

INCENTIVE MECHANISMS IN STRATEGIC LEARNING

Kun Jin¹, Xueru Zhang², Mohammad Mahdi Khalili³, Parinaz Naghizadeh⁴, and Mingyan Liu¹

¹EECS, University of Michigan, ²CSE, Ohio State University, ³CS, University of Delaware, ⁴ISE & ECE, Ohio State University

ABSTRACT

We study the design of a class of incentive mechanisms that can effectively improve algorithm robustness in strategic learning. A conventional strategic learning problem is modeled as a Stackelberg game between an algorithm designer (a principal, or decision maker) and individual agents subject to the algorithm’s decisions, potentially from different demographic groups. While the former benefits from the decision accuracy, the latter may have an incentive to game the algorithm into making favorable but erroneous decisions by merely changing their observable features without affecting their true labels. While prior works tend to focus on how to design decision rules robust to such strategic maneuvering, this study focuses on an alternative, which is to design incentive mechanisms to shape the utilities of the agents and induce improvement actions that genuinely improve their skills and true labels and thus in turn benefit both parties in the Stackelberg game. Specifically, the principal and the mechanism provider (could be the principal itself) move together in the first stage, publishing and committing to a classifier and an incentive mechanism. The agents are second movers and best respond to the published classifier and incentive mechanism. We study how the mechanism can induce improvement actions, positively impact a number of social well-being metrics, such as the overall skill levels of the agents (efficiency) and positive or true positive rate differences between different demographic groups (fairness).

1 INTRODUCTION

This work presents the design of a subsidy mechanism and its impacts in strategic classification.¹ Conventional strategic classification model the interaction between a principal and agents who are subject to the decision outcomes. While the former benefits from the classification accuracy, the latter may have an incentive to *game* the classifier into making favorable but erroneous decisions. Recognizing the potential issue, prior works focus on designing algorithms that are more robust to such strategic maneuvering, see e.g., (Hardt et al., 2016a; Milli et al., 2019; Hu et al., 2019; Brückner & Scheffer, 2011; Brückner et al., 2012; Dong et al., 2018; Braverman & Garg, 2020; Chen et al., 2020; Miller et al., 2020). Equally important, however, is the possibility for a mechanism designer to *incentivize* effort by the users who genuinely improve their true label; this benefits the users by increasing utilities and the principal by preserving the algorithm performance at the same time.

Toward this end, we present a strategic classification problem augmented by a subsidy mechanism (augmented strategic learning problem) modeled as a Stackelberg game between the principal, the mechanism designer (which could be the principal itself or a third party) and individuals from different demographic groups who are subject to the classifiers’ decisions (the agents). The principal and the mechanism designer move in the first stage by publishing and committing to a binary classifier and an incentive mechanism. The published classifier takes as input the agents’ *observable* features and outputs decision outcomes that impact the agents’ utilities. The agents are (simultaneous) second movers and best respond to the published decision rule and incentive mechanism. To capture the agent’s ability to both game the decision rule and make real changes, we assume each agent has an endowed pre-response attribute (endowed private information), that is causal (Miller et al. (2020)) to a set of observable features as well as its true label, also referred to as its *qualification state* in the context of the strategic learning problem. An agent can exert effort to improve this causal state,

¹Our study on strategic regression has similar results and is covered in the appendix

thereby improving its features and its underlying attributes, or choose to game the classifier by employing non-causal schemes to improve only its features without changing its underlying attributes, or using a combination of them. Both choices of action, referred to as *improvement* and *gaming*, respectively, come at a cost to the agent. The principal derives its utility from the prediction accuracy, thus even a selfish principal may have an incentive to motivate the agents to choose improvements over gaming to preserve the classifier’s performance. When the principal is also the mechanism designer, one such incentive mechanism is for the principal to subsidize the agents’ improvement costs, thereby making improvement more appealing compared to gaming.

The main contribution of our work is summarized as follows. We analyze the design of the subsidy mechanism and characterize the Stackelberg equilibrium in the augmented strategic classification. The design of an optimal mechanism requires solving non-convex problems but we have polynomial time and even constant time solutions in some realistic special cases. We also study the case where the mechanism designer is a third party (e.g., a government) with social well-being metrics (can be efficiency or fairness oriented) as its objective. The third party designs a mechanism that incentivizes agents’ improvement action and charges a tax from the principal for this *improvement service* to ensure budget balance, while also making sure that incentive compatibility and individual rationality constraints are satisfied for both the agents and the principal. In addition, we study the impact of the mechanism designer’s objective and the corresponding mechanism on the fairness and qualification status, when agents come from different demographic groups which differ in endowment or action cost. We show with analytical and numerical results that the outcomes of the augmented strategic learning with a fairness oriented third party are the most well-rounded since it improves the fairness, the qualification status, and the robustness of the classifier.

2 MODEL FOR AUGMENTED STRATEGIC CLASSIFICATION

We model the augmented strategic classification problem as a single-round, two-stage Stackelberg game, where the principal and the mechanism designer move first to design, publish, and commit to a classifier $f = \mathbf{1}\{\mathbf{w}^T \mathbf{x} \geq \tau\}$, $\mathbf{w} \geq 0$, $\tau \geq 0$ combined with an incentive mechanism G ; the agents then best respond to both f and G in the second stage. Figure 1 illustrates the augmented strategic classification problem where the principal is the mechanism designer.

Attributes, Features, and Labels. An agent has an N -dimensional *pre-response attribute* $\mathbf{x} \in \mathcal{X}$, $\mathcal{X} \subseteq \mathbb{R}_{\geq 0}^N$, which is its private information. Its probability density function (pdf) is $p(\mathbf{x})$, which is public information. In the response phase, an agent takes an M -dimensional action $\mathbf{a} := (\mathbf{a}_+, \mathbf{a}_-)$, where $\mathbf{a}_+ \in \mathbb{R}_{\geq 0}^{M_+}$ denotes an *improvement action* profile while $\mathbf{a}_- \in \mathbb{R}_{\geq 0}^{M_-}$ is a *gaming action* profile, with $M_+ + M_- = M$ (with action indices ordered such that $\forall i \leq M_+$ is an improvement action).

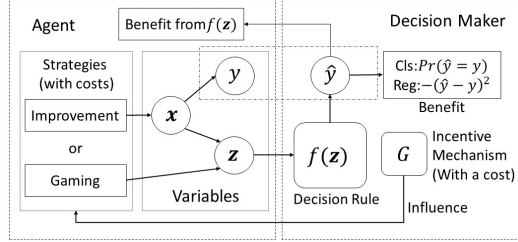


Figure 1: Augmented Strategic Classification

The agent’s action impacts its attribute as well as feature through a *projection matrix* $P = [P_+, P_-]$, $P \geq 0$, where $P_+ \in \mathbb{R}^{N \times M_+}$ (resp. $P_- \in \mathbb{R}^{N \times M_-}$) is the improvement (resp. gaming) projection in the following sense. The action results in the agent having a *post-response attribute* $\mathbf{x}' = \mathbf{x} + P_+ \mathbf{a}_+ = \mathbf{x} + \hat{P} \mathbf{a}$, where $\hat{P} = [P_+, \mathbf{0}] \in \mathbb{R}^{N \times M}$, and a *post-response observable feature* (simply feature for brevity) $\mathbf{z} = \mathbf{x} + P \mathbf{a} = \mathbf{x} + P_+ \mathbf{a}_+ + P_- \mathbf{a}_-$. Crucially, \mathbf{x}' is the agent’s private information, whereas \mathbf{z} is observable by the principal. This model captures the fact that improvement actions can improve an agent’s underlying attribute as well as observable feature, while gaming actions only affect the outward feature without changing its underlying attribute. The projection matrix P , the available action dimensions, and the quality coefficients θ are assumed to be public information. We discuss cases when these parameters are initially unknown in the appendix.

An agent with pre- (resp. post-) response attribute \mathbf{x} (resp. \mathbf{x}') has a pre- (resp. post-) response *true label* y (resp. y') which indicates the quality of an agent. $y, y' \in \{0, 1\}$, and we use a similar setting as in Hu et al. (2019): $Pr(y = 1) = l(\theta^T \mathbf{x})$, $Pr(y' = 1) = l(\theta^T \mathbf{x}')$, $y' \geq y$, where we can interpret $l : \mathbb{R} \mapsto [0, 1]$ as a likelihood function which is weakly increasing (l is a step-function in

Hu et al. (2019)). We assume that $y' \geq y$ holds for every agent, with improvement actions weakly improving the agent’s true label, and gaming actions leaving it unchanged.

Strategic Learning Problems. We will consider two different strategic learning systems/game settings: (1) the *conventional strategic (CS)* problem where the agents and the principal play the standard Stackelberg game without any added incentive mechanism, both being fully strategic, and (2) the *augmented strategic (AS)* problem, where the agents and the principal play the Stackelberg game with a subsidy mechanism.

Utilities and Optimal Strategies in CS Learning. In a CS learning problem, it is assumed that an agent has the following utility function $u_C(\mathbf{x}, \mathbf{a}) = f(\mathbf{x} + P\mathbf{a}) - h(\mathbf{a})$, where the agent benefits from the decision outcome $f(\mathbf{z})$ and incurs a cost of $h(\mathbf{a}) := \mathbf{c}^T \mathbf{a}$. Denote by $\mathbf{a}_C^*(\mathbf{x}) := \arg \max_{\mathbf{a}} u_C(\mathbf{x}, \mathbf{a})$ the agent’s *CS best response*, with ties broken in favor of its qualification status $\theta^T \mathbf{x}'$. In the same problem, denote y'_C as the *CS post-response label*. The principal’s utility is the classification accuracy $U_C(f) = \int_{\mathcal{X}} Pr(f(\mathbf{x} + P\mathbf{a}_C^*(\mathbf{x})) = y'_C) p(\mathbf{x}) d\mathbf{x}$.

