
Linear Models are Robust Optimal Under Strategic Behavior

Wei Tang[†], Chien-Ju Ho[†], and Yang Liu^{*}

[†]Washington University in St. Louis, ^{*}UC Santa Cruz
{w.tang, chienju.ho}@wustl.edu, yangliu@ucsc.edu

Abstract

There is an ubiquitous use of algorithms to inform decisions nowadays, from student evaluations, college admissions, to credit scoring. These decisions are made by applying a decision rule to individual’s observed features. Given the impacts of these decisions on individuals, decision makers are increasingly required to be transparent on their decision making to offer the “right to explanation.” Meanwhile, being transparent also invites potential manipulations, also known as gaming, that individuals can utilize the knowledge to strategically alter their features in order to receive a more beneficial decision.

In this work, we study the problem of *robust* decision-making under strategic behavior. Prior works often assume that the decision maker has full knowledge of individuals’ cost structure for manipulations. We study the robust variant that relaxes this assumption: The decision maker does not have full knowledge but knows only a subset of the individuals’ available actions and associated costs. To approach this non-quantifiable uncertainty, we define robustness based on the worst-case guarantee of a decision, over all possible actions (including actions unknown to the decision maker) individuals might take. A decision rule is called *robust optimal* if its worst case performance is (weakly) better than that of all other decision rules. Our main contributions are two-fold. First, we provide a crisp characterization of the above robust optimality: For any decision rules under mild conditions that are robust optimal, there exists a linear decision rule that is equally robust op-

timal. Second, we explore the computational problem of searching for the robust optimal decision rule and demonstrate its connection to distributionally robust optimization. We believe our results promote the use of simple linear decisions with uncertain individual manipulations.

1 Introduction

Algorithms have been increasingly engaged in making consequential decisions across a variety of sectors in our society. Examples include judges using defendant risk scores to set bail decisions and banks evaluating individuals’ profiles to make loan decisions. In these scenarios, the decision maker aims to determine a decision rule (or a model), which takes a set of individual’s observed behavior or features as input, and output decisions that maximize some given utility function¹.

Given the consequential impacts to individuals, there is an increasing demand to make the decision rule transparent to offer “right to explanation” (Goodman and Flaxman, 2017). Transparency not only allows the public to audit models to mitigate potential fairness concerns but also enables the participants to understand what decisions they might receive if they have different features (See, for example, “right to recourse” (Ustun et al., 2019)). However, on the flip side, transparency simultaneously creates opportunities for individuals to strategically respond to the deployed model. Specifically, if individuals understand how their observed features affect decisions, they may strategically alter their features to obtain a more favorable decision.

In response to this strategic behavior, there has been a recent flurry of work in studying decision making under strategic behavior (Brückner et al., 2012; Brückner and Scheffer, 2011; Hardt et al., 2016; Kleinberg and Raghavan, 2019; Alon et al., 2020). To make the analysis

Proceedings of the 24th International Conference on Artificial Intelligence and Statistics (AISTATS) 2021, San Diego, California, USA. PMLR: Volume 130. Copyright 2021 by the author(s).

¹Throughout the work, we address the decision maker as “she” and the individual as “he”. We also use the terms individual and agent interchangeably.

tractable, almost all the works explicitly assume the decision maker has the *full knowledge* of agents’ action space and the corresponding costs for agents to manipulate their features. The above knowledge enables a game theoretic analysis that characterizes agents’ best responses when offered a particular decision rule.

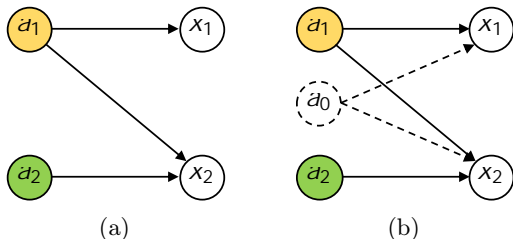


Figure 1: An instance of student evaluation problem.

However, the “full information” assumption is often not true in practice. Consider an example of student evaluation in Fig. 1 (Kleinberg and Raghavan, 2019; Alon et al., 2020). The student’s observed features are their exam score (x_1) and homework score (x_2). The student can choose to either study (a_1) or copy homework answers (a_2) to alter their features. Studying improves both exam score (x_1) and homework score (x_2), while copying homework only improves homework score. The teacher evaluates the student through a final score, which is a function of x_1 and x_2 , and students are assumed to aim to maximize their final score minus the cost of the actions. If the teacher knows the actions a_1 and a_2 , and they are indeed the only actions the student can take, the teacher can design a decision rule (a final score as a function of x_1 and x_2) that maximizes some given objective by considering students’ best responses. However, in practice, the teacher might not be aware of the full set of actions the student can take. For instance, the student might consider taking action a_0 unknown to the teacher (in Fig. 1b), such as hiring a tutor or working with other students. With this incomplete knowledge of the student’s actions, how should the teacher design her evaluation rule?

In this work, we answer the above question by studying the design of *robust* optimal decision rules with strategic agent, where we relax the assumption of complete knowledge over agent actions. We define the robustness notion as used in robust contract design (Carroll, 2015): Evaluate the worst-case guarantee of a decision, over all possible actions (including actions unknown to the decision maker) agents might take. More formally, the decision maker only knows a subset of actions (denoted by A_d) among all the actions available to the agent (denoted by A_a). Let $V_d(f/A_a)$ be the utility the decision maker obtains with decision rule f when the agent’s action space is A_a . The decision maker’s goal is to maximize her *worst-case* performance $V_d(f)$ over all possible actions the agent may have access to

($A_a \supseteq A_d$):

$$\max_f V_d(f) = \max_f \inf_{A_a \supseteq A_d} V_d(f/A_a): \quad (1)$$

A decision rule f^* is robust optimal if it achieves the maximum of the above worst-case utility.

Our contribution Our contributions are two-fold. First, we formalize the problem of robust strategic decision-making and characterize the robust optimal decision rules. We show that under mild conditions, for any robust optimal decision rule, there exists a linear one that is equally robust optimal. Our result implies that, to find robust optimal decision rules, it suffices to search over the space of linear decision rules. Second, we explore the computational problem of searching for the robust optimal f^* . While the problem is NP-hard in general (since non-robust strategic decision-making is only solvable in restricted settings but is generally NP-hard (Kleinberg and Raghavan, 2019)), we investigate the additional complexity introduced by our robustness desiderata, through adapting techniques from distributionally robust optimization (Delage and Ye, 2010). Our results inform efficient algorithms especially in settings when non-robust strategic decision-making problem is efficiently solvable.

1.1 Related Work

Our problem closely connects to the recent literature in machine learning in the presence of strategic manipulation (Hardt et al., 2016; Brückner et al., 2012; Brückner and Scheffer, 2011). Hardt et al. (2016) study the design of optimal classification when the agents can incur costs to manipulate their features. Motivated by fairness concerns, Hu et al. (2019) and Milli et al. (2019) consider settings in which the costs for manipulation differ for different groups and explore the societal impacts. There are also works directly utilizing the decision rule as an incentive device to induce desired behavior (Kleinberg and Raghavan, 2019; Alon et al., 2020; Haghtalab et al., 2020; Ball, 2020; Dong et al., 2018; Tabibian et al., 2019; Miller et al., 2019). Among these works, Kleinberg and Raghavan (2019) is closest to our work: they introduce a graphic model to capture the known agent’s available actions and show that simple linear mechanisms suffice for a single known agent. Alon et al. (2020) then extend the discussion to multiple agents. Our work departs from the above works in the sense that the decision maker only has incomplete knowledge of the agent’s cost structure or his available actions.

Our formulation resembles the principal-agent problem in contract theory (Grossman and Hart, 1992; Shavell, 1979; Holmstrom and Milgrom, 1987), which studies the

strategic interplay between two parties with misaligned interests. Our characterization of robust decision rule follows the works on robust contract design (Carroll, 2015; Dai and Toikka, 2017; Miao and Rivera, 2016; Carroll and Segal, 2019; Carroll, 2017; Diamond, 1998; Hansen and Sargent, 2012; Chassang, 2013) in which robustness is defined as the worst-case optimal mechanisms. Our work differs from this line of research in that the decision maker determines a decision rule (instead of a “contract” in contract theory) that is multi-dimensional and could take arbitrary forms. Moreover, we do not restrict the decision maker’s utility to be in additive form (reward minus the payment). We generalize the utility to be arbitrary function that satisfies some mild conditions. Other computational approaches to contract design in computer science community can be found in the work by Dütting et al. (2019); Babaioff et al. (2006); Ho et al. (2016); Babaioff et al. (2010). Our work also shares similar flavor for max-min analysis in worst-case algorithmic analysis (Azar et al., 2013; Bandi and Bertsimas, 2014). In all of these works, the setting and the formulation are different from the ones we consider in the present work.

Our work complements a recent literature on discussing the effects of linear models in social stratification. For example, Wang et al. (2018) extend the notion of interpretability to credibility and discuss the credibility in a linear setting. Fawzi et al. (2018) analyze the robustness of linear classifiers to adversarial perturbations. Ustun and Rudin (2014) and Ustun et al. (2019) discuss the interpretability and right to recourse in linear classification. Our work promotes the usage of linear models: In addition to interpretability and good generalization, linear models are also robust to unknown strategic manipulation.

