

1 We would like to thank the reviewers for their comments and suggestions. We will incorporate their feedback in the
2 revised version of the paper.

3 **Reviewers #1 and #2.** As discussed in lines 125-135, our algorithms will look for the counterfactual explanations \mathcal{A}
4 and decision policies $\pi(\mathbf{x})$ that maximize the decision maker’s utility rather than the individuals’ best interest¹. As a
5 consequence, it is true that, compared to the non-strategic setting, the policies that are optimal in the strategic setting
6 may induce some of the subjects to pay an additional (immediate) cost to change features in order to receive a positive
7 decision, as shown in Figure 6c in Appendix E. However, we would like to point out that any subject who would have
8 received a positive decision under the decision policy that is optimal in the non-strategic setting will still receive a
9 positive decision under the decision policy that is optimal in the strategic setting after they best respond. Moreover,
10 subjects who do pay an additional cost to change features will always increase their outcomes $P(y | \mathbf{x})$, as shown in
11 Figure 7 in Appendix F.3, and this is likely to increase their individual utility in the long term. In the revised version of
12 the paper, we will expand our discussion regarding the potential of our algorithms to favor the decision-maker at the
13 expense of the decision subjects, in light of the results shown in Figures 6c and 7.

14 **Reviewer #2 and #3.** If our submission is accepted, we will make use of the ninth content page of the camera-ready
15 version to bring the definition of α , the algorithmic boxes for Algorithms 1 and 2, and the discussion of the cost function
16 estimation to the main text.

17 **Reviewer #1.** Under our problem formulation, a decision d is beneficial to the individuals who are subject to (semi)-
18 automated decision making if $d = 1$ (e.g., an individual receives a loan) and a prediction \hat{y} made by a machine learning
19 model is positive if $\hat{y} = 1$ (e.g., an individual repays a loan). In this context, note that, rather than explaining predictions
20 by machine learning models as in previous work, we pursue the development of methods to find counterfactual
21 explanations for the decisions, as argued in lines 38-42. We will clarify this in the revised version of the paper.

22 We will expand our comparison with the existing literature and further discuss the necessity to distinguish between
23 decisions and predictions, as argued by several authors in a series of recent papers [23, 25-27, 47, 48].

24 **Reviewer #2.** As noted by the reviewer, some of our assumptions are quite strong, however, we still think they do not
25 nullify our contributions, especially given the paucity of work in the area. That being said, we are hopeful to relax some
26 of these assumptions in future work.

27 We will fix the statement of Proposition 4.

28 **Reviewer #3.** The “strategic setting” refers to a scenario in which individuals who are subject to (semi)-automated
29 decision making use knowledge, gained by explainability, to change their own features to maximize their chances of
30 receiving a beneficial decision. In our work, we formally characterize this setting mathematically for counterfactual
31 explanations.

32 A counterfactual is a statement of how the world would have to be different for a desirable outcome to occur [13]. In our
33 problem formulation, the world are the feature values \mathbf{x} , the desirable outcome is the positive decision $d = 1$, and the
34 statement is the counterfactual explanation $\mathcal{E}(\mathbf{x})$. Given an individual with initial feature values \mathbf{x} who would receive a
35 negative decision $d = 0$, the counterfactual explanation provides her with an example of a feature value $\mathcal{E}(\mathbf{x})$ under
36 which she is guaranteed to receive a positive decision $d = 1$. We will clarify this in the revised version of the paper.

37 The problem formulation is new. Previous work on counterfactual explanations [13-15, 36-37, 39] has focused on
38 explaining predictions, rather than decisions, and has not investigated the connection between strategic machine learning
39 and explanations. The most closely related work is by Tabibian et al. [23] in the strategic machine learning literature,
40 however, they have considered a setting where decision makers share their entire policies with the individuals subjects to
41 their decisions rather than counterfactual explanations. In this context, please note that we have included further related
42 work in Appendix A. If our submission is accepted, we will make use of the ninth content page of the camera-ready
43 version to bring that content to the main.

¹We did not explicitly use the wording social welfare or decision subject’s utility, however, it is in the individuals’ best interest to maximize their utility.