Principal as Mechanism Designer in Augmented Strategic Learning. We focus on *discount mechanisms* that are based on providing a *discount on actions*, where the mechanism provider has the ability to lower the cost of agents’ actions, e.g., making the cost of getting tutoring or exam preparation cheaper during the school admission process. We use G to denote the discount mechanism where the designer chooses a *cost rate discount* value on each action dimension $\Delta \mathbf{c} = (\Delta c_i)_{i=1}^M$, $\Delta c_i < c_i$, and set a *discount amount range* $[\underline{c}, \bar{c}]$. Then with G , the agent’s utility function in the augmented strategic learning becomes $u_A(\mathbf{x}, \mathbf{a}) = f(\mathbf{x} + P\mathbf{a}) - h_A(\mathbf{a})$, where $h_A(\mathbf{a}) = h(\mathbf{a}) - \Delta \mathbf{c}^T \mathbf{a} \cdot \mathbf{1}\{\Delta \mathbf{c}^T \mathbf{a} \in [\underline{c}, \bar{c}]\}$. With G , $\mathbf{a}_A^*(\mathbf{x}) := \arg \max_{\mathbf{a}} u_A(\mathbf{x}, \mathbf{a})$ denotes the agent’s *augmented best response* or *AS best response*, with ties broken in favor of maximizing $\theta^T \mathbf{x}'$ unless otherwise suggested by the mechanism designer. The designer incurs a *subsidy cost* $H(G) = \int_{\mathcal{X}} \Delta \mathbf{c}^T \mathbf{a}_A^*(\mathbf{x}) \cdot \mathbf{1}\{\Delta \mathbf{c}^T \mathbf{a}_A^*(\mathbf{x}) \in [\underline{c}, \bar{c}]\} p(\mathbf{x}) d\mathbf{x}$. Denote by y'_A the *AS post-response label*. The AS utility of the principal is then: $U_A(f) = \int_{\mathcal{X}} Pr(f(\mathbf{x} + P\mathbf{a}_A^*(\mathbf{x})) = y'_A) p(\mathbf{x}) d\mathbf{x} - H(G)$.

Third Party Mechanism in Augmented Strategic Learning. The subsidy cost is the same as above, and the third party charges the principal a tax $\mathcal{T}(G)$ for improved decision performance. The principal’s AS utility becomes $U_A(f) = \int_{\mathcal{X}} Pr(f(\mathbf{x} + P\mathbf{a}_A^*(\mathbf{x})) = y'_A) p(\mathbf{x}) d\mathbf{x} - \mathcal{T}(G)$. Figure 2 illustrates the three party AS classification problem. In designing G , we will consider three commonly studied properties in the mechanism design literature: **Individual rationality (IR)**: The participants are better off in the mechanism than opting out. **Incentive compatibility (IC)**: The participants act in self-interest. **Budget balance (BB)**: $\mathcal{T}(G) \geq H(G)$. We will cover the third party objective in Section 4.

3 OPTIMAL MECHANISM IN AUGMENTED STRATEGIC CLASSIFICATION

In this section, we consider agents from a single demographic group, with the principal as the mechanism designer. Please see the appendix for pictorial interpretations of our results. When designing an incentive mechanism, we consider the classifier f as given. Before analyzing the optimal mechanism design, we define the *subsidy surplus*, $S(f, G) := \int_{\mathcal{X}} [Pr(f(\mathbf{x} + P\mathbf{a}_A^*(\mathbf{x})) = y'_A) - Pr(f(\mathbf{x} + P\mathbf{a}_C^*(\mathbf{x})) = y'_C)] p(\mathbf{x}) d\mathbf{x} - H(G)$. The integral part in $S(f, G)$ is the *benefit gain* of the decision maker and the value in the square bracket is the *individual subsidy benefit*. The decision maker’s problem is equivalent to maximizing $S(f, G)$ under IC and IR.

Theorem 3.1. For general $f(\mathbf{z}) = \mathbf{1}\{\mathbf{w}^T \mathbf{z} \geq \tau\}$, p , and l functions, finding the optimal IC and IR discount mechanism requires solving non-convex optimization problems and thus is NP-hard.

While finding the optimal mechanism under IC and IR constraints is NP-hard, we can develop polynomial time solutions that find IC, IR and $S(f, G) > 0$ when l is convex (see Algorithm E in appendix). Moreover, we have the following result for an important special case $\mathbf{w} = \theta$.

Theorem 3.2. If $\mathbf{w} = \theta$ in f , f incentivizes gaming, and l is convex on $[0, \tau]$, then we have a constant time algorithm that finds a G that is IC, IR, and $S(f, G) > 0$, which is also the optimal mechanism under certain conditions (please see appendix for algorithm and condition).

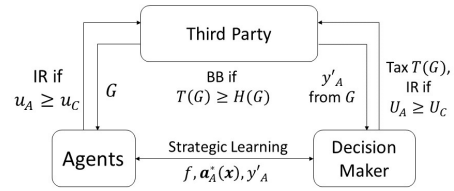


Figure 2: Three-party AS learning.

The convexity requirement of l on a low range is satisfied in real-world datasets such as the FICO credit score dataset, in which the likelihood function l frequently has an S-shape. $\mathbf{w} = \boldsymbol{\theta}$ holds in the optimal CS classifier when it is impossible to incentivize improvement with f , as shown below.

Theorem 3.3. *Let κ_i denote the substitutability of action dimension i (Kleinberg & Raghavan (2020); Jin et al. (2021)). Formally, $\kappa_i := \min_{\mathbf{a} \in \mathbb{R}^M, \mathbf{a} \geq 0} \frac{\mathbf{c}^T \mathbf{a}}{c_i}$, s.t. $P\mathbf{a} - \mathbf{p}_i \geq 0$, where \mathbf{p}_i is the i -th column of P . If $\kappa_i < 1, \forall i \leq M_+$, then is no f that can incentivize improvement ($\mathbf{a}_C^*(\mathbf{x})$ is always gaming), and the decision maker’s CS optimal strategy f_C^* satisfies $\mathbf{w} = \boldsymbol{\theta}$.*

4 DEMOGRAPHIC GROUPS, SOCIAL WELL-BEING

Consider now the case where agents come from two demographic groups distinguished by a *sensitive attribute* $d \in \{1, 2\}$ (e.g., gender, race), which is not a part of the N skill-related attributes (not in \mathbf{x}) and is never influenced by an agent’s action \mathbf{a} . Suppose the decision rule is *not allowed* to use the sensitive attribute as input but that it can be used to design *group specific* subsidies, where the agents voluntarily reveal their sensitive attributes.

Without loss of generality, we will refer to group 1 as the *advantaged* group and 2 as the *disadvantaged* group, and consider group 2 to have the same distribution on \mathbf{x} but disadvantaged in cost², i.e., $h^{(2)}(\mathbf{a}) > h^{(1)}(\mathbf{a}), \forall \mathbf{a} \neq \mathbf{0}$, where $h^{(d)}$ denotes the cost function for group d .

Social Well-being Metrics. We will use the equilibrium qualification status $\mathbb{E}[y'_t], t \in \{C, A\}$ as an *efficiency* oriented social well-being metric. We also introduce *fairness* oriented well-being metrics.

Quality gain measures the increase in agents’ expected qualification status (positive rate) in the response phase $\Delta Q_t^d := \mathbb{E}[Y'_t | D = d] - \mathbb{E}[Y | D = d]; \forall d \in \{1, 2\}, \forall t \in \{A, C\}$. The *quality gain gap* $\gamma_t^Q(f, G) := \Delta Q_t^{(1)} - \Delta Q_t^{(2)}$ measures the group-wise improvement difference.

We also consider two commonly used fairness criteria in classification, Equal Opportunity (EO) (equalized true positive rates) Hardt et al. (2016b) and Demographic Parity (DP) (equalized positive decision rates), and define their respective group differences:

$$\begin{aligned} \gamma_t^{EO}(f, G) &:= Pr(f(\mathbf{z}_t) = 1 | Y'_t = 1, D = 1) - Pr(f(\mathbf{z}_t) = 1 | Y'_t = 1, D = 2), t \in \{A, C\}; \\ \gamma_t^{DP}(f, G) &:= Pr(f(\mathbf{z}_t) = 1 | D = 1) - Pr(f(\mathbf{z}_t) = 1 | D = 2). \end{aligned}$$

An efficiency oriented third party maximizes the *efficiency* oriented metric and a fairness oriented third party minimizes a linear combination of the three fairness gaps above.

Theorem 4.1. *Denote AS-fair (resp. AS-eff), and AS-dm as the equilibrium with a fairness (resp. efficiency) oriented third party, and the principal. 1. The DP gap in weakly ascending order is: AS-fair, CS, AS-dm, AS-eff; 2. The EO gap (or quality gain gap) in weakly ascending order is: AS-fair, CS, AS-dm, AS-eff; 3. The equilibrium social qualification status in weakly descending order is: AS-eff, AS-dm, CS. (Please see the appendix for conditions that make the ranking strict)*

Figure 3 illustrates our numerical results on the FICO dataset Hardt et al. (2016b).³ It shows that while AS-eff reaches the highest efficiency metric, the AS-fair is the best in fairness. Moreover, we note that AS-fair achieves a higher efficiency metric than CS and AS-dm in this experiment.

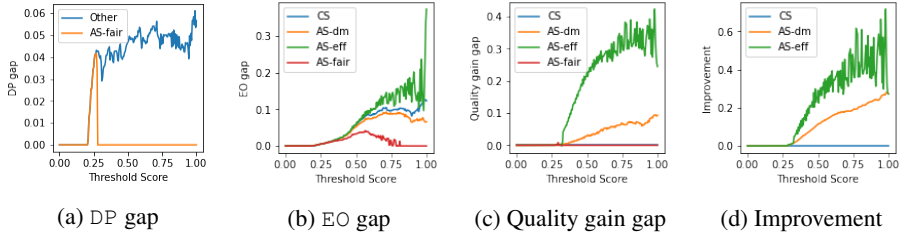


Figure 3: AS Equilibrium Outcomes in FICO with Cost Disadvantages

²Please see the appendix for another case where group 2 is disadvantaged in attribute distribution.