2 A Model of Robust Strategic Decision-making

In this section, we formalize our model for robust strategic decision-making. Agent features are represented by a vector $\mathbf{x} = (x_1, \dots, x_n)$, which takes value in a bounded compact set $X \subseteq \mathbb{R}^n$. The agent can take actions to alter the features. An action of the agent can be represented by the outcome (i.e., the distribution of agent features after the action) and the cost of the action. We use a pair $(P; c) \in (X) \times \mathbb{R}_+$ to denote an action, where P is the outcome, i.e., the distribution of the agent features after action, and c is the associated cost. The decision maker cannot observe the agent’s action but can only observe the features, the realized outcome of the action.

Action set We define two important action sets A_a and A_d . In particular, $A_a \subseteq (X) \times \mathbb{R}_+$ is the set of all possible actions that the agent can take, and A_d is the set of action that the decision maker is aware of. While the decision maker only knows A_d and not A_a , she knows that $A_d \subseteq A_a$. The decision maker’s unquantifiable uncertainty of A_a is the key conceptual element of this work. Informally, using the student evaluation example, the available actions to the student A_a could be (studying, cheating, hiring tutors). The teacher only knows A_d , (studying, cheating), a subset of A_a but aims to design a decision rule that is robust to this uncertainty.

Decision rule A decision rule $f: X \rightarrow \mathbb{R}_{\geq 0}$ is a mapping from the agent’s features to a decision, where the decision domain of f is normalized to be non-negative and directly represents the value of the decision to the agent. The decision rule f is contingent only on the observable features, but not on the actions that are not observable to the decision maker.

The decision maker aims to maximize her utility function $h: X \rightarrow \mathbb{R}_{\geq 0}$. This function characterizes the utility that the agent brings to the designer. For example, it could be a qualification function, assuming the decision maker aims to increase the chance that the agent passes the qualification, and the agent’s effort in changing their features may lead to self-improvement, thus in their true qualifications. Assume that there’s an upper bound $C > 0$ of $f(\mathbf{x})$ for any $\mathbf{x} \in X$. In addition, we define the following simple class of decision rules:

Definition 1 (Linear decision rule). *A decision rule f is linear if f is a linear function of the feature², i.e., $f(\mathbf{x}) = \omega^\top \mathbf{x} + \beta$ for $\omega \in \mathbb{R}^n$ and $\beta \in \mathbb{R}$. Let $G^{lin} = \{(\omega; \beta) \in \mathbb{R}^n \times \mathbb{R} : f(\mathbf{x}) = \omega^\top \mathbf{x} + \beta \in [0; C]; \forall \mathbf{x} \in X\}$ be the space of parameter pair $(\omega; \beta)$.*

The interaction between the decision maker and the agent goes as follows: (1) the decision maker publishes a decision rule f based on the knowledge of A_d ; (2) the agent, knowing A_a , chooses action $(P; c) \in A_a$ to respond to f ; (3) the agent features are then moved to $\mathbf{x} \in P$; (4) the decision maker derives utility of $h(\mathbf{x})$ and the agent derives utility of $f(\mathbf{x}) - c$.

Robustness of decision making under strategic behavior We first characterize the agent’s behavior. Given the decision rule f and his available action set A_a , the agent obtains expected utility $E_P[f(\mathbf{x})] - c$ for taking action $(P; c)$. Let $A_a^*(f/A_a)$ be the set of actions that maximize the agent’s utility, and $V_a(f/A_a)$ be the

²More precisely, it is an affine decision rule with the form of $\omega^\top \mathbf{x} + \beta$.

corresponding utility:

$$A_a^*(f|A_a) = \arg \max_{(P;c) \in \mathcal{A}_a} (E_P[f(\mathbf{x})] - c);$$

$$V_a(f|A_a) = \max_{(P;c) \in \mathcal{A}_a} (E_P[f(\mathbf{x})] - c);$$

When there are multiple maximizers of the agent's objective, the agent may choose the action that is the most beneficial to the decision maker. The expected utility of the decision maker, given decision rule f and the action set available to the agent A_a is

$$V_d(f|A_a) = \max_{(P;c) \in \mathcal{A}_a(f|A_a)} E_P[h(\mathbf{x})];$$

The decision maker only knows A_d but not A_a . Therefore, she cannot optimize $V_d(f|A_a)$ directly. To address this nonquantifiable uncertainty, we define $V_d(f)$ as the worst case utility the decision maker obtains over all possible actions sets A_a that are supersets of A_d :

$$V_d(f) = \inf_{A_a \supseteq A_d} V_d(f|A_a); \quad (2)$$

We define the robust optimal decision rule f^* as the one that maximizes $V_d(f)$, since it is robust to any action (even unknown to the decision maker) the agent might take:

$$f^* \triangleq \arg \max_f V_d(f) = \arg \max_f \inf_{A_a \supseteq A_d} V_d(f|A_a); \quad (3)$$

3 Linear Model is Robust Optimal Under Strategic Behavior

In this section, we establish our main result that there exists a linear decision rule that is robust optimal.

Theorem 1. *There exists a decision rule f that maximizes $V_d(f)$ and is linear, namely: $f \triangleq \arg \max V_d(f)$; where $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$; for some $\mathbf{w} \in \mathbb{R}^n$; $\mathbf{w} \in \mathbb{R}$.*

The above theorem characterizes the robust optimal decision rule defined in (3). The key implication of the theorem is that, when aiming to find the robust optimal model against strategic responses, it suffices to only consider linear models.

In the following, we provide the proof sketch and use an example to demonstrate our results. Our result and analysis extend the work of robust contract design (Carroll, 2015) to deal with situations in which both the decision rule and the utility of the decision maker can take more general function forms (instead of restricting to one-dimensional contract as decision rule, and additive utility for decision maker). The proof consists of three main steps. We first characterize the properties of the worst case utility $V_d(f)$ for a given decision rule f ; we then show that any nonlinear decision

rule can be (weakly) improved by a linear decision rule in terms of the worst case utility. Finally, we wrap up by showing the existence of an optimal linear decision rule in the linear decision space.

3.1 Characterize the worst-case utility $V_d(f)$

Before we move to the main analysis, consider a trivial case that the decision maker chooses to post no decision rule (i.e., $f(\mathbf{x}) = 0; \mathcal{B}\mathbf{x}$). Since this is also a linear decision rule, if the robust optimal decision rule is to post no decision rule, Theorem 1 is trivially correct. In the following discussion, we focus the discussion on the cases in which the decision maker can benefit from posting some decision rule (otherwise, she can choose to post no decision rules). In particular, we define rational decision rules as follows.

Definition 2 (Rational decision rule). *A decision rule f is rational for the decision maker if $V_d(f) > V_d(0)$, where $V_d(0)$ represents the utility of the decision maker when she publishes no decision rules.*

We first characterize the worst-case utility guarantee for any given rational decision rule.

Lemma 1. *Let f be any rational decision rule. Define a set $\mathcal{P} = \{P \in \mathcal{X} : E_P[f(\mathbf{x})] - V_a(f|A_d) \geq 0\}$. Then one of the following two cases occurs:*

$$(i) \quad V_d(f) = \min_{P \in \mathcal{P}} E_P[h(\mathbf{x})]; \quad (4)$$

or

$$(ii) \quad \max_{P \in \mathcal{X}} E_P[f(\mathbf{x})] - V_a(f|A_d) = 0; \quad (5)$$

Moreover, for P attaining the minimum in (4), the inequality in (5) will reduce to equality at P .

The key message of this lemma is that, we can replace the definition of $V_d(f)$ in (2), that depends on unknown A_a , with an expression that depends only on variables known to the decision maker. In particular, in case (i), this is given by identifying \mathcal{P} which is constrained by $V_a(f|A_d)$ using the designer's knowledge A_d . In case (ii), we know that the best response from the agent is indeed in A_d , so again the designer can focus on the action space she is aware of.

Proof Sketch. The full proof is in Appendix 6, and we provide a sketch here. For any action set $A_a \supseteq A_d$ the agent has, and any optimal action $(P; c)$ he chooses under A_a and the rational decision rule f , the expected utility the agent gets from f must satisfy:

$$E_P[f(\mathbf{x})] - c = V_a(f|A_a) - V_a(f|A_d);$$

Here the second inequality holds because A_a contains A_d , and having more actions available can only make the agent better off. Thus, for any decision rule f ,

the agent will only take the actions that guarantee himself a utility that is at least $V_a(f/A_d)$, these action actually formulates the set S . Furthermore, the decision maker's utility $V_d(f/A_d) = E_P[h(\mathbf{x})]$ is at the least the minimum given by Eqn. (4). Thus, we have $V_d(f) \geq \min_{P \in \mathcal{P}} E_P[h(\mathbf{x})]$. To show this is actually tight, we then prove the other direction. To achieve that, we construct some worst case action set A_d to guarantee that $V_d(f)$ cannot exceed $\min_{P \in \mathcal{P}} E_P[h(\mathbf{x})]$. Case (i) is simply the boundary case in which the agent's best action under any possible actions sets is already included in A_d . \square

3.2 Improve nonlinear rule to a linear one

Having characterized the worst-case utility guarantee of decision maker, we can now show that any nonlinear decision rule can be (weakly) improved by a linear decision rule in terms of its $V_d(f)$.

