

1 We thank the reviewers for the constructive comments. We will revise the paper accordingly. Below are the responses  
2 to the main concerns.

3 **Reviewer #1.**

4 - The assumption that the causal graph is given is common in the fairness research based on Pearl's structural causal  
5 models. In practice, there are quite a number of algorithms to build causal graphs from the data (and possibly some  
6 background knowledge), such as the PC algorithm, the GES algorithm, the FCI algorithm, and their variants.

7 - We admit that it is a limitation of the proposed method that requires all variables are discrete with finite domains.  
8 The barrier to continuous variables lies in how to parameterize a causal model for continuous variables with arbitrary  
9 distributions and how to solve the infinite-dimensional optimization problem when we estimate the bounds. For the first  
10 problem, there are some related work on causal graph learning and inference with continuous variables, under some  
11 model/distribution assumptions, e.g. the additive noise model. But relaxing those assumptions is challenging. For the  
12 second problem, there will be infinite response variables to parameterize the continuous causal model, thus the optimal  
13 solution to  $P(\mathbf{r})$  is infinite dimensional. Hence, Eq. 4 (estimate the tight bound) is an infinite-dimensional optimization  
14 problem, which is also challenging. How to address these two challenges will be a future direction for our research.

15 - Constructing fair predictive models is another future research direction. One possible method would be to incorporate  
16 the bounding formulation into a post-processing method. The new formulation will be a min-max optimization problem,  
17 where the optimization variables will include response variables  $P(\mathbf{r})$  as well as a post-processing mapping  $P(\hat{y}|\hat{y}, \text{pa}_Y)$ .  
18 The inner optimization is to maximize the path-specific counterfactual effect to find the upper bound, and the outer  
19 optimization is to minimize both the loss function and the upper bound. We plan to explore this method in future work.

20 - The proposed method can provide the tightest bounds because the response variables cover all possible domains of  $\mathbf{U}$   
21 so that we can explicitly traverse all possible causal models. We will add more explanations about how the proposed  
22 method works in the revised version.

23 - In Table 3, the results of the proposed method are either equivalent to or tighter than previous methods. The bold  
24 lines are to highlight the situation where the tighter bounds make differences in detecting discrimination, showing the  
25 practical meaning of the proposed method.

26 **Reviewer #2.**

27 - To the best of our knowledge, the notion of path-specific counterfactual effect has not been proposed in previous  
28 works. It is worthy to point out that a similar term has been used in paper "Path-Specific Counterfactual Fairness"  
29 (AAAI'19), but with a different meaning. The paper studied the causal effect along some specific pathways without  
30 conditioning on any observed values, which is equivalent to path-specific fairness, a special case of our proposed  
31 fairness notion where  $\mathbf{O} = \emptyset$ . In paper "A Potential Outcomes Calculus for Identifying Conditional Path-Specific  
32 Effects" (AISTATS'19), the conditional path-specific effect is different from our notion in that, for the former the  
33 condition is on the post-intervention distribution, and for the latter, the condition is on the pre-intervention distribution.  
34 We will add more references and discussions in the revised version.

35 - Our proposed notion is definitely practical. It can not only unify the previous notions but also resolve new types of  
36 fairness that the previous notions cannot do. A typical example is individual indirect discrimination, which means  
37 discrimination along the indirect paths for a particular individual. Individual indirect discrimination has not been  
38 studied yet in the literature, probably due to the difficulty in definition and identification. However, it can be directly  
39 defined and analyzed using our proposed notion by letting  $\mathbf{O} = \{S, \mathbf{X}\}$  and  $\pi = \pi_i$ . Note that the condition here is  
40 on the pre-intervention distribution, i.e., we focus on a particular individual with certain observed values, and want to  
41 estimate the change of these values after the intervention is performed. Thus, individual indirect discrimination cannot  
42 be defined using the above conditional path-specific effect. We will add the above discussions and make our motivation  
43 clearer in the revised version.

44 **Reviewer #3.**

45 Thanks for the comments. We will incorporate all the comments into the revised version. We will add more deriving  
46 details for Section 4.2, reorganize this manuscript accordingly, and move some discussions into the supplementary file  
47 if necessary.