³Please see appendix for experiment settings.

REFERENCES

- Mark Braverman and Sumegha Garg. The role of randomness and noise in strategic classification. In *1st Symposium on Foundations of Responsible Computing*, 2020.
- Michael Brückner and Tobias Scheffer. Stackelberg games for adversarial prediction problems. pp. 547–555, 08 2011. doi: 10.1145/2020408.2020495.
- Michael Brückner, Christian Kanzow, and Tobias Scheffer. Static prediction games for adversarial learning problems. *The Journal of Machine Learning Research*, 13:2617–2654, 09 2012.
- Yatong Chen, Jialu Wang, and Yang Liu. Strategic recourse in linear classification. *arXiv preprint arXiv:2011.00355*, 2020.
- 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*, pp. 55–70, 2018.
- Moritz Hardt, Nimrod Megiddo, Christos Papadimitriou, and Mary Wootters. Strategic classification. pp. 111–122, 01 2016a. doi: 10.1145/2840728.2840730.
- Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. In *Advances in neural information processing systems*, pp. 3315–3323, 2016b.
- Keegan Harris, Hoda Heidari, and Zhiwei Steven Wu. Stateful strategic regression, 2021.
- Lily Hu, Nicole Immorlica, and Jennifer Vaughan. The disparate effects of strategic manipulation. pp. 259–268, 01 2019. doi: 10.1145/3287560.3287597.
- Kun Jin, Tongxin Yin, Charles A. Kamhoua, and Mingyan Liu. Network games with strategic machine learning. In Branislav Bošanský, Cleotilde Gonzalez, Stefan Rass, and Arunesh Sinha (eds.), *Decision and Game Theory for Security*, pp. 118–137, Cham, 2021. Springer International Publishing. ISBN 978-3-030-90370-1.
- Jon Kleinberg and Manish Raghavan. How do classifiers induce agents to invest effort strategically? *ACM Transactions on Economics and Computation*, 8:1–23, 11 2020. doi: 10.1145/3417742.
- John Miller, Smitha Milli, and Moritz Hardt. Strategic classification is causal modeling in disguise. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 6917–6926. PMLR, 13–18 Jul 2020.
- Smitha Milli, John Miller, Anca Dragan, and Moritz Hardt. The social cost of strategic classification. pp. 230–239, 01 2019. doi: 10.1145/3287560.3287576.
- US Federal Reserve. Report to the congress on credit scoring and its effects on the availability and affordability of credit. *Board of Governors of the Federal Reserve System*, 2007.
- Yonadav Shavit, Benjamin Edelman, and Brian Axelrod. Causal strategic linear regression, 2020.

A SUPPLEMENTARY MATERIALS FOR SECTION 2

Remark 1. *The projection matrix P , the available action dimensions, and the quality coefficients θ are assumed to be public information for the remainder of the paper. We discuss in the appendix when these parameters are initially unknown for the principal. Parameter learning requires multi-round online learning Shavit et al. (2020); Harris et al. (2021), which is different from the model settings in this paper. However, we show that our incentive mechanisms help with parameter learning in the multi-round online strategic learning models.*

A.1 DISCUSSION ON REMARK 1

In this part, we discuss the case where game parameters like θ and P are unknown and the principal need to be learn them. Unlike the single round, two stage game in the main article, the learning process requires online learning with multiple rounds, each containing two stages.

We note that the quality coefficients θ can be learned in one round by setting $f = 0$ and then we have $(z, y') = (x, y)$, and running any suitable learning algorithm can get an estimate of θ .

However, P can not always be learned in the conventional learning problem. We can use an example can from the impossibility conditions in Theorem 3.3, given those conditions, only the columns whose index has substitutability 1 can be learned, the other columns are always unknown. Below we show how the discount mechanism help with learning the the projection matrix P .

In the regression problem with L1 cost, we can use the following procedures to learn the projection matrices,

- Choose f such that $w > 0$ (without loss of generality, assume that $w > 0 \Rightarrow P^T w > 0$, otherwise some action dimensions are meaningless)
- For each time step $t = 1, \dots, M$, get a sufficiently large sample of agents with their observable features z
- At $t = 0$, $G_d = 0$, let $\bar{z}_0 = \mathbb{E}[z]$
- At $t = 1, \dots, M$, let G_d induce the best response along action dimension t by lowering the cost to \tilde{c}_t , and let $\bar{z}_t = \mathbb{E}[z]$
- Compute $v_t = (\bar{z}_t - \bar{z}_0)\tilde{c}_t/B$, which is an estimate of $Pe_t = p_t$, i.e., the t -th column of P .

Discount mechanisms can enable best responses to in action dimensions that are impossible to be incentivized with the decision rule itself, and this is true for both classification and regression, both L1 cost and other types of costs like L2 or squared.

A.2 THE LIMITED STRATEGIC (LS) LEARNING PROBLEM

The *limited strategic (LS)* problem where the agents are fully strategic but the principal does not anticipate the agents' strategic behavior and applies the optimal non-strategic decision rule, e.g., $f(z) = \theta^T z$ in regression as a sub-optimal option.⁴

Lemma A.1. *The LS optimal decision rule is $f_L^*(z) = \mathbf{1}\{\theta^T z \geq \tau_L\}$, $\tau_L = \arg \min_{\tau} l(\tau) \geq 0.5$.*

Proof. This is because it is optimal for the principal to accept every agent with $l(\theta^T x) \geq 0.5$, since rejecting this agent results in a decrease in the expected individual prediction outcome $1 - l(\theta^T x) \leq l(\theta^T x)$. Similarly, the principal wants to reject every agent with $l(\theta^T x) < 0.5$. \square

A.3 OTHER COST FUNCTIONS

⁴The agents perform the same as in CS problems. The LS problem is reasonable since the CS problem is in general NP-hard for the principal Kleinberg & Raghavan (2020).

We will use the L2 cost $h(\mathbf{a}) = \|\mathbf{a}\|_2$ for demonstration purpose, and we note that higher orders of cost functions $h(\mathbf{a}) = \frac{1}{2}\|\mathbf{a}\|_2^2$ are very similar in classification but different in regression. In regression, higher order costs are convex and the marginal cost grows, and thus there is no need to be a budget constraint $B \geq h(\mathbf{a})$, other than that, $h(\mathbf{a}) = \|\mathbf{a}\|_2$ is very representative.

For all other cost functions, we can equivalently have a set of “equal cost contour” i.e., $\{\mathbf{a} | h(\mathbf{a}) = C\}$ for some constant C is a contour. Most cost functions used in economic and computer science literature have contours with different sizes but a constant “shape” (the surface of norm balls, since the cost functions are norm based), like the L1 cost, L2 cost, tilted L2 cost $h(\mathbf{a}) = \sqrt{\mathbf{a}^T C \mathbf{a}}$ and squared cost $h(\mathbf{a}) = \frac{1}{2}\|\mathbf{a}\|_2^2$. The constant shape of contours made it possible to have a concise (closed-form in most cases) representation of the best responses’ directional and magnitude properties.

For example, when $h(\mathbf{a}) = \|\mathbf{a}\|_2$, the best responses satisfy $\rho(\mathbf{a}_t^*, P^T \mathbf{w}) = 1$ where $\rho(\mathbf{v}_1, \mathbf{v}_2) = \frac{\mathbf{v}_1^T \mathbf{v}_2}{\|\mathbf{v}_1\|_2 \|\mathbf{v}_2\|_2}$ is the cosine similarity. We still have properties in Lemma B.1 and in classification and regression, the best responses are

$$\mathbf{a}_C^*(\mathbf{x}) = \frac{\tau - \mathbf{w}^T \mathbf{x}}{\|P^T \mathbf{w}\|_2} P^T \mathbf{w}, \quad \mathbf{a}_C^*(\mathbf{x}) = \frac{B}{\|P^T \mathbf{w}\|_2} P^T \mathbf{w},$$

and we can similarly write out the expressions of the AS best responses for other cost functions.

For L2 cost $h(\mathbf{a}) = \|\mathbf{a}\|_2$, we can think of discounts with minimum effective discount value as giving certain action directions a fixed discount rate or incentivizing agents to play a different action and pay the cost differences.

Therefore, the implementer will try to incentivize some of the agents to take an AS best response that also reach the boundary, this can be done by making the discount amount equal the cost difference between the AS and the CS/LS best response. The implementer wants to maximize the subsidy surplus on a given agent, which is the quality gain $l(\mathbf{a}) - l(\mathbf{a}_C^*(\mathbf{x}))$ minus the subsidy cost $\|\mathbf{a}\|_2 - \|\mathbf{a}_C^*(\mathbf{x})\|_2$ and thus $\mathbf{a}_A^*(\mathbf{x})$ is the solution to the optimization problem

$$\begin{aligned} & \text{minimize}_{\mathbf{a}} \|\mathbf{a}\|_2 - l(\boldsymbol{\theta}^T (\mathbf{x} + \hat{P} \mathbf{a})) \\ & \text{subject to } \mathbf{w}^T P \mathbf{a} = \tau - \mathbf{w}^T \mathbf{x} \end{aligned} \quad (1)$$

However, the above problem is in general not convex and can be NP hard to find the optimal solution. But the below assumption guarantees a solution.