Lemma 2. *Fix any h and any (nonlinear) rational decision rule f , there exists a linear one f' such that: $V_d(f') \geq V_d(f)$.*

Proof Sketch. The full proof is in Appendix 7. At a very high-level, we show that for every decision rule f , we can construct two convex sets, with one containing information about the agent and one about the decision maker. We then show that the two convex sets are disjoint, and therefore there exists a hyperplane that separates the two convex sets. Then it turns out that separating hyperplane is the linear decision rule that weakly improves on f .

In more detail, given a decision rule f , consider a point $(E_P[\mathbf{x}]; E_P[f(\mathbf{x})])$ generated by any possible action $(P; c)$. This point will be in the convex hull of $(\mathbf{x}; f(\mathbf{x}))$. We define S to be the convex hull of all pairs $(\mathbf{x}; f(\mathbf{x}))$, for $\mathbf{x} \in X$. To construct another convex set, we separately consider the two cases in Lemma 1. For case (i), we define $t(\mathbf{x}) = \max_{f \in V_a(f/A_d)} h(\mathbf{x}) + f(\mathbf{x}) - V_d(f)g$. Intuitively, $t(\mathbf{x})$ is constructed to accommodate the constraint in the set S for Eqn. (4). We define T as the convex hull of all pairs $(\mathbf{x}; z)$ that \mathbf{x} lies in the convex hull of X , and $z > t(\mathbf{x})$. By utilizing the results in Lemma 1, we can show that the two convex sets are disjoint (details in Appendix). By hyperplane separation theorem, we can find a hyperplane f' separating S and T . f' has two advantages: First it gives the agent the same incentive as f . Second, it gives a weakly greater guarantee to the decision maker. For case (ii), we change the set T to be the set of all $(\mathbf{x}; z)$ with \mathbf{x} in the convex hull of X and $z > V_a(f/A_d)$. Similar arguments in case (i) still apply here. \square

3.3 Wrapping up

We have shown that any rational decision rule f can be (weakly) improved to a linear one. We now wrap up our analysis by showing the existence of an optimum within the class of linear decision rules.

Lemma 3. *There exists a robust optimal linear decision rule.*

Recall our definition of G^{lin} in Definition 1. The proof reduces to show that $V_d(f)$ is upper semi-continuous w.r.t. $(!; g) \in G^{\text{lin}}$, this guarantees that $V_d(f)$ has a maximum over the compact set G^{lin} . We defer the proof to Appendix 8.

3.4 Illustrating example: Student evaluation

We now use the example of student evaluation to demonstrate the intuitions of our results and analysis. We first illustrate the application of Lemma 2: For a particular nonlinear decision rule, we show how to find an improved linear decision rule. Then, we compute the worst case utility for both decision rules according to Lemma 1. Finally, we return to the environment with student being able to take actions unknown to the teacher, as depicted in Fig. 1b to discuss how these two decision rules perform.

We first specify the environment details of our example. Suppose each feature is a binary variable in $\{0, 1\}$ (e.g., x_1 : pass or fail the exam, x_2 : whether the homework is qualified or not). Assume the cost of actions are the same, the student needs to decide a distribution over the actions. Using the terminology by Kleinberg and Raghavan (2019), we say the student needs to allocate their effort budget of 1 to two actions, with e_j denoting the effort of (i.e., the probability of choosing) action a_j . The effort-feature conversion obeys the following rule: $\Pr(x_i = 1) = \sum_j w_{j,i} e_j$, where $w_{j,i} \in [0, 1]$ is the weight on how the student's effort $e_j \in [0, 1]$ on action a_j contributes to the value of feature x_i . For example, a student may study for the exam and still fail with some (small) probability. The effort-feature conversion weights are detailed in Fig. 2a.

Suppose for a moment the student's available actions are $\{a_1; a_2\}$. The teacher wants to incentivize the student to invest all their efforts on studying (namely, the action a_1). This could correspond to the teacher setting her utility function as $h(\mathbf{x}) = !^T \mathbf{x} + \eta$, where $! = (1; 0)$ and η is a small positive value³. One (nonlinear) decision rule that maximizes $h(\mathbf{x})$ is $f(\mathbf{x}) = \max\{x_1; x_2\}$. It is easy to verify that this decision rule results in the student to invest all his effort to action a_1 (i.e., $e_1 = 1$), and leads to the teacher's utility of ρ .

³ η can be used to guarantee the decision rules are rational.

Note that f is a (non-robust) optimal decision rule for maximizing h . We now show that by leveraging the constructive proof in Lemma 2, we can find a linear rule that weakly improves the worst-case utility. In particular, upon defining the convex sets S and T for f , we can find one hyperplane $f'(\mathbf{x}) = x_1 + x_2$ that separates these two sets, as illustrated in Fig. 2b. From Lemma 2, f' weakly improves over f in terms of worst-case utility. Below we compute the worst-case utility for f' and f for confirmation. For f' , by Eqn. (4) in Lemma 1, $V_d(f') = \min_{P \in \mathcal{G}} \mathbb{E}_P[h(\mathbf{x})] = \min_{P \in \mathcal{G}} \mathbb{E}_P[x_1] + \mathbb{E}_P[x_2]$, where P satisfies $\mathbb{E}_P[f'(\mathbf{x})] = \mathbb{E}_P[x_1 + x_2] = V_a(f' | A_d)$. Observe that, when the student's available action set A_d is depicted as in Fig. 2a, $V_a(f' | A_d) = 2\rho$. Since when P attains the minimum of $V_d(f')$, the inequality must bind. Thus, we have $\min_{P \in \mathcal{G}} \mathbb{E}_P[x_1] = 2\rho - 1$, which gives us $V_d(f') = 2\rho - 1 + h$. However, follow the same analysis, one can compute that $V_d(f) = h$, which is smaller than $V_d(f')$.

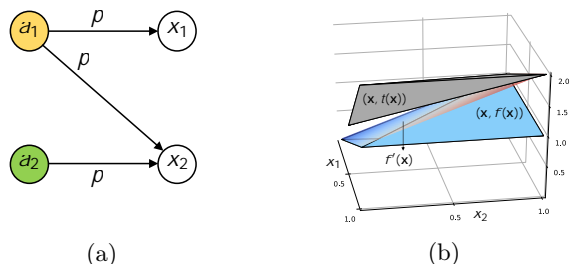


Figure 2: (a): $\rho \in (0.5; 1)$ is the weight parameter. (b): Construct S and T when $f(\mathbf{x}) = \max\{x_1, x_2\}$. The gray shaded region is the set T (we actually plot the convex hull of all points $(\mathbf{x}; t(\mathbf{x}))$), while the light sky blue is the set S . The color hyperplane is exactly $f'(\mathbf{x}) = x_1 + x_2$.

Moreover, f' does outperform f for our example with the student being able to take one action unknown to the teacher, as introduced in Fig. 1b. Suppose the student has one more action a_0 available to accomplish his course responsibilities (where $w_{0,1} = \rho$ and $w_{0,2} = \rho + \epsilon$ for some $\epsilon \in (0; 1 - \rho)$). The teacher is not informed by this change and may only be aware of the original student's available actions (which is $\{a_1, a_2\}$) and has to design her decision rule based on this restricted knowledge (see A_d and A_a in Table 1) Facing this uncertainty, it is easy to see that the linear one f' can guarantee teacher's maximal utility ρ , while f can only ensure a utility of $\rho - \epsilon$ to the teacher (since in this case, the student will deviate to invest all effort to action a_0), which is smaller than ρ .

4 The Complexity for Computing Robust Optimal Decision Rule

Having shown that a robust optimal decision rule f^* is linear, one may wonder whether it is possible to efficiently compute such f^* . Note that our analysis for robust optimality is constructive, and it establishes an algorithmic procedure to compute the optimal f^* . Below we show that computing f^* is generally hard.

Theorem 2. *We state the computation complexity for computing f^* :*

1. *Computing the linear f^* is at least as hard as solving the corresponding strategic decision making problem without robustness concern (under the linear decision space G^{lin}).*
2. *In general, computing f^* is NP-hard since its corresponding strategic decision making problem without robustness concern (under the linear decision space G^{lin}) is generally NP-hard.*
3. *When X is finite, if there is a polynomial-time algorithm for solving the corresponding strategic decision making problem without robustness concern (under the linear decision space G^{lin}), then there is a polynomial-time algorithm for computing f^* .*

The proof and the description of a procedure for computing f^* are included in Appendix 10. The key idea is to first formulate the problem of computing f^* as an optimization problem. We then demonstrate that it can be further decomposed into two optimization problems, with one to be the same as solving (non-robust) optimal decision rule with strategic behavior (under the linear decision space G^{lin}), and the other being a linear program with equality constraint.

More formally, let a linear decision rule be in the form of $f_{(l; c)}(\mathbf{x}) = \mathbf{l}^\top \mathbf{x} + c$, where $(l; c) \in G^{\text{lin}}$ (see Definition 1). We use SO to denote the corresponding (non-robust) strategic decision making problem (under linear decision space G^{lin}) where the agent's available action set is exactly A_d (matching the knowledge of the decision maker):

$$\begin{aligned} & \arg \max_{(l; c) \in G^{\text{lin}}} \mathbb{E}_P[h(\mathbf{x})]; & (\text{S0}) \\ \text{s.t. } & (P; c) \geq \arg \max_{(P; c) \in A_d} \mathbb{E}_P[f_{(l; c)}(\mathbf{x})] - c; & (6) \end{aligned}$$

where $\mathbb{E}_P[\cdot]$ is the expectation taken with respect to the random vector \mathbf{x} given that it follows the probability distribution P . Note that while there exist efficient algorithms to solve this (non-robust) decision making under uncertainty (SO) in restricted cases, the problem is known to be NP-hard in general (Hansen et al., 1992; Kleinberg and Raghavan, 2019). Therefore, instead of

investigating the complexity of solving the robust variant, we aim to understand the *additional* complexity of requiring robustness in (non-robust) decision making under uncertainty.

