **References**

23 Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness. In *Advances in neural*
24 *information processing systems*, pages 4066–4076, 2017.