Assumption 1. $\mathbf{w} = \boldsymbol{\theta}$, and the implementer limit the AS best response to be gaming free, i.e., $[\mathbf{a}_A^*(\mathbf{x})]_j = \mathbf{1}\{j \leq M_i\} \Leftrightarrow P \mathbf{a}_A^*(\mathbf{x}) = \hat{P} \mathbf{a}_A^*(\mathbf{x})$,

Under Assumption 1, the problem becomes convex since $l(\boldsymbol{\theta}^T (\mathbf{x} + \hat{P} \mathbf{a})) = l(\tau)$ is constant

$$\begin{aligned} & \text{minimize}_{\mathbf{a}} \|\mathbf{a}\|_2 \\ & \text{subject to } \boldsymbol{\theta}^T \hat{P} \mathbf{a} = \tau - \boldsymbol{\theta}^T \mathbf{x} \end{aligned} \quad (2)$$

and the solution (AS best response to incentivize) is

$$\mathbf{a}_A^*(\mathbf{x}) = (\tau - \boldsymbol{\theta}^T \mathbf{x}) \frac{\hat{P}^T \boldsymbol{\theta}}{\|\hat{P}^T \boldsymbol{\theta}\|_2} \quad (3)$$

We can then similarly define the individual subsidy surplus in the L2 case and find sufficient conditions that guarantees an IC, IR and BB mechanism $G \neq 0$ or even find the optimal solutions with the same assumptions made in the Theorems B.4 and 3.2.

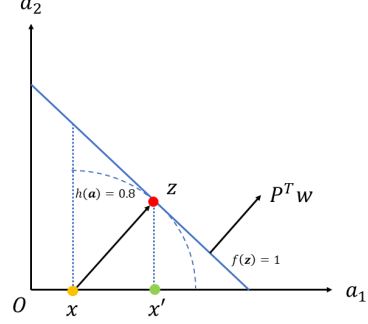


Figure 4: An illustration of a CS best response in classification with L2 cost, where the blue dashed curve (quarter circle) is an equal cost contour, $P = [1, 1]$, $\mathbf{w} = (1, 1)$, a_1 is improvement and a_2 is gaming.

One interesting difference in the L2 cost case is that the decision rule can incentivize *partial improvement*, which can also be called *partial gaming*, which means $\theta^T P\mathbf{a} > \theta^T \hat{P}\mathbf{a} > 0$, and the corresponding theorems in L1 case still applies when f incentivizes pure gaming $\theta^T P\mathbf{a} > \theta^T \hat{P}\mathbf{a} = 0$. An example of pure gaming happens when for every improvement action j , there is a corresponding gaming action k with the an exaggerated effect $\mathbf{p}_k = \alpha_j \mathbf{p}_j$, $\alpha_j > 1$, which can model problems like multi-subject exams where an agent has an improvement and gaming action for each of the subject and cheating is always more cost efficient than working hard without an incentivize mechanism.

A.4 AN ALTERNATIVE INCENTIVE MECHANISM

An alternative mechanism to consider, the *transfer mechanisms* is based on *monetary transfer*, where the mechanism designer provides reimbursement or bonus payment when the agent meets certain feature criteria, e.g., rewards for high scores. We use G_t to denote the transfer mechanism, where the designer chooses a bonus amount $b(\mathbf{z})$, $b: \mathbb{R}^N \mapsto \mathbb{R}$, effectively revising the agent’s utility to

$$u_A(\mathbf{x}, \mathbf{a}) = f(\mathbf{x} + P\mathbf{a}) - h(\mathbf{a}) + b(\mathbf{x} + P\mathbf{a}). \quad (4)$$

In transfer mechanisms, knowing the actual \mathbf{x} seems to help the designer reduce the subsidy cost on agents with high endowment and low improvement, but we will show below that this extended version with bonus amount $\tilde{b}(\tilde{\mathbf{x}}, \mathbf{z})$ is equivalent as the bonus $b(\mathbf{z})$ that only uses features as input, where $\tilde{\mathbf{x}}$ is the reported pre-response attribute. This is because $\tilde{b}(\tilde{\mathbf{x}}, \mathbf{z})$ either can not incentivize agents to truthfully report $\tilde{\mathbf{x}} = \mathbf{x}$, or it can not generate more benefit for the mechanism designer.

With the alternative version of the monetary transfer mechanism, the agent’s utility now becomes

$$\tilde{u}_A(\mathbf{x}, \mathbf{a}, G_t) = f(\mathbf{x} + P\mathbf{a}) - h(\mathbf{a}) + \max_{\tilde{\mathbf{x}}} \tilde{b}(\tilde{\mathbf{x}}, \mathbf{x} + P\mathbf{a}),$$

and we can find the corresponding $\mathbf{a}_A^*(\mathbf{x})$, and only if

$$\mathbf{x} \in \arg \max_{\tilde{\mathbf{x}}} \tilde{b}(\tilde{\mathbf{x}}, \mathbf{x} + P\mathbf{a}_A^*(\mathbf{x})),$$

truth-reporting is incentivized. If truth reporting is not incentivized, $\tilde{b}(\tilde{\mathbf{x}}, \mathbf{z})$ and $b(\mathbf{z}) = \max_{\tilde{\mathbf{x}}} \tilde{b}(\tilde{\mathbf{x}}, \mathbf{z})$ are equivalent for both the agents and the mechanism designer. Meanwhile, for $\forall \mathbf{x}_1 \neq \mathbf{x}_2$, truth telling requires either

$$\mathbf{x}_1 + P\mathbf{a}_A^*(\mathbf{x}_1) \neq \mathbf{x}_2 + P\mathbf{a}_A^*(\mathbf{x}_2),$$

indicating that backward induction from $\mathbf{x} + P\mathbf{a}_A^*(\mathbf{x})$ to \mathbf{x} is achievable, or

$$\mathbf{x}_1 + P\mathbf{a}_A^*(\mathbf{x}_1) = \mathbf{x}_2 + P\mathbf{a}_A^*(\mathbf{x}_2), \text{ and } \mathbf{x}_1, \mathbf{x}_2 \in \arg \max_{\tilde{\mathbf{x}}} \tilde{b}(\tilde{\mathbf{x}}, \mathbf{x} + P\mathbf{a}_A^*(\mathbf{x})).$$

In either case, $b(\mathbf{z})$ is sufficient.

However, the computational complexity is very high in the backward induction step for a general $b(\mathbf{z})$ bonus function. Recall that the AS utility of an agent is

$$u_A(\mathbf{x}, \mathbf{a}) = f(\mathbf{x} + P\mathbf{a}) - h(\mathbf{a}) + b(\mathbf{x} + P\mathbf{a}),$$

and thus computing $\mathbf{a}_A^*(\mathbf{x}) = \arg \max_{\mathbf{a}} u_A(\mathbf{x}, \mathbf{a})$ is non-convex for a non-concave $b(r)$ bonus function.

On one hand, we can’t guarantee concave $b(r)$ is the optimal solution. On the other hand, for a concave $b(\mathbf{z})$, the computation of $\mathbf{a}_A^*(\mathbf{x}) = \arg \max_{\mathbf{a}} u_A(\mathbf{x}, \mathbf{a})$ is convex and but the individual subsidy surplus

$$s(\mathbf{x}, f, G_t) = l(\theta^T(\mathbf{x} + \hat{P}\mathbf{a}_A^*(\mathbf{x}))) - l(\theta^T(\mathbf{x} + \hat{P}\mathbf{a}_C^*(\mathbf{x}))) - b(\mathbf{x} + P\mathbf{a}_A^*(\mathbf{x}))$$

on the agents are not concave unless l is convex (we are supposing $\mathbf{x} \in \mathcal{M}(f)$ here, otherwise more non-convexity is introduced). Moreover, the overall objective depends on the integration on a subset of $\hat{\mathcal{X}} \subseteq \mathcal{X}$

$$S(f, G_t) = \int_{\hat{\mathcal{X}}} s(\mathbf{x}, f, G_t) p(\mathbf{x}) d\mathbf{x},$$

and a general probability density function p , and the convexity of set $\hat{\mathcal{X}}$ can make the mechanism designer’s objective non-convex even if l is convex.

We also note that when changing the value $b(\mathbf{z})$ for a certain \mathbf{z} , the AS best response for all agents with pre-response attribute \mathbf{x} in the cone $\mathbf{x} - \mathbf{z} \leq 0$ (element wise non-positive) might change, and this also makes the analysis hard.

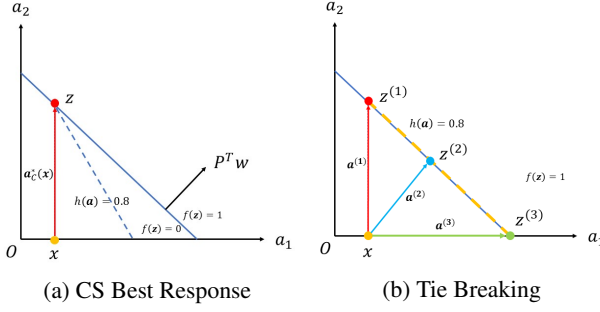


Figure 5: An illustration of a CS best response in classification, where $P = [1, 1]$, $w = 1$, $P^T w = (1, 1)$. Solid blue is the decision boundary. In (a) blue dashed line is an equal cost contour; $c_2 < c_1$, thus gaming is cheaper than improving leading to best response shown in red. (b) illustrates tie breaking in best responses, where $c_1 = c_2$, with equal cost contour shown in yellow dashed line.

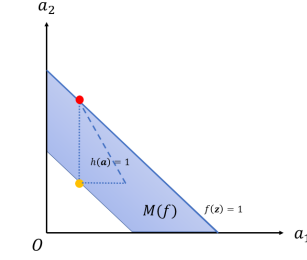


Figure 6: An illustration of the manipulation margin in classification, given by the shaded region; every agent inside can reach the boundary with an action cost no more than 1.

B SUPPLEMENTARY MATERIALS FOR SECTION 3

This part considers agents from a single demographic group. Throughout our analysis, we will provide pictorial interpretations of our results, using an example with 2 action dimensions: a_1 is an improvement action and a_2 a gaming action.

We begin with some preliminaries. The next two Lemmas characterize the magnitude and direction of the agents' best responses $\mathbf{a}_t^*(\mathbf{x})$ ($t \in \{C, A\}$) in the conventional and augmented strategic games.