With slight abuse of notation, for any $(I; \cdot) \in \mathcal{G}^{\text{lin}}$, we use $(P_I; c_I) \in \mathcal{A}_d$ to denote the solution to the constraint (6) in SO and let $C_I = \mathbb{E}_{P_I} [f_{(I; \cdot)}(\mathbf{x})] - c_I$. Then according to Lemma 1, we can compute f^* by solving:

$$\arg \max_{(I; \cdot) \in \mathcal{G}^{\text{lin}}} \min_{P \in \mathcal{P}} \mathbb{E}_P[h(\mathbf{x})]; \quad (\text{R-SO})$$

where \mathcal{P} can be expressed as follows:

$$\begin{aligned} \mathcal{P} = \{ & \mathbf{x} \in \mathcal{X} \mid \mathbb{E}_P[f_{(I; \cdot)}^\top \mathbf{x}] = C_I, \\ & \Pr(\mathbf{x} \in \mathcal{X}) = 1 \} \end{aligned} \quad (7)$$

Different from the problem in SO , after identifying the agent's best response $(P_I; c_I) \in \mathcal{A}_d$ under $f_{(I; \cdot)}$, our problem in R-SO will have an additional layer of optimization over the set \mathcal{P} . It is easy to see that this is a linear program with equality constraint, where the decision variables are a probability simplex over \mathcal{X} . Therefore, the computation of R-SO can be decomposed into the computation of SO and a linear program. This decomposition enables us to complete the proof.

4.1 Solving R-SO when \mathcal{X} is Infinite

So far, we demonstrate that the additional complexity of solving robust decision making under uncertainty (R-SO) compared with the non-robust version (SO) can be characterized by a linear program. When the space of agent features \mathcal{X} is finite, this additional complexity is polynomial. However, in some applications, the space of agent features could be infinite, e.g., with real-valued features. In this subsection, we investigate the situation when \mathcal{X} is infinite, through adapting the techniques from distributional robust optimization (Delage and Ye, 2010; Jiang and Guan, 2016; Ben-Tal et al., 2013).

To highlight the additional complexity of requiring robustness, in the following discussion, we assume that there exists an oracle that can provide agent's best response $(P_I; c_I)$ and compute the value C_I for any $(I; \cdot) \in \mathcal{G}^{\text{lin}}$ in time polynomial in n . Equipped with such an oracle, the problem defined in R-SO resembles the spirit of distributionally robust optimization (in short DRO), which aims to evaluate optimal solutions under the worst-case expectation with respect to a family of probability distributions of the uncertain parameters. The key concept in DRO is the ambiguity set, a family of measures consistent with the prior knowledge about uncertainty. In our formulation, the ambiguity set \mathcal{P} is specified via a hyperplane (see Eqn. (7)).

While our discussion so far applies for arbitrary utility functions $h(\cdot)$, analyzing the additional complexity of R-SO is challenging when the agent feature space \mathcal{X} is infinite. Therefore, in the rest of this subsection, we focus on a general set of concave and piecewise utility functions as defined below.

$$h(\mathbf{x}) = \min_{k \in [K]} h_k(\mathbf{x}); \quad (8)$$

Note that this set of utility functions is general since many commonly-seen utility functions are concave and can usually be approximated using simple piecewise functions, such as the piecewise linear functions: $h_k(\mathbf{x}) = \mathbf{a}_k^\top \mathbf{x} + b_k$, where for all $k \in [K]$, $\mathbf{a}_k \in \mathbb{R}^n$ and $b_k \in \mathbb{R}$ and the piecewise quadratic functions: $h_k(\mathbf{x}) = \min_{k \in [K]} \mathbf{x}^\top \mathbf{A}_k \mathbf{x} + \mathbf{b}_k^\top \mathbf{x} + c_k$ where for all $k \in [K]$, $\mathbf{A}_k \succeq 0$ and $\mathbf{A}_k \in \mathbb{R}^{n \times n}$; $\mathbf{b}_k \in \mathbb{R}^n$ and $c_k \in \mathbb{R}$.

Theorem 3. *Given that \mathcal{X} is ellipsoidal, i.e., $\mathcal{X} = \{\mathbf{x} \mid (\mathbf{x} - \mathbf{x}_0)^\top (\mathbf{x} - \mathbf{x}_0) \leq 1\}$, where \mathbf{x}_0 has at least one strictly positive eigenvalue, the objective of the problem R-SO is the same as the optimal value of the following optimization problem:*

When $h(\cdot)$ is a piecewise linear function:

$$\begin{aligned} \arg \min_{(I; \cdot) \in \mathcal{G}^{\text{lin}}; \cdot; \cdot} & C_I + \sum_{k \in [K]} \lambda_k (C_I - \mathbf{a}_k^\top \mathbf{x}_0) \\ \text{s.t.} & \sum_{k \in [K]} \lambda_k = 1 \\ & \sum_{k \in [K]} \lambda_k \mathbf{a}_k = \mathbf{0} \\ & \lambda_k \geq 0, \lambda_k \in \mathbb{R}; \quad k \in [K]; \end{aligned} \quad (9)$$

When $h(\cdot)$ is piecewise quadratic function, the first constraint will be replaced by the following

$$\begin{aligned} \arg \min_{(I; \cdot) \in \mathcal{G}^{\text{lin}}; \cdot; \cdot} & C_I + \sum_{k \in [K]} \lambda_k (C_I - \mathbf{b}_k^\top \mathbf{x}_0 - \mathbf{x}_0^\top \mathbf{A}_k \mathbf{x}_0) \\ \text{s.t.} & \sum_{k \in [K]} \lambda_k = 1 \\ & \sum_{k \in [K]} \lambda_k \mathbf{A}_k = \mathbf{0} \\ & \lambda_k \geq 0, \lambda_k \in \mathbb{R}; \quad k \in [K]; \end{aligned} \quad (10)$$

Proof. To solve R-SO , we first reformulate it as a minimization problem:

$$\min_{(I; \cdot) \in \mathcal{G}^{\text{lin}}} \max_{P \in \mathcal{P}} \mathbb{E}_P \max_{k \in [K]} h_k(\mathbf{x}); \quad (10)$$

For every $(I; \cdot) \in \mathcal{G}^{\text{lin}}$, let $(I; \cdot)_P$ denote the inner supremum problem in (10) over \mathcal{P} :

$$(I; \cdot)_P = \max_{P \in \mathcal{P}} \mathbb{E}_P \max_{k \in [K]} h_k(\mathbf{x}); \quad (11)$$

We can now recast the inner supremum problem $(I; \cdot)_P$ as a minimization problem, which can be performed jointly with the outer minimization over \mathcal{G}^{lin} . Introducing dual variables λ_k that correspond to the respective probability and expectation constraints in (7),

we have the following dual of (11):

$$\begin{aligned} \text{dual}(\beta; \gamma) \quad & \min_{\beta, \gamma} \quad \beta + (C_I \quad \gamma) \\ \text{s.t.} \quad & \beta + \beta^\top \mathbf{x} \leq h(\mathbf{x}); \delta \mathbf{x} \geq X \\ & \beta \geq 2R; \gamma \geq 2R; \end{aligned} \quad (12)$$

which provides an upper bound on (11). Indeed, consider any $P \geq P$ and any feasible solution $(\beta; \gamma)$ in problem (12); the robust counterpart in the dual implies that

$$E_P[h(\mathbf{x})] \leq E_P[\beta + \beta^\top \mathbf{x}] = \beta + (C_I \quad \gamma):$$

Thus, we have that weak duality holds: $\text{dual}(\beta; \gamma) \leq \text{val}(\text{P})$. Furthermore, the strong duality also holds since the problem (11) is a linear optimization problem (with respect to P). Having established the dual of (11) and its strong duality, we can formulate the problem (11) via a min-min operation that can be performed jointly over the constraint involving $h(\mathbf{x})$ decomposes.

$$\begin{aligned} \min_{(\beta; \gamma) \in \mathcal{G}^{\text{lin}}; \beta; \gamma} \quad & \beta + (C_I \quad \gamma) \\ \text{s.t.} \quad & \beta + \beta^\top \mathbf{x} + h_k(\mathbf{x}) \leq 0; \delta \mathbf{x} \geq X; k \in [K] \\ & \beta \geq 2R; \gamma \geq 2R; \end{aligned} \quad (13)$$

Note that when X has infinite elements, i.e., P is a measure with infinite support over X , there will be infinite-many constraints in (14). However, with our assumption on function $h(\cdot)$ and leveraging the geometry of X , we can reduce the above optimization problem with infinite-many constraints to the problem with tractable finite number of constraints. In particular, when X is ellipsoidal, i.e., $X = \{\mathbf{x} : (\mathbf{x} - \mathbf{x}_0)^\top (\mathbf{x} - \mathbf{x}_0) \leq 1\}$, and \mathbf{A}_k has at least one strictly positive eigenvalue, we can apply S-Lemma (cf., Theorem 2.2 in Pólik and Terlaky (2007)) for any given $k \in [K]$ to replace Constraint (14), which enforces that

$$\begin{aligned} \exists \lambda_k \geq 0 \text{ s.t.} \quad & \beta + \beta^\top \mathbf{x} + h_k(\mathbf{x}) \leq 0 \\ & \lambda_k (\mathbf{x} - \mathbf{x}_0)^\top (\mathbf{x} - \mathbf{x}_0) \leq 1 \end{aligned} \quad (15)$$

with the equivalent constraint that

$$\begin{aligned} \exists \lambda_k \geq 0 \text{ s.t.} \quad & \delta \mathbf{x} \geq 2R; \beta + \beta^\top \mathbf{x} + h_k(\mathbf{x}) \\ & \lambda_k (\mathbf{x} - \mathbf{x}_0)^\top (\mathbf{x} - \mathbf{x}_0) \leq 1 \end{aligned} \quad (16)$$

When $h_k(\mathbf{x}) = \mathbf{a}_k^\top \mathbf{x} + b_k$, then one can further use Schur's complement to replace Constraint (14) by an equivalent linear matrix inequality for any $k \in [K]$:

$$\begin{aligned} \text{"} \quad & \frac{\beta + \mathbf{a}_k^\top \mathbf{x}}{2} \leq \lambda_k \mathbf{x}_0^\top \mathbf{x}_0 \\ & \lambda_k \mathbf{x}_0^\top \mathbf{x}_0 + b_k + \lambda_k (\mathbf{x}_0^\top \mathbf{x}_0 - 1) \leq 0 \end{aligned} \quad \text{"}$$

The problem can therefore be reformulated as:

$$\begin{aligned} \min_{(\beta; \gamma) \in \mathcal{G}^{\text{lin}}; \beta; \gamma} \quad & \beta + (C_I \quad \gamma) \\ \text{s.t.} \quad & \frac{\beta + \mathbf{a}_k^\top \mathbf{x}}{2} \leq \lambda_k \mathbf{x}_0^\top \mathbf{x}_0 \\ & \lambda_k \mathbf{x}_0^\top \mathbf{x}_0 + b_k + \lambda_k (\mathbf{x}_0^\top \mathbf{x}_0 - 1) \leq 0; \delta \mathbf{x} \geq 2R; \lambda_k \geq 0; \delta \lambda_k \geq [K] \end{aligned} \quad \text{"}$$

where $\mathbf{x}_0 = (\mathbf{x}_1; \dots; \mathbf{x}_K)$. When $h_k(\mathbf{x}) = \min_{k \in [K]} \mathbf{x}^\top \mathbf{A}_k \mathbf{x} + \mathbf{b}_k^\top \mathbf{x} + c_k$ where for all $k \in [K]$, $\mathbf{A}_k \succeq 0$ and $\mathbf{A}_k \in \mathbb{R}^{n \times n}$; $\mathbf{b}_k \in \mathbb{R}^n$ and $c_k \in \mathbb{R}$. We then have following equivalent linear matrix inequality for Constraint (14) for any $k \in [K]$:

$$\begin{aligned} \text{"} \quad & \frac{\beta + \mathbf{a}_k^\top \mathbf{x}}{2} \leq \lambda_k \mathbf{x}_0^\top \mathbf{x}_0 \\ & \lambda_k \mathbf{x}_0^\top \mathbf{x}_0 + c_k + \lambda_k (\mathbf{x}_0^\top \mathbf{x}_0 - 1) \leq 0 \end{aligned} \quad \text{"}$$

The problem can therefore be reformulated as:

$$\begin{aligned} \min_{(\beta; \gamma) \in \mathcal{G}^{\text{lin}}; \beta; \gamma} \quad & \beta + (C_I \quad \gamma) \\ \text{s.t.} \quad & \frac{\beta + \mathbf{a}_k^\top \mathbf{x}}{2} \leq \lambda_k \mathbf{x}_0^\top \mathbf{x}_0 \\ & \lambda_k \mathbf{x}_0^\top \mathbf{x}_0 + c_k + \lambda_k (\mathbf{x}_0^\top \mathbf{x}_0 - 1) \leq 0; \delta \mathbf{x} \geq 2R; \lambda_k \geq 0; \delta \lambda_k \geq [K] \end{aligned} \quad \text{"}$$

where $\mathbf{x}_0 = (\mathbf{x}_1; \dots; \mathbf{x}_K)$. \square

In both linear and quadratic $h_k(\mathbf{x})$, we show that we can reformulate the original problem R-SO to the problem with a finite number of tractable linear matrix inequalities, instead of infinitely many constraints with X . This formulation provides a tool for us to analyze the additional complexity of R-SO compared with SO. For example, the problem in (9) could be possibly efficiently solvable when C_I of the agent's best response exhibits nice behaviors to retain a convexity of the objective in (9), e.g., when C_I is a linear form of β , then (9) is a semi-definite program. Then it is known that an interior point algorithm can be used to solve the above SDP with the polynomial time, i.e., the above problem can be solved to any precision ϵ in time polynomial in $\log(1/\epsilon)$ and the sizes of the problem. We leave the full characterization of conditions for the problem to be efficiently solvable for future work.

5 Discussions and Future Work

Linear models, one of the "white-box" models (contrary to the black-box models such as neural networks), have several desired properties such as nice generalizability, interpretability, transparency, and right to recourse. In this work, we further show that it is *robust* to unknown strategic manipulations when being used for making decisions. This is another dimension that is worth taking into account when deciding on which models to

deploy. While we demonstrate that finding the robust optimal decision rule is generally hard, our analysis in decomposing the problem could provide directions in figuring out efficient solvers in special cases.

There are still a number of open questions. In particular, our robustness notion could be overly pessimistic, considering the worst-case scenario over all possible unknown actions. One natural future direction is to explore Bayesian approaches, i.e., incorporating prior beliefs over all possible agent’s action sets, to model and quantify these uncertainties. Secondly, our work has focused on dealing with a single agent (or more broadly, a set of homogeneous agents: The decision-maker knows the *common* subset of all agents’ available actions). It would be interesting to extend the discussion to heterogeneous agents or a distribution of agents.

Acknowledgements

This work is supported in part by the Office of Naval Research Grant N00014-20-1-2240, National Science Foundation grant IIS-1939677, and Amazon.

References

- Tal Alon, Magdalen Dobson, Ariel D Procaccia, Inbal Talgam-Cohen, and Jamie Tucker-Foltz. Multiagent evaluation mechanisms. 2020.
- Pablo Azar, Silvio Micali, Constantinos Daskalakis, and S Matthew Weinberg. Optimal and efficient parametric auctions. In *Proceedings of the twenty-fourth annual ACM-SIAM symposium on Discrete algorithms*, pages 596–604. SIAM, 2013.
- Moshe Babaioff, Michal Feldman, and Noam Nisan. Combinatorial agency. In *Proceedings of the 7th ACM conference on Electronic commerce*, pages 18–28, 2006.
- Moshe Babaioff, Michal Feldman, and Noam Nisan. Mixed strategies in combinatorial agency. *Journal of Artificial Intelligence Research*, 38:339–369, 2010.
- Ian Ball. Scoring strategic agents. 2020.
- Chaithanya Bandi and Dimitris Bertsimas. Optimal design for multi-item auctions: a robust optimization approach. *Mathematics of Operations Research*, 39(4):1012–1038, 2014.
- Aharon Ben-Tal, Dick Den Hertog, Anja De Waegenare, Bertrand Melenberg, and Gijs Rennen. Robust solutions of optimization problems affected by uncertain probabilities. *Management Science*, 59(2):341–357, 2013.
- Michael Brückner and Tobias Scheffer. Stackelberg games for adversarial prediction problems. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 547–555, 2011.
- Michael Brückner, Christian Kanzow, and Tobias Scheffer. Static prediction games for adversarial learning problems. *Journal of Machine Learning Research*, 13(Sep):2617–2654, 2012.
- Gabriel Carroll. Robustness and linear contracts. *American Economic Review*, 105(2):536–63, 2015.
- Gabriel Carroll. Robustness and separation in multidimensional screening. *Econometrica*, 85(2):453–488, 2017.
- Gabriel Carroll and Ilya Segal. Robustly optimal auctions with unknown resale opportunities. *The Review of Economic Studies*, 86(4):1527–1555, 2019.
- Sylvain Chassang. Calibrated incentive contracts. *Econometrica*, 81(5):1935–1971, 2013.
- Tianjiao Dai and Juuso Toikka. Robust incentives for teams. *Unpublished manuscript, Mass. Inst. of Technology, Cambridge, MA*, 2017.
- Erick Delage and Yinyu Ye. Distributionally robust optimization under moment uncertainty with application to data-driven problems. *Operations research*, 58(3):595–612, 2010.
- Peter Diamond. Managerial incentives: on the near linearity of optimal compensation. *Journal of Political Economy*, 106(5):931–957, 1998.
- Jinshuo Dong, Aaron Roth, Zachary Schutzman, Bo Waggoner, and Zhiwei Steven Wu. Strategic classification from revealed preferences. In *Proceedings of the 2018 ACM Conference on Economics and Computation*, pages 55–70, 2018.
- Paul Dütting, Tim Roughgarden, and Inbal Talgam-Cohen. Simple versus optimal contracts. In *Proceedings of the 2019 ACM Conference on Economics and Computation*, pages 369–387, 2019.
- Allussein Fawzi, Omar Fawzi, and Pascal Frossard. Analysis of classifiers’ robustness to adversarial perturbations. *Machine Learning*, 107(3):481–508, 2018.
- Bryce Goodman and Seth Flaxman. European union regulations on algorithmic decision-making and a “right to explanation”. *AI magazine*, 38(3):50–57, 2017.
- Sanford J Grossman and Oliver D Hart. An analysis of the principal-agent problem. In *Foundations of Insurance Economics*, pages 302–340. Springer, 1992.
- Nika Haghtalab, Nicole Immorlica, Brendan Lucier, and Jack Wang. Maximizing welfare with incentive-aware evaluation mechanisms. In *29th International Joint Conference on Artificial Intelligence*, 2020.
- Lars Peter Hansen and Thomas J Sargent. Three types of ambiguity. *Journal of Monetary Economics*, 59(5):422–445, 2012.