Lemma B.1. *For CS and AS classification, $\mathbf{w}^T(\mathbf{x} + P\mathbf{a}_t^*(\mathbf{x})) = \tau \Leftrightarrow \mathbf{a}_t^*(\mathbf{x}) \neq \mathbf{0}, \forall t$.*

Proof. For $\forall \mathbf{a}$ such that $\mathbf{w}^T(\mathbf{x} + P\mathbf{a}) < \tau$, $f(\mathbf{z}) = 0$ and thus it is dominated by $\mathbf{0}$ due to $h(\mathbf{a}) \geq h(\mathbf{0}) = 0$ and $h_A(\mathbf{a}) \geq h_A(\mathbf{0}) = 0$. On the other hand, for $\forall \mathbf{a}$ such that $\mathbf{w}^T(\mathbf{x} + P\mathbf{a}) > \tau$, there exists an $\alpha \in (0, 1)$ such that $\mathbf{w}^T(\mathbf{x} + P\alpha\mathbf{a}) = \tau$. Both \mathbf{a} and $\alpha\mathbf{a}$ guarantees $f(\mathbf{z}) = 1$, and thus \mathbf{a} is dominated by $\alpha\mathbf{a}$ due to $h(\mathbf{a}) > h(\alpha\mathbf{a})$ and $h_A(\mathbf{a}) > h_A(\alpha\mathbf{a})$ if $\mathbf{a} \neq \mathbf{0}$. \square

Lemma B.1 describes the magnitude of the best response in CS and AS classification: it is such that the feature \mathbf{z} reaches the decision boundary but not beyond, as going beyond the boundary only increases the cost without affecting the decision. This is illustrated by the red arrow in Figure 5a.

Lemma B.2. *For CS and AS classification,*

$$\begin{aligned} (\mathbf{a}_C^*(\mathbf{x}))_i &\geq 0, \text{ if } i \in \{\arg \max_j (P^T \mathbf{w})_j / c_j\}; \quad (\mathbf{a}_t^*(\mathbf{x}))_i = 0, \text{ o.w., } \forall \mathbf{x}. \\ (\mathbf{a}_A^*(\mathbf{x}))_i &\geq 0, \text{ if } i \in \{\arg \max_j (P^T \mathbf{w})_j / (c_j - \Delta c_j)\}; \quad (\mathbf{a}_A^*(\mathbf{x}))_i = 0, \text{ o.w., } \forall \mathbf{x}. \end{aligned} \quad (5)$$

Proof. Assume by contradiction $a_j^* > 0, j \neq i_C = \arg \max_k \frac{(P^T \mathbf{w})_k}{c_k}$. By Lemma B.1, as $\mathbf{a}^* \neq \mathbf{0}$ we have $\mathbf{w}^T(\mathbf{x} + P\mathbf{a}^*) = \tau$. Denote $\tilde{\mathbf{a}} = \mathbf{a}^* - a_j^* \mathbf{e}_j + \frac{a_j^* (P^T \mathbf{w})_j}{(P^T \mathbf{w})_{i_C}} \mathbf{e}_{i_C}$, where \mathbf{e}_i is the i -th orthonormal base vector of \mathbb{R}^M . It is easy to see that $\mathbf{w}^T(\mathbf{x} + P\tilde{\mathbf{a}}) = \tau$ and thus $f(\mathbf{z}) = 1$, while $h(\tilde{\mathbf{a}}) < h(\mathbf{a}^*)$, indicating that $\tilde{\mathbf{a}}$ achieves a higher utility than \mathbf{a}^* , contradicting the optimality of \mathbf{a}^* . The proof for AS classification is similar. \square

Lemma B.2 describes the directional properties of the best response: the agent will invest in the action dimension(s) with the highest *return on investment* $(P^T \mathbf{w})_j / c_j$ (in CS) or $(P^T \mathbf{w})_j / (c_j - \Delta c_j)$ (in AS). Without loss of generality, we assume that the optimal CS action dimension $i_C := \arg \max_j (P^T \mathbf{w})_j / c_j$ is unique. This property is also shown in Figure 5a, where $i_C = 2$ is the action with the highest return on investment.

We note that there may be multiple actions that are tied in their return on investment. In such cases, we assume the agent follows the algorithm designer's recommendation if any, and otherwise

chooses the one that leads to the maximum improvement (i.e., the one maximizing $\theta^T \hat{P}\mathbf{a}$). Figure 5b explains this tie breaking: here $c_1 = c_2$ and every point on the yellow contour has equal cost and benefit to the agent, making the agent indifferent between $\mathbf{a}^{(1)}, \mathbf{a}^{(2)}, \mathbf{a}^{(3)}$. We take $\mathbf{a}^{(3)}$, the largest improvement, to be the agent's choice.

Using Lemma B.1 and B.2, we have

$$\mathbf{a}_C^*(\mathbf{x}) = \frac{\tau - \mathbf{w}^T \mathbf{x}}{(P^T \mathbf{w})_{i_C}} \mathbf{e}_{i_C}, \text{ if } \mathbf{x} \in \mathcal{M}(f); \quad \mathbf{a}_C^*(\mathbf{x}) = \mathbf{0}, \text{ o.w.}, \quad (6)$$

where $\mathcal{M}(f) := \left\{ \mathbf{x} \mid \frac{(\tau - \mathbf{w}^T \mathbf{x})_{i_C}}{(P^T \mathbf{w})_{i_C}} \in (0, 1] \right\}$ denotes the *manipulation margin* of f : every agent in the manipulation margin has non-zero best response to improve their decision outcome to 1. This is illustrated in Figure 6.

In classification, if $i_C \leq M_+$, we say the decision rule *incentivizes improvement*, otherwise we say the decision rule *incentivizes gaming*.

Lemma B.3. *To induce an agent to take an AS best response with non-zero investment in action dimension $j \leq M_+$, i.e., $[\mathbf{a}_A^*(\mathbf{x})]_j > 0$, the discount value Δc_j should satisfy $(P^T \mathbf{w})_j / (c_j - \Delta c_j) \geq (P^T \mathbf{w})_{i_C} / (c_{i_C})$, i.e., $\Delta c_j \geq c_j - \frac{(P^T \mathbf{w})_j}{(P^T \mathbf{w})_{i_C}} c_{i_C}$.*

Based on Lemma B.3, we denote the *minimum effective discount value* as

$$\Delta c_j^* := c_j - \frac{(P^T \mathbf{w})_j}{(P^T \mathbf{w})_{i_C}} c_{i_C}. \quad (7)$$

Intuitively, Lemma B.3 states that to induce a best response in action j , the discount has to make j the action with the highest (potentially tied) return on investment. Figure 7a illustrates an example of how the discount mechanism with minimum effective discount value works. By choosing $\Delta c_1 = \Delta c_1^*$, the two actions have the same return on investment; the agents choose $\mathbf{a}_A^*(\mathbf{x})$, the maximum improvement action, in this case. In contrast, the CS action $\mathbf{a}_C^*(\mathbf{x})$ is a gaming action.

Below we present the proof of Theorem 3.3

Proof. The proofs of the claims

1. If $\kappa_j = 1$, then there exists a \mathbf{w} in f that can incentivize action dimension j , and the \mathbf{w} can be found in polynomial time;
2. if $\kappa_j < 1$, meaning there always are linear combinations of gaming actions weakly dominate every action j , then there is no f that can incentivize best response on action j .

are covered in Kleinberg & Raghavan (2020); Jin et al. (2021). Intuitively, if $\kappa_j < 1, \forall j \leq M_+$, the corresponding \mathbf{a} is the combination that strictly dominates \mathbf{e}_j for any f and thus there is no f that can incentivize improvement.

We will proceed to show the principal's CS optimal strategy satisfy $\mathbf{w} = \theta$. The main idea is that when f always incentivizes gaming, then the CS decision outcomes with $f_C(\mathbf{z}) = \mathbf{1}\{\mathbf{w}_C^T \mathbf{z} \geq \tau_C\}$ always have an equivalent LS decision outcomes with $f_L(\mathbf{z}) = \mathbf{1}\{\mathbf{w}_L^T \mathbf{z} \geq \tau_L\}$, where the $\mathbf{w}_C = \mathbf{w}_L$, and τ_C, τ_L satisfy

$$\tau_L = \min \left\{ 0, \tau_C - \frac{(P^T \mathbf{w})_k}{c_k} \right\}.$$

In other words, we can show that $\forall \mathbf{x}, f_L(\mathbf{x}) = f_C(\mathbf{x} + P\mathbf{a}^*)$, and thus is equivalent for the principal to find an optimal f_L which guarantees $\mathbf{w}_L = \theta$ as the Lemma A.1 suggests. \square

Recall that Theorem 3.1 states that finding the optimal mechanism requires solving non-convex optimization problem and thus is NP-hard. Below we present the proof.

Proof. We will first show the problem is non-convex when discount is placed on multiple actions, then show even the discount is only on one action, the problem is still non-convex.

When the discount is on multiple actions, providing the optimal tie breaking strategy for an agent with \mathbf{x} requires solving

$$\text{maximize}_{\mathbf{a}} l(\boldsymbol{\theta}^T(\mathbf{x} + \hat{P}\mathbf{a})) - \Delta\mathbf{c}^T\mathbf{a},$$

which is non-convex for a general l function. This is for individual subsidy surplus for a fixed $\Delta\mathbf{c}$, and it has to be integrated over \mathcal{X} to compute the overall subsidy surplus $S(f, G)$. So finding the optimal mechanism will only have higher computational complexity when the principal has to optimize over $\Delta\mathbf{c}$, \underline{c} , \bar{c} , and take into account the influence of $p(\mathbf{x})$.

When the discount is only on one action, from Lemma B.3, the mechanism designer need to choose $\Delta\mathbf{c}$ such that

$$\Delta\mathbf{c}_j \geq \Delta\mathbf{c}_j^* = c_j - \frac{(P^T\mathbf{w})_j}{(P^T\mathbf{w})_{i_C}} c_{i_C},$$

for some improvement action dimension $j \leq M_+$ that it wants to incentivize the agents.