- Pierre Hansen, Brigitte Jaumard, and Gilles Savard. New branch-and-bound rules for linear bilevel programming. *SIAM Journal on scientific and Statistical Computing*, 13(5):1194–1217, 1992.
- Moritz Hardt, Nimrod Megiddo, Christos Papadimitriou, and Mary Wootters. Strategic classification. In *Proceedings of the 2016 ACM conference on Innovations in Theoretical Computer Science*, pages 111–122, 2016.
- Chien-Ju Ho, Aleksandrs Slivkins, and Jennifer Wortman Vaughan. Adaptive contract design for crowdsourcing markets: Bandit algorithms for repeated principal-agent problems. *Journal of Artificial Intelligence Research*, 55:317–359, 2016.
- Bengt Holmstrom and Paul Milgrom. Aggregation and linearity in the provision of intertemporal incentives. *Econometrica: Journal of the Econometric Society*, pages 303–328, 1987.
- Lily Hu, Nicole Immorlica, and Jennifer Wortman Vaughan. The disparate effects of strategic manipulation. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 259–268, 2019.
- Ruiwei Jiang and Yongpei Guan. Data-driven chance constrained stochastic program. *Mathematical Programming*, 158(1-2):291–327, 2016.
- Jon Kleinberg and Manish Raghavan. How do classifiers induce agents to invest effort strategically? In *Proceedings of the 2019 ACM Conference on Economics and Computation*, pages 825–844, 2019.
- Jianjun Miao and Alejandro Rivera. Robust contracts in continuous time. *Econometrica*, 84(4):1405–1440, 2016.
- John Miller, Smitha Milli, and Moritz Hardt. Strategic adaptation to classifiers: A causal perspective. *arXiv preprint arXiv:1910.10362*, 2019.
- Smitha Milli, John Miller, Anca D Dragan, and Moritz Hardt. The social cost of strategic classification. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 230–239, 2019.
- Imre Pólik and Tamás Terlaky. A survey of the s-lemma. *SIAM review*, 49(3):371–418, 2007.
- Steven Shavell. Risk sharing and incentives in the principal and agent relationship. *The Bell Journal of Economics*, pages 55–73, 1979.
- Behzad Tabibian, Stratis Tsirtsis, Moein Khajehnejad, Adish Singla, Bernhard Schölkopf, and Manuel Gomez-Rodriguez. Optimal decision making under strategic behavior. *arXiv preprint arXiv:1905.09239*, 2019.
- Berk Ustun and Cynthia Rudin. Methods and models for interpretable linear classification. *arXiv preprint arXiv:1405.4047*, 2014.
- Berk Ustun, Alexander Spangher, and Yang Liu. Actionable recourse in linear classification. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 10–19, 2019.
- Jiaxuan Wang, Jeeheh Oh, Haozhu Wang, and Jenna Wiens. Learning credible models. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2417–2426, 2018.

Linear Models are Robust Optimal Under Strategic Behavior: Supplementary Materials

6 Proof for Lemma 1

Proof. Let us first fix an arbitrary action set $A_a \subseteq A_d$, and a rational decision rule f . We must have that the agent's utility is at least $V_a(fjA_d)$, that is, any action $(P; c)$ the agent would chose under the decision rule f must satisfy:

$$E_P[f(\mathbf{x})] \geq E_P[h(\mathbf{x})] \quad c = V_a(fjA_a) \geq V_a(fjA_d):$$

Thus, the decision maker's utility $V_d(fjA_a) = E_P[h(\mathbf{x})]$ is at the least the minimum given by the (4). This implies the following guarantee of worst-case utility $V_d(f)$:

$$V_d(f) \geq \min_{P \in \mathcal{X}} E_P[h(\mathbf{x})] \quad \text{s.t.} \quad E_P[f(\mathbf{x})] \geq V_a(fjA_d): \quad (17)$$

We now show that (17) is tight. Let $\text{supp}(P)$ denote the support of distribution P . Let P_0 be a distribution attaining the minimum in (4) and also satisfying the constraint. We consider following two cases:

Case 1: $\text{supp}(P_0) \not\subseteq \arg \max_{\mathbf{x}} f(\mathbf{x})$. Then let P_1 be a distribution which achieves a higher value of $E_P[f(\mathbf{x})]$. Let P' be a mixture distribution $P' = (1 - \epsilon)P_0 + \epsilon P_1$, with a small positive ϵ . Then we have $E_{P_0}[f(\mathbf{x})] = (1 - \epsilon)E_{P_0}[f(\mathbf{x})] + \epsilon E_{P_1}[f(\mathbf{x})] > E_{P_0}[f(\mathbf{x})]$. Now take $A'_a = A_d \setminus [f(P'; 0)g]$, then the agent's unique optimal action under A'_a is $(P'; 0)$. This brings the decision maker with utility of $V_d(fjA'_a) = (1 - \epsilon)E_{P_0}[h(\mathbf{x})] + \epsilon E_{P_1}[h(\mathbf{x})]$. Since $V_d(fjA'_a) \geq V_d(f)$, we further have

$$V_d(f) \leq V_d(fjA'_a) = (1 - \epsilon)E_{P_0}[h(\mathbf{x})] + \epsilon E_{P_1}[h(\mathbf{x})]: \quad (18)$$

When $\epsilon \rightarrow 0$, the RHS in (18) will converge to $E_{P_0}[h(\mathbf{x})]$. This implies $V_d(f) \leq E_{P_0}[h(\mathbf{x})]$ when $\epsilon \rightarrow 0$. Recall our definition of P_0 , and together with the lower bound we have shown for $V_d(f)$ in (17), we can conclude our results in (4) for this case.

Case 2: $\text{supp}(P_0) \subseteq \arg \max_{\mathbf{x}} f(\mathbf{x})$. For this case, we discuss following two situations.

(i): $E_{P_0}[f(\mathbf{x})] > V_a(fjA_d)$, we now consider action set $A'_a = A_d \setminus [f(P_0; 0)g]$. Since $E_{P_0}[f(\mathbf{x})] > V_a(fjA_d)$, then the agent will uniquely chose action $(P_0; 0)$ for f under the action set A'_a . This brings the decision maker with the utility of $V_d(fjA'_a) = E_{P_0}[h(\mathbf{x})]$. Again, with the fact that $V_d(fjA'_a) \geq V_d(f)$ and the definition of P_0 , we have now proved (4).

(ii): $E_{P_0}[f(\mathbf{x})] = V_a(fjA_d) = \max f(\mathbf{x})$, this situation can only be satisfied when A_d contains some action of the form $(P'; 0)$ with $\text{supp}(P') \not\subseteq \arg \max f(\mathbf{x})$. Thus, we define

$$G := \{f(P'; 0) \geq V_a(fjA_d) : \text{supp}(P') \not\subseteq \arg \max f(\mathbf{x})\} \neq \emptyset:$$

Then, under action set A_d , the agent will choose an action in G which would benefit decision maker (according to the tie-breaking assumption, when there are multiple optimal actions for agent, agent will choose the one which maximizes decision maker's utility.), leading the decision maker's utility $V_d(fjA_d) = \max_{(P; 0) \in G} E_P[h(\mathbf{x})] \geq V_d(f)$. In this scenario, the unique optimal action for the agent under any action set $A \subseteq A_d$ is some $(P; 0) \in G$. However, the agent would stick to the same action even under zero decision rule (recall our tie-breaking assumption), leading the decision maker's utility $V_d(0jA) = \max_{(P; 0) \in G} E_P[h(\mathbf{x})] = V_d(0)$. This implies $V_d(0) \geq V_d(f)$, which contradicts our rationality assumption.

Now we establish the equality claims. Without loss of generality, we may assume the agent has a costless action $(\underline{x}; 0)$ in A_d where $h(\underline{\mathbf{x}}) = 0$.⁴ Recall that we have $E_{P_0}[h(\mathbf{x})] = V_d(f) > V_d(0) > 0$ by our assumption on P_0 and DM's rationality. If we have $E_{P_0}[f(\mathbf{x})] > V_a(fjA_d)$ strictly, then replace P_0 by a mixture distribution $P' = (1 - \epsilon)P_0 + \epsilon \delta_{\underline{\mathbf{x}}}$ for small ϵ . Consider $A'_a = A_d \setminus [f(P'; 0)g]$, then the agent's utility by taking the action $(P'; 0)$ is given by $V_a(fjA'_a) = (1 - \epsilon)E_{P_0}[f(\mathbf{x})] + \epsilon f(\underline{\mathbf{x}})$, then one can always find a small ϵ such that $V_a(fjA'_a)$ is strictly larger than $V_a(fjA_d)$. As a result, this brings the decision maker with a utility of $V_d(fjA'_a) = (1 - \epsilon)E_{P_0}[h(\mathbf{x})] + \epsilon h(\underline{\mathbf{x}}) = (1 - \epsilon)E_{P_0}[h(\mathbf{x})]$. Since $V_d(fjA'_a) > V_d(f)$, given any positive ϵ , this implies that $V_d(f) = (1 - \epsilon)E_{P_0}[h(\mathbf{x})] < E_{P_0}[h(\mathbf{x})]$, which contradicts the minimality of P_0 . Thus we have $E_{P_0}[f(\mathbf{x})] = V_a(fjA_d)$. Finally, if $P_0 \not\geq \arg \max_{P \in \mathcal{P}(\mathcal{X})} E_P[f(\mathbf{x})]$, and $E_{P_0}[f(\mathbf{x})] = V_a(fjA_d)$, then we have (5). \square

After finishing the proof, we would like to give following explanation on our construction of worst-case action set in the proof.