Then for the principal, maximizing its AS utility is equivalent as maximizing the subsidy surplus, so the principal solves

$$\begin{aligned} & \text{maximize}_{j, \Delta\mathbf{c}_j, \underline{c}, \bar{c}} \int_{\mathcal{X}} [Pr(f(\mathbf{x} + P\mathbf{a}_A^*(\mathbf{x})) = y'_A) - \mathbf{1}\{\Delta\mathbf{c}^T\mathbf{a}_A^*(\mathbf{x}) \in [\underline{c}, \bar{c}]\}] p(\mathbf{x}) d\mathbf{x} \\ & \text{subject to } \Delta\mathbf{c} \in [\Delta\mathbf{c}^*, c_j], j \leq M_+ \end{aligned}$$

where the problem can be non-convex and not monotone for general p and l . Specifically, when j has the highest return of investment after the discount, the backward induction that anticipates the agent's AS best response is,

$$\mathbf{a}_A^*(\mathbf{x}) = \begin{cases} \frac{\tau - \mathbf{w}^T\mathbf{x}}{(P^T\mathbf{w})_j} \mathbf{e}_j, & \text{if } \frac{\Delta c_j (\tau - \mathbf{w}^T\mathbf{x})}{(P^T\mathbf{w})_j} \in [\underline{c}, \bar{c}], \\ \frac{\tau - \mathbf{w}^T\mathbf{x}}{(P^T\mathbf{w})_{i_C}} \mathbf{e}_{i_C}, & \text{o.w.} \end{cases}$$

This indicates that agents with \mathbf{x} in a belt shape subset of \mathcal{X} will be incentivized to improve, but the overall subsidy surplus is in general not convex, not concave and not monotone in either the upper bound (determined by f and \bar{c}) or the lower bound (determined by f and \underline{c}) of the belt even when the other is fixed. Moreover, the minimum effective discount value $\Delta\mathbf{c}_j^*$ is not always the optimal solution, adding more complexity to the problem. This is because sometimes the principal wants to put more discount on the action dimension and incentivize some agents outside of the manipulation margin to improve and accept them rather than reject them. For example, if 80 percent agent has attribute that makes their likelihood 0.49, the minimum effective discount value still makes them rejected and take 0 AS best response, but a slightly higher discount can incentivize them all to improve to the threshold value, say 0.7, the $0.7 - (1 - 0.49) \cdot 0.8 = 0.152$ amount of improvement may largely outweigh the extra subsidy cost.

Overall speaking, the difference between \mathbf{w} and $\boldsymbol{\theta}$ in f , the global properties of p , l and their local properties influenced by τ all makes the problem hard to solve. \square

Theorem B.4. *Algorithm 1 runs in polynomial time, and if l is convex on $[0, \max_{\mathbf{x}: \mathbf{w}^T\mathbf{x}=\tau} l(\mathbf{x})]$, then any $G \neq 0$ returned by Algorithm 1 is IC, IR, and satisfies $S(f, G) \geq 0$.*

Proof. We will show that any $G \neq 0$ returned by Algorithm 1 is IC, IR and satisfies $S(f, G) \geq 0$.

The IC part follows that the participants act in self-interest. Also, as previously discussed, the minimum effective discounted value $\Delta\mathbf{c}_j = \Delta\mathbf{c}_j^* = c_j - \frac{(P^T\mathbf{w})_j}{(P^T\mathbf{w})_{i_C}} c_{i_C}$ makes sure the agents are weakly better off in the AS game than the CS game (given the same f).

We note that for all f that incentivizes gaming, the principal would prefer $\mathbf{w} = \boldsymbol{\theta}$ and we can use Theorem 3.2 to find G , so below we have $i_C \leq M_+$.

The basic logic of ensuring $S(f, G) \geq 0$ is that the algorithm finds a specific agent that is incentivized, and if this specific agent has a non-negative individual subsidy surplus, it is sufficient

that all the other incentivized agents also have non-negative individual subsidy surplus and thus $S(f, G) \geq 0$.

In Algorithm 1, the designer finds (a convex problem and easy to solve)

$$\underline{\mathbf{x}} = \arg \min_{\mathbf{x}: \mathbf{w}^T \mathbf{x} = \tau - \delta_j (P^T \mathbf{w})_j} \boldsymbol{\theta}^T \mathbf{x},$$

which is the attribute of the specific agent. From the upper bound set on δ_j in the algorithm, we assume the specific agent is in $\mathcal{M}(f)$, and then uses

$$\underline{s} = l_+ - \delta_j \Delta c_j = l(\boldsymbol{\theta}^T(\underline{\mathbf{x}} + \delta_j \hat{P} \mathbf{e}_j)) - l(\boldsymbol{\theta}^T(\underline{\mathbf{x}} + \delta_{i_C} \hat{P} \mathbf{e}_{i_C})) - \delta_j \Delta c_j,$$

as a benchmark, where δ_j is the δ in the algorithm and $\delta_{i_C} = \frac{(P^T \mathbf{w})_j}{(P^T \mathbf{w})_{i_C}} \delta_j$. $\delta_j \mathbf{e}_j$ and $\delta_{i_C} \mathbf{e}_{i_C}$ help the agent achieve the same $\mathbf{w}^T \mathbf{z}$, $\underline{c} = 0$, $\bar{c} = \delta_j \Delta c_j$ here.

Then \underline{s} is the specific agent’s individual subsidy surplus, i.e.,

$$s(\underline{\mathbf{x}}, f, G) = l(\boldsymbol{\theta}^T(\underline{\mathbf{x}} + \hat{P} \mathbf{a}_A^*(\underline{\mathbf{x}}))) - l(\boldsymbol{\theta}^T(\underline{\mathbf{x}} + \hat{P} \mathbf{a}_C^*(\underline{\mathbf{x}}))) - \mathbf{1}\{\Delta \mathbf{c}^T \mathbf{a}_A^*(\underline{\mathbf{x}}) \in [\underline{c}, \bar{c}]\} = \underline{s}.$$

We start with agents with CS best response $\mathbf{a}_C^*(\mathbf{x}) = \delta_{i_C} \mathbf{e}_{i_C}$, i.e., $\mathbf{w}^T \mathbf{x} = \mathbf{w}^T \underline{\mathbf{x}}$. For them, the AS best response is $\mathbf{a}_A^*(\mathbf{x}) = \delta_j \mathbf{e}_j$, the individual subsidy surplus is then

$$s(\mathbf{x}, f, G) = l(\boldsymbol{\theta}^T(\mathbf{x} + \delta_j \hat{P} \mathbf{e}_j)) - l(\boldsymbol{\theta}^T(\mathbf{x} + \delta_{i_C} \hat{P} \mathbf{e}_{i_C})) - \delta_j \Delta c_j,$$

since (1) $\boldsymbol{\theta}^T \hat{P}(\delta_j \mathbf{e}_j - \delta_{i_C} \mathbf{e}_{i_C})$ is constant, (2) $\boldsymbol{\theta}^T \mathbf{x} \geq \boldsymbol{\theta}^T \underline{\mathbf{x}}$ and (3) l is convex on this range, we have $s(\mathbf{x}, f, G) \geq \underline{s} \geq 0$.

For agents with “higher endowment” $\mathbf{w}^T \mathbf{x} > \mathbf{w}^T \underline{\mathbf{x}}$, i.e., with CS best response $\mathbf{a}_C^*(\mathbf{x}) = \alpha_{i_C} \mathbf{e}_{i_C}$, $\alpha_{i_C} < \delta_{i_C}$, we denote $\alpha_j = \alpha_{i_C} (P^T \mathbf{w})_{i_C} / (P^T \mathbf{w})_j$, then the (sub-optimal) AS best response is $\mathbf{a}_A^*(\mathbf{x}) = \alpha_j \mathbf{e}_j$, and the individual subsidy surplus is

$$\begin{aligned} s(\mathbf{x}, f, G) &= l(\boldsymbol{\theta}^T(\mathbf{x} + \alpha_j \hat{P} \mathbf{e}_j)) - l(\boldsymbol{\theta}^T(\mathbf{x} + \alpha_{i_C} \hat{P} \mathbf{e}_{i_C})) - \alpha_{i_C} \bar{c} / \delta_{i_C} \\ &\geq \frac{\alpha_{i_C}}{\delta_{i_C}} [l(\boldsymbol{\theta}^T(\mathbf{x} + \delta_j \hat{P} \mathbf{e}_j)) - l(\boldsymbol{\theta}^T(\mathbf{x} + \delta_{i_C} \hat{P} \mathbf{e}_{i_C})) - \bar{c}] \\ &\geq \frac{\alpha_{i_C}}{\delta_{i_C}} \underline{s} \geq 0, \end{aligned}$$

where the second inequality comes from the convexity of l .

For agents with “lower endowment” i.e., with CS best response $\mathbf{a}_C^*(\mathbf{x}) = \beta_{i_C} \mathbf{e}_{i_C}$, $\beta_{i_C} < \delta_{i_C}$, the mechanism designer suggest that they break tie choosing $\mathbf{a}_A^*(\mathbf{x}) = \beta_{i_C} \mathbf{e}_{i_C}$ as the AS best response and thus the individual subsidy surplus is 0. For $\beta_j = \beta_{i_C} (P^T \mathbf{w})_{i_C} / (P^T \mathbf{w})_j$, we note that $\mathbf{a}_A^*(\mathbf{x}) = \beta_j \mathbf{e}_j$ is a dominated strategy since $\Delta \mathbf{c}^T \mathbf{a}_A^*(\mathbf{x}) > \bar{c}$. \square

From Algorithm 1 we see that the principal prefers subsidizing agents that are “closer” to the boundary when l is convex on $[0, \max_{\mathbf{x}: \mathbf{w}^T \mathbf{x} = \tau} l(\mathbf{x})]$. This is because when l is convex, the subsidy benefit becomes concave while the subsidy cost is linear in the “distance to the boundary”; thus the agents close enough to the boundary can have positive individual subsidy surplus; Figure 8 provides an illustration of this.

The convexity requirement of l on a low range is satisfied in real-world datasets such as the FICO credit score dataset, in which the likelihood function l frequently has an S-shape (see Appendix E). We discuss the case of other likelihood function types (including concave) in the appendix.

Also note that in Algorithm 1 the mechanism designer places discount on only one dimension. This is because even though it technically can set the discount $\Delta c_i > 0$ for multiple improvement actions, ultimately the agent either finds the dimension with the highest return on investment or breaks ties in favor of the largest improvement ⁵.

⁵When placing discounts on multiple actions, finding the optimal tie-breaking rule is a non-convex problem.

The optimal mechanism can be found more efficiently for the special case when $\mathbf{w} = \boldsymbol{\theta}$ in f (this happens, e.g., in the optimal LS strategy as shown in Lemma A.1 in the appendix, or in the optimal CS strategy when $\kappa_i < 1, \forall i \leq M_+$ in Theorem 3.3). This can be done in a fixed number of steps (faster than polynomial) using Algorithm 2.

Below is a refined version of Theorem 3.2.

Theorem B.5. *If $\mathbf{w} = \boldsymbol{\theta}$ in f , f incentivizes gaming, and l is convex on $[0, \tau]$, then Algorithm 2 finds a G that is IC, IR, and satisfies $S(f, G) \geq 0$. In addition, algorithm 2 finds the optimal G if $l(\tau) - l(\underline{r}_f) \leq \frac{(\tau - \underline{r}_f)\Delta c_{i_A}^*}{(P^T \boldsymbol{\theta})_{i_A}}$, where $\underline{r}_f = \min_{\mathbf{x} \in \mathcal{M}(f)} \boldsymbol{\theta}^T \mathbf{x}$.*

Proof. When $\mathbf{w} = \boldsymbol{\theta}$, the mechanism designer is indifferent about AS best responses along any improvement action dimension.

The mechanism designer find the ‘‘cheapest to incentivize’’ target action dimension

$$i_A = \arg \max_{j \leq M_+} (P^T \boldsymbol{\theta})_j / c_j \Leftrightarrow i_A = \arg \min_{j \leq M_+} \frac{\Delta c_j^* (\tau - \boldsymbol{\theta}^T \mathbf{x})}{(P^T \boldsymbol{\theta})_j}, \Delta c_{i_A} \leftarrow c_{i_A} - \frac{(P^T \boldsymbol{\theta})_{i_A}}{(P^T \boldsymbol{\theta})_{i_C}} c_{i_C};$$

and set $\Delta \mathbf{c}$ so that $\Delta c_{i_A} \geq \Delta c_{i_A}^* = c_{i_A} - \frac{(P^T \boldsymbol{\theta})_{i_A}}{(P^T \boldsymbol{\theta})_{i_C}} c_{i_C}$.

The choice of \bar{c} depends on the individual subsidy surplus, which is the quality improvement of an incentivized agent minus the subsidy cost, denote $r_{\mathbf{x}} = \boldsymbol{\theta}^T \mathbf{x}$, then ⁶

$$s(r_{\mathbf{x}}, f, G_d) := l(\tau) - l(r_{\mathbf{x}}) \mathbf{1}\{\mathbf{x} \in \mathcal{M}(f)\} - [1 - l(\underline{r}_f)] \mathbf{1}\{\mathbf{x} \notin \mathcal{M}(f)\} - \frac{(\tau - r_{\mathbf{x}})\Delta c_{i_A}^*}{(P^T \boldsymbol{\theta})_{i_A}}, \quad (8)$$

which is because when agents break tie choosing the action with the largest improvement, we have

$$\boldsymbol{\theta}(\mathbf{x} + \hat{P} \mathbf{a}_A^*(\mathbf{x})) = \tau.$$

When the minimum effective discount value is chosen, and the condition

$$l(\tau) - l(\underline{r}_f) \leq \frac{(\tau - \underline{r}_f)\Delta c_{i_A}^*}{(P^T \boldsymbol{\theta})_{i_A}} = (\tau - \underline{r}_f) \left[\frac{c_{i_A}}{(P^T \boldsymbol{\theta})_{i_A}} - \frac{c_{i_C}}{(P^T \boldsymbol{\theta})_{i_C}} \right] \quad (9)$$

holds, all incentivized agents satisfy $\mathbf{x} \in \mathcal{M}(f)$ and $s(r_{\mathbf{x}}, f, G_d) = l(\tau) - l(r_{\mathbf{x}}) - \frac{(\tau - r_{\mathbf{x}})\Delta c_{i_A}^*}{(P^T \boldsymbol{\theta})_{i_A}}$, which is concave in $r, \forall r \leq \tau$ since l is convex on $[0, \tau]$. A rational principal will make sure that an agent with $r_{\mathbf{x}} \geq 0.5$ is in $\mathcal{M}(f)$, and $r_{\mathbf{x}} < 0.5$ is not. And similar to the case in Theorem B.4, agents that fully spends \bar{c} but still need $(\mathbf{a}_A)_{i_C} > 0$ are suggested to stick with their CS best responses.

⁶If $\mathbf{x} \in \mathcal{M}(f)$, incentivizing this agent will result in the same decision outcome and an improvement equilibrium qualification status and thus the subsidy benefit is $l(\tau) - l(r_{\mathbf{x}})$; if $\mathbf{x} \notin \mathcal{M}(f)$, subsidizing this agent will change the decision outcome from 0 to 1, and the subsidy benefit is $l(\tau) - [1 - l(r_{\mathbf{x}})]$. When applying the minimum effective discount value, the agent’s equilibrium action cost is the same in AS and CS outcomes, and thus $\mathbf{x} \in \mathcal{M}(f)$ are incentivized to improve.

ALGORITHM 1: Find a $G \neq 0$ that is IC, IR and $S(f, G) > 0$ for Classification

```

 $\mathbf{x}_1 \leftarrow \arg \min_{\mathbf{x}: \mathbf{w}^T \mathbf{x} = \tau} \boldsymbol{\theta}^T \mathbf{x};$ 
 $i_C \leftarrow \arg \max_j (P^T \mathbf{w})_j / c_j;$ 
for  $j = 1 : M_+$  do
   $\Delta \mathbf{c} \leftarrow \mathbf{0}; \bar{c} \leftarrow 0; l_+ \leftarrow 0;$ 
   $\Delta c_j \leftarrow c_j - \frac{(P^T \mathbf{w})_j}{(P^T \mathbf{w})_{i_C}} c_{i_C};$ 
  Define function
   $\mathbf{a}(\delta) := \delta \mathbf{e}_j - \delta \frac{(P^T \mathbf{w})_j}{(P^T \mathbf{w})_{i_C}} \mathbf{e}_{i_C};$ 
   $l_+(\delta) := l(\boldsymbol{\theta}^T \mathbf{x}_1) - l(\boldsymbol{\theta}^T \mathbf{x}_1 - \mathbf{a}(\delta));$ 
   $\delta^* \leftarrow \arg \max_{\delta} \text{ s.t. } l_+(\delta) \geq \delta \Delta c_j;$ 
  if  $\delta^* = 0$  then
    | Go back to for loop
  end
   $\bar{c} \leftarrow \min\{\delta^*, 1/(c_j - \Delta c_j)\} \cdot \Delta c_j;$ 
  Return  $(\Delta \mathbf{c}, 0, \bar{c})$ 

```

end
Return $(0, 0, 0)$

ALGORITHM 2: A G that is IC, IR and $S(f, G) \geq 0$ for Classification when $\mathbf{w} = \boldsymbol{\theta}$

```

 $i_A \leftarrow \arg \max_{j \leq M_+} (P^T \boldsymbol{\theta})_j / c_j;$ 
 $\Delta c_{i_A} \leftarrow c_{i_A} - \frac{(P^T \boldsymbol{\theta})_{i_A}}{(P^T \boldsymbol{\theta})_{i_C}} c_{i_C};$ 
Define functions
 $s_1(r) := l(\tau) - l(r) - \frac{(\tau - r)\Delta c_{i_A}}{(P^T \boldsymbol{\theta})_{i_A}};$ 
 $s_2(r) := l(\tau) + l(r) - 1 - \frac{(\tau - r)\Delta c_{i_A}}{(P^T \boldsymbol{\theta})_{i_A}};$ 
 $r \leftarrow \arg \min_r \text{ s.t. } s_1(r) \geq 0;$ 
if  $l(r) < 0.5$  then
  |  $r \leftarrow \arg \min_r \text{ s.t. } s_2(r) \geq 0;$ 
end
 $\bar{c} = (\tau - r)\Delta c_{i_A} / (P^T \boldsymbol{\theta})_{i_A};$ 
Return  $(\Delta \mathbf{c}, 0, \bar{c})$ 

```

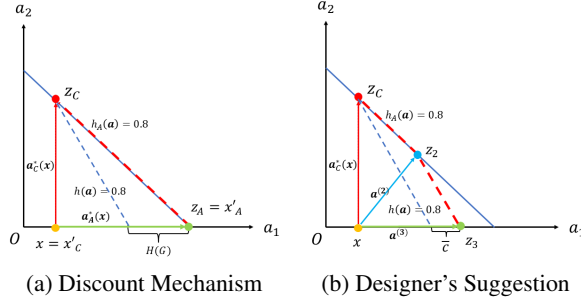


Figure 7: An illustration of the discount mechanism in classification, $P = [1, 1]$, $w = 1$, $P^T w = (1, 1)$, $c_2 < c_1$, the red dashed line is the discounted equal cost contour with a minimum effective discount. In Figure 7b, the \bar{c} is of a smaller value, and the equal cost contour has a different shape. The principal suggests the agents choose $\mathbf{a}_C^*(\mathbf{x})$ instead of $\mathbf{a}^{(3)}$ in tie breaking in Algorithm 1 and 2 when l is convex.