Remark 1. *The above proof relies on a construction of agent's worst case action set by adding an arbitrary action of the form $(P; 0)$. It may seem unrealistic to allow the agent to arbitrarily manipulate himself at zero cost. However, we note that the zero cost is not a substantive assumption: the logic can be carried over to more realistic models that can explicitly incorporate the effort costs as a function of expected manipulated feature. Then the equivalent step consists of adding an action to the action set that produces P at the lowest allowable cost.*

7 Proof for Lemma 2

Proof. Our proof structure is similar to [Carroll \(2015\)](#), with the key difference on how to define the two disjoint convex sets. Suppose that the convex hull of X is a full-dimensional set in \mathbb{R}^n . Now fix any nonlinear decision rule f , our proof will hinge on the discussion of two cases we have shown in Lemma 1.

Case 1. We first define

$$t(\mathbf{x}) = \max_{f \in V_a(fjA_d)} h(\mathbf{x}) + f(\mathbf{x}) - V_d(f)g;$$

Now we define two sets in $\mathbb{R}^{n+1} = \mathbb{R}^n \times \mathbb{R}$: Let S be the convex hull of all pairs $(\mathbf{x}; f(\mathbf{x}))$, for $\mathbf{x} \in X$, let T be the convex hull of all pairs $(\mathbf{x}; z)$ that \mathbf{x} lies in the convex hull of X , and $z > t(\mathbf{x})$. We note that T is then a convex set. A graph illustration of our proof is presented in Figure 3.

We now claim that S and T are disjoint. To see this, suppose S and T are not disjoint, then there exists a distribution $P \in \mathcal{P}(X)$ such that $E_P[f(\mathbf{x})] > E_P[t(\mathbf{x})]$. In particular, we have

$$E_P[f(\mathbf{x})] > V_a(fjA_d);$$

and also

$$\begin{aligned} E_P[f(\mathbf{x})] &> E_P[h(\mathbf{x})] + E_P[f(\mathbf{x})] - V_d(f) \\ &\implies V_d(f) > E_P[h(\mathbf{x})]: \end{aligned}$$

This is a direct contradiction to our statement of (4) in Lemma 1.

The disjointness and convexity of S and T enable us to apply the separating hyperplane theorem: There exists a

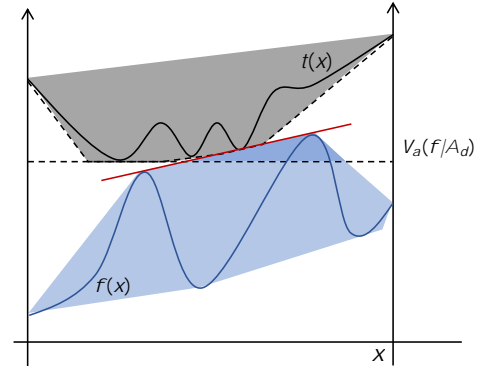


Figure 3: Illustrate S and T when $n = 1$. The blue line is $f(\mathbf{x})$ and its associated convex hull in blue shaded region (the top blue triangle is the set S). Black line is $t(\mathbf{x})$. The black shaded region is the convex hull for all points $(\mathbf{x}; z)$ where $z > t(\mathbf{x})$. The red line is the hyperplane to separate S and T .

⁴This assumption is merely an additive normalization of the decision maker's utility and it can be relaxed to a more general scenario where our results still hold (see our discussion at the end of the Appendix 7). Earlier works also make similar assumption ([Carroll, 2015](#); [Dütting et al., 2019](#)): The agent can always exert no effort, namely, the zero-cost action, to produce a minimum output (denote by 0); this corresponds to assuming $(\underline{o}; 0) \in A_d$.

vector $\mathbf{x} = (x_1, \dots, x_n)$ and constants $v_i > 0$ such that

$$\sum_i x_i + z \geq v_i; \quad \delta(\mathbf{x}; z) \geq S \quad (19)$$

$$\sum_i x_i + z' \geq v_i; \quad \delta(\mathbf{x}; z') \geq T \quad (20)$$

and \mathbf{x} is a non-zero vector. Note that (19) and (20) implies $v_i > 0$. To see this, fix a point $\mathbf{x} \geq X$, then for $(\mathbf{x}; z) \geq S$ and $(\mathbf{x}; z') \geq T$ we have

$$\sum_i x_i + z' \geq \sum_i x_i + z \geq z' - z;$$

by earlier argument on the disjointness of S and T , we can conclude that $v_i > 0$. We now also show that v_i is a positive constant. Suppose $v_i = 0$, then (19) gives $\sum_j x_j \geq v$ and (20) gives $\sum_j x_j \geq v$, which leads to $\sum_j x_j = v$. Since not all x_j are zero, this contradicts the full-dimensionality of X .

Now we can rewrite (19) as following

$$f(\mathbf{x}) = \frac{v}{\sum_i x_i}; \quad \delta \mathbf{x} \geq X;$$

This motivates us to define following linear decision rule

$$f'(\mathbf{x}) = \frac{v}{\sum_i x_i}; \quad \delta \mathbf{x} \geq X; \quad (21)$$

Note that we have $f'(\mathbf{x}) \geq f(\mathbf{x})$ pointwise.

Now we are ready to check that $V_d(f') \geq V_d(f)$. Let $(P_0; c_0)$ be the action that the agent would like to choose under f and action set A_d . Consider any action set $A_a \subseteq A_d$, as we have shown before, we must have

$$V_a(f' | A_a) \geq V_a(f' | A_d) \geq V_a(f | A_d); \quad (22)$$

Let $(P; c)$ be the action that the agent chooses under f' and action set A_a . Then (20) implies

$$\begin{aligned} E_P[t(\mathbf{x})] &= \frac{v}{\sum_i E_P[x_i]} \\ &= E_P[f'(\mathbf{x})] \\ &= V_a(f' | A_a) + c \\ &\geq V_a(f' | A_a) && (c \geq R_+) \\ &\geq V_a(f | A_d); && \text{(by (22))} \end{aligned} \quad (23)$$

It is worthy noting that if above inequality is strict, then according to our definition of $t(\mathbf{x})$, we must have

$$E_P[t(\mathbf{x})] = E_P[h(\mathbf{x})] + E_P[f(\mathbf{x})] \geq V_d(f); \quad (24)$$

So we have

$$\begin{aligned} V_d(f' | A_a) &= E_P[h(\mathbf{x})] = E_P[t(\mathbf{x})] - E_P[f(\mathbf{x})] + V_d(f) \\ &\geq E_P[t(\mathbf{x})] - E_P[f'(\mathbf{x})] + V_d(f) && \text{(by definition of } f') \\ &\geq V_d(f); && \text{(by (23))} \end{aligned}$$

On the other hand, if $E_P[t(\mathbf{x})] = V_a(f | A_d)$. This implies all the inequalities in the stacked chain above are equalities. In particular, we will have

$$V_a(f' | A_a) = V_a(f' | A_d) = V_a(f | A_d);$$

Since the agent now does at least as well as $V_a(f | A_d)$ by taking action $(P_0; c_0)$, this action is in his choice set under f' and A_a , as a result, the decision maker gets at least the corresponding utility: $V_d(f' | A_a) \geq E_{P_0}[h(\mathbf{x})] =$

$V_d(fjA_d) \geq V_d(f)$, where the first inequality is due to the tie-breaking assumption of the agent (when there are multiple maximizers, the agent will chose the most beneficial one for the decision maker).

Thus, in either case, we have $V_d(f'jA_a) \geq V_d(f)$, this holds for any $A_a \in A_d$, thus we have $V_d(f') \geq V_d(f)$.

Case 2. In this case, we define S to be the convex hull of all pairs $(\mathbf{x}; f(\mathbf{x}))$, and T to be the set of all $(\mathbf{x}; z)$ with \mathbf{x} in the convex hull of X and $z > V_a(fjA_d)$. We still claim both of S and T are convex, and disjoint: otherwise, there exists P such that

$$E_P[f(\mathbf{x})] > V_a(fjA_d):$$

This contradicts our statement (5) in Lemma 1. Using the same arguments as in case 1, we find a vector $\mathbf{c} = (c_1; \dots; c_n)$ and constants $\epsilon; \nu$ such that (19) and (20) hold, and we can still guarantee that $\epsilon > 0$. Again, we define a linear decision rule f' by (21); from (19) we know that $f' \geq f$ pointwise. Consider the agent's behavior under decision rule f' , for any action $(P; c)$ chosen by the agent under any possible action set, we have

$$E_P[f'(\mathbf{x})] - c = f'(E_P[\mathbf{x}]) - c \geq V_a(fjA_d): \quad (\text{by (20)})$$

This means that the agent cannot earn a higher expected utility than $V_a(fjA_d)$. On the other hand, the agent can always earn at least this much, since $V_a(f'jA_a) \geq V_a(f'jA_d) \geq V_a(fjA_d)$. This means we have equality $V_a(f'jA_a) = V_a(f'jA_d) = V_a(fjA_d)$. From here, the argument finishes just as at the end of case 1, and we have $V_d(f') \geq V_d(f)$. \square

Extensions: General cost lower bounds As mentioned in Remark 1, our analysis relies on the construction of worst case action sets, using actions, that produce an undesirable distribution P , at costs of zero. This zero-cost action assumption (together with the assumption in Footnote 6) is not substantial and one natural relaxation is that the decision maker knows a lower bound on the cost of any available actions, or of producing any given level of expected output. Our analysis and results will go through for this scenario. Specifically, suppose the known lower bound cost is denoted by $\underline{c} > 0$, then our Lemma 1 can be accordingly changed to: $V_d(f) = \min_{P \in \mathcal{C}(X)} E_P[h(\mathbf{x})]$; s.t. $E_P[f(\mathbf{x})] \leq \underline{c} + V_a(fjA_d)$ or $\max_{P \in \mathcal{C}(X)} E_P[f(\mathbf{x})] \leq \underline{c} + V_a(fjA_d)$. To get the analogous result in Lemma 2, one can change the function $t(\mathbf{x})$ as $t(\mathbf{x}) = \max_{f \in \mathcal{F}} V_a(fjA_d) + \underline{c}; h(\mathbf{x}) + f(\mathbf{x}) - V_d(f)g$, then all the analysis can be carried over here.

8 Proof for Lemma 3

Proof. We prove Theorem 1 via showing the existence of an optimum within the class of linear decision rules, and this decision rule will then be optimal among all decision rules. Note that for any rational decision rule $f(\mathbf{x})$, the value of $f(\mathbf{x})$ that it assigns to \mathbf{x} is bounded within $(0; C]$. Let a linear decision rule be the form of $f_{(I; \gamma)}(\mathbf{x}) = \mathbf{I}^\top \mathbf{x} + \gamma$. Then it suffices to show that the guaranteed worst-case utility $V_d(f)$ is an upper semi-continuous function of $(\mathbf{I}; \gamma) \in G^{\text{lin}}$. Now fix a sequence $(\mathbf{I}^1; \gamma^1); (\mathbf{I}^2; \gamma^2); \dots$ in G^{lin} converging to some $(\mathbf{I}^\infty; \gamma^\infty)$ in G^{lin} . Then it suffices to show that $V_d(f_{(I^1; \gamma^1)}) \leq \limsup_k V_d(f_{(I^k; \gamma^k)})$. To prove this, first note that by replacing the sequence $((\mathbf{I}^k; \gamma^k))$ with a subsequence along which $V_d(f_{(I^k; \gamma^k)})$ converges to its lim sup on the original sequence, thus, we can assume that $V_d(f_{(I^k; \gamma^k)})$ converges to $\limsup_k V_d(f_{(I^k; \gamma^k)})$. Now for any action set A_a , and let $(P^k; c^k)$ be the agent's chosen action under A_a and the decision rule $f_{(I^k; \gamma^k)}$. Then if necessary, by extracting a further subsequence, we can assume that the sequence $(P^k; c^k)$ converges to some $(P^\infty; c^\infty) \in A_a$. Since the agents' utility are continuous in $(\mathbf{I}; \gamma)$, then $(P^\infty; c^\infty)$ is an optimal action for the agent under $f_{(I^1; \gamma^1)}$, and its utility to the decision maker is the limit of the corresponding utility of $(P^k; c^k)$ under $f_{(I^k; \gamma^k)}$. We thus have

$$V_d(f_{(I^1; \gamma^1)}jA_a) = E_{P^1}[h(\mathbf{x})] = \lim_k E_{P^k}[h(\mathbf{x})] = \lim_k V_d(f_{(I^k; \gamma^k)}jA_a) = \lim_k V_d(f_{(I^k; \gamma^k)}):$$

Since $A_a \in A_d$ is arbitrary, then we have $V_d(f_{(I^1; \gamma^1)}) \leq \lim_k V_d(f_{(I^k; \gamma^k)})$. \square

9 Missing Table in Section 3.4

Given the student's efforts \mathbf{e} invested to each action, we can enumerate all possible induced distributions over X in A_d and A_a (see Table 1). Note that since the student can now also invest efforts to action a_0 , A_a contains more availabilities compared to A_d .

$\mathbf{x} = (x_1; x_2)$	P in A_d	P in A_a
$\Pr(\mathbf{x} = (1;1))$	$e_1 p^2$	$(e_1 p + (p - e_0)(p + e_0))$
$\Pr(\mathbf{x} = (1;0))$	$e_1 p(1 - p)$	$(e_1 p + (p - e_0)(1 - p - e_0))$
$\Pr(\mathbf{x} = (0;1))$	$(1 - e_1 p)p$	$(1 - e_1 p - (p - e_0)(p + e_0))$
$\Pr(\mathbf{x} = (0;0))$	$(1 - e_1 p)(1 - p)$	$(1 - e_1 p - (p - e_0)(1 - p - e_0))$

Table 1: All possible distributions P in A_d and A_a induced by student's effort $\mathbf{e} = (e_0; e_1; 1 - e_0 - e_1)$. e_1, e_0 are the efforts decided by the student for actions a_1 and a_0 , and $e_1 + e_0 \geq [0; 1]$.

10 Missing proof and the Algorithm for Theorem 2

Algorithm 1 Find the optimal robust decision rule

- 1: Input: Decision maker's knowledge A_d , linear decision space G^{lin} , objective function h .
- 2: Initial $f^* \in G^{\text{lin}}$ arbitrarily and $V_d(f^*) = 0$.
- 3: **for** every $(f; c) \in G^{\text{lin}}$ **do**
- 4: Let $(P_0; c_0) \in \arg \max_{(P; c) \in A_d} E_P[f^\top \mathbf{x} + c]$;
- 5: Solve the set $P = \{P : E_{P_0}[\mathbf{x}] - E_P[\mathbf{x}] = c_0 - c; P \in (X)\}$;
- 6: Compute $V_d(f; c) = \min_{P \in P} E_P[h(\mathbf{x})]$;
- 7: **if** $V_d(f; c) > V_d(f^*)$ **then**
- 8: $f^* = f^\top \mathbf{x} + c$.
- 9: **end if**
- 10: **end for**
- 11: **Output** Robust optimal decision: f^* .

Proof. According Lemma 1, given $f_{(f; c)}$, for any distribution P attaining the minimum in (4), we know that the inequality in (4) must bind at P. Let $(P_f; c_f) \in A_d$ be the solution to the constraint in S0. Then we can compute f^* by solving:

$$\begin{aligned} & \arg \max_{(f; c) \in G^{\text{lin}}} \min_{P \in \mathcal{P}} E_P[h(\mathbf{x})]; & (\text{R-S0}) \\ \text{s.t. } & P = P' : E_{P_0}[f_{(f; c)}(\mathbf{x})] = f^\top E_{P'}[\mathbf{x}] - c_f = c_f; P' \in (X); & (25) \end{aligned}$$

where we refer to the set \mathcal{P} , as the *worst-action set*, since we choose the worst action among it to minimize the expected utility $E_P[h(\mathbf{x})]$. Different from the problem in S0, after identifying the agent's best response $(P_f; c_f) \in A_d$ under $f_{(f; c)}$, our problem in R-S0 first turns to characterizing a worst-action set \mathcal{P} . Then the searching of f^* will hinge on maximizing $E_P[h(\mathbf{x})]$ in each P over G^{lin} . This implies that to make our problem tractable, one may first need to guarantee the corresponding strategic decision-making problem tractable. Furthermore, given a linear $f_{(f; c)}$, the additional computational complexity in R-S0 is due to the robustness concern in minimizing $E_P[h(\mathbf{x})]$ over set \mathcal{P} . It is easy to see that this is a linear programming with equality constraint, where the decision variables are a probability simplex over X .

$$\min_{P \in \mathcal{P}} E_P[h(\mathbf{x})]; \quad \text{s.t. } P = P' : f^\top E_{P'}[\mathbf{x}] = c_f; P' \in (X); \quad (26)$$

Inside the optimization, for every $(f; c) \in G^{\text{lin}}$, our problem R-S0 has one more induced Linear programming to solve compared with the standard problem S0.

As it will in general be hard to optimize arbitrary non-concave functions, we may consider assuming a concave h . However, as pointed out by other studies (Kleinberg and Raghavan, 2019; Alon et al., 2020), there exist concave functions h that are NP-hard to solve the problem S0 (via a reduction from the maximum independent set problem), which naturally leads the hardness of our problem. In particular, back to our student evaluation setting, let $P(\mathbf{e})$ be the induced feature distribution if the agent's effort profile is \mathbf{e} . As a result, the decision maker's goal on maximizing $h(\mathbf{x})$ can be reduced to maximizing $h(P(\mathbf{e}))$. When $h(P(\mathbf{e})) = k e k_0$, solving the problem S0 is then NP-hard. \square