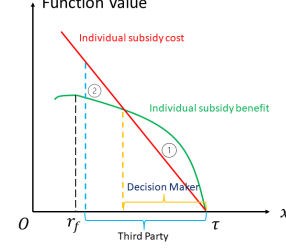


Figure 8: A simplified illustration of the individual subsidy benefit and cost in the mechanism. Region 1/2 corresponds to agents with subsidy surplus/deficit. The third party incentivize region 2 agents for social well-being objectives. r_f represents the lower boundary of $\mathcal{M}(f)$.

The principal chooses \bar{c} by

$$\bar{c} = (\tau - \underline{r}) \Delta c_{i_A}^* / (P^T \boldsymbol{\theta})_{i_A}, \quad \text{where } \underline{r} = \arg \min_r s(r, f, G) \geq 0,$$

intuitively, it incentivizes every agent with non-negative individual subsidy surplus.

Here we highlight some of the key reasons why the mechanism is still IC, IR and satisfies $S(f, G)$ if the condition in equation 9 does not hold.

In fact, when $\mathbf{w} = \boldsymbol{\theta}$ in f , we can assume that a rational principal makes sure if $\mathbf{x} \notin \mathcal{M}(f)$, then $l(r_{\mathbf{x}}) < 0.5 \Leftrightarrow 1 - l(r_{\mathbf{x}}) > l(r_{\mathbf{x}})$. As a result, we know that

$$\bar{s}(r, f, G) = l(\tau) - l(r) - \frac{(\tau - r) \Delta c_{i_A}}{(P^T \boldsymbol{\theta})_{i_A}} \geq s(r, f, G),$$

is concave in r and

$$\underline{s}(r, f, G) = l(\tau) + l(r) - 1 - \frac{(\tau - r) \Delta c_{i_A}}{(P^T \boldsymbol{\theta})_{i_A}} \leq s(r, f, G),$$

is increasing in r , $\forall r$ s.t. $l(r) < 0.5$. Therefore, if the condition in equation 9 does not hold, we have $l(\underline{r}) < 0.5$, where $\underline{r} = \arg \min_r s(r, f, G) \geq 0$ we can also conclude that $r_{\mathbf{x}}$ in $[\underline{r}, \tau]$ satisfies $s(r_{\mathbf{x}}, f, G) \geq 0$, i.e., every agent incentivized has non-negative individual subsidy surplus. \square

Intuitively, the condition $l(\tau) - l(\underline{r}_f) \leq \frac{(\tau - \underline{r}_f) \Delta c_{i_A}^*}{(P^T \boldsymbol{\theta})_{i_A}}$ indicates the subsidy cost is larger than the subsidy gain for an agent on the “far side” boundary of $\mathcal{M}(f)$ in equation 6. This holds when improvement costs are much larger than gaming costs, so that the discount payment is higher than the resulting benefit from the agent’s improvement. Such a condition is needed to enable the efficient calculation of the optimal mechanism for the following reason. If the condition does not hold, the mechanism can further increase the cost discount rate on the actions and let agents with a pre-response attribute such that $\boldsymbol{\theta}^T \mathbf{x} < \underline{r}_f$ to also take improvement actions. However, this would again make the problem hard for the principal, since it has to jointly optimize Δc_j and \bar{c} , and such optimization is non-convex.

We note that the s_1 and s_2 functions in Algorithm 2 capture the following properties of individual subsidy surplus: for agents in $\mathcal{M}(f)$, these agents’ qualification status improvement equals the individual subsidy benefit $l(\boldsymbol{\theta}^T \mathbf{x}'_A) - l(\boldsymbol{\theta}^T \mathbf{x}'_C)$, but for agents not in $\mathcal{M}(f)$, the individual subsidy benefit is not the qualification status improvement, but instead $l(\boldsymbol{\theta}^T \mathbf{x}'_A) - [1 - l(\boldsymbol{\theta}^T \mathbf{x}'_C)]$ since these agents are supposed to receive 0 decision outcomes (rejections) in the CS problem. The green curve in Figure 8 also illustrates the above.

C AUGMENTED STRATEGIC REGRESSION

An agent with pre- (resp. post-) response attribute \mathbf{x} (resp. \mathbf{x}') has a pre- (resp. post-) response *true label* y (resp. y') which indicates the quality of an agent. For strategic regression, we use the same setting as in Shavit et al. (2020):

$$y = q(\mathbf{x}) := \boldsymbol{\theta}^T \mathbf{x} + \eta, \quad y' = q(\mathbf{x}') = \boldsymbol{\theta}^T \mathbf{x}' + \eta, \quad (10)$$

where $\boldsymbol{\theta} \geq 0$ is the quality coefficient vector, and η is a subgaussian noise with 0 mean and variance σ .

The principal's utility is $U_C^{(reg)}(f) = \int_{\mathcal{X}} \mathbb{E}_{\eta} [- (f(\mathbf{x} + P\mathbf{a}_C^*(\mathbf{x})) - y'_C)^2] p(\mathbf{x}) d\mathbf{x}$. Here the principal aims to minimize the mean squared error in regression, respectively. We will use $f_C^* := \arg \max_f U_C(f)$ to denote the principal's optimal conventional strategic decision rule.

For CS and AS regression, the best response directions are the same as CS and AS classification, as given in Lemma B.2.

However, different from the strategic classification problem, the agents can have best responses with infinite magnitude. For example, if $(P^T \mathbf{w})_{i_C} \geq c_{i_C}$, the agent will invest an infinite amount in action i_C . To handle this issue, we assume that the agents' actions are bounded by an action budget $h(\mathbf{a}) \leq B$ in CS (and LS) regression, and $h_A(\mathbf{a}) \leq B$ in AS regression.⁷

Given these bounds on the agent's budgets, the agents' best responses can be characterized as follows: if $(P^T \mathbf{w})_{i_C} \geq c_{i_C}$, then $\mathbf{a}_C^*(\mathbf{x}) = \frac{B}{c_{i_C}} \mathbf{e}_{i_C}$; otherwise $\mathbf{a}_C^*(\mathbf{x}) = \mathbf{0}$. Similarly, let $i_A = \arg \max_j (P^T \mathbf{w})_j / (c_j - \Delta c_j)$, if $(P^T \mathbf{w})_{i_A} \geq c_{i_A} - \Delta c_{i_A}$. Then, the AS-discount best response is $\mathbf{a}_A^*(\mathbf{x}) = \frac{B}{c_{i_A} - \Delta c_{i_A}} \mathbf{e}_{i_A}$; otherwise $\mathbf{a}_A^*(\mathbf{x}) = \mathbf{0}$.

An interesting difference to highlight is that the agents' best responses in strategic classification depend on both the pre-response attributes of the agents and the decision rule, whereas in strategic regression, the best responses are the same for all agents and only depend on the decision rule.

In this strategic regression setting, we will say f incentivizes 0 responses if $\mathbf{a}_C^*(\mathbf{x}) = \mathbf{0}$. Otherwise, if $i_C \leq M_+$ (resp. $i_C > M_+$), we say f incentivizes improvement (resp. gaming).

If f incentivizes non-zero responses (improvement or gaming), the cost discount rates will again follow Lemma B.3, with the minimum effective discount rate is still the same as in equation 7; otherwise, the minimum effective cost discount rate on action j will be such that $(P^T \mathbf{w})_j = (c_j - \Delta c_j)$, $\Delta c_j^* = \max\{c_j - c_{i_C} (P^T \mathbf{w})_j / (P^T \mathbf{w})_{i_C}, c_j - (P^T \mathbf{w})_j\}$.

Using this, the error incurred by the designer on an agent with pre-response attributes \mathbf{x} will consist of two parts, an *equilibrium coefficient error* and an inevitable error due to noises,

$$\mathcal{E}(f, \mathbf{a}, \mathbf{x}) = [\mathbf{w}^T (\mathbf{x} + P\mathbf{a}) - \boldsymbol{\theta}^T (\mathbf{x} + \hat{P}\mathbf{a})]^2 + \text{err}(\eta). \quad (11)$$

Note that although the agents' best responses are independent of \mathbf{x} , the equilibrium individual errors depend on \mathbf{x} for any $\mathbf{w} \neq \boldsymbol{\theta}$.

We next consider the problem of designing an incentive (discount) mechanism.

Theorem C.1. *For general $f(z) = \mathbf{w}^T z$ and $p(\mathbf{x})$, finding the optimal IC, IR, and discount mechanism requires solving non-convex optimization problems and thus is NP-hard.*

Proof. Recall that the AS utility of the principal is

$$U_A^{(reg)}(f) = \int_{\mathcal{X}} \mathbb{E}_{\sigma} [- (f(\mathbf{x} + P\mathbf{a}_A^*(\mathbf{x})) - y'_A)^2] p(\mathbf{x}) d\mathbf{x} - H(G),$$

which if we rewrite the equilibrium individual error as

$$\mathcal{E}(f, \mathbf{a}, \mathbf{x}) = [\mathbf{w}^T (\mathbf{x} + P\mathbf{a}) - \boldsymbol{\theta}^T (\mathbf{x} + \hat{P}\mathbf{a})]^2 + \text{err}(\sigma),$$

the objective becomes

$$U_A^{(reg)}(f) = \int_{\mathcal{X}} -\mathcal{E}(f, \mathbf{a}, \mathbf{x}) p(\mathbf{x}) d\mathbf{x} - H(G).$$

⁷Such bound was not needed in the classification setting, as the fact that $f(z) \leq 1$ naturally provided this.

If f incentivizes non-zero responses, then Algorithm 3 sets Δc_j at the minimum effective discount value, and sets no discount on other actions. Then, it chooses $\underline{c} = 0, \bar{c} = \frac{\alpha B \Delta c_j}{c_j - \Delta c_j}$ so that it incentivizes all agents to take an AS best response $\mathbf{a}_A^*(\mathbf{x}) = \alpha \frac{B}{c_j - \Delta c_j} \mathbf{e}_j + (1 - \alpha) \frac{B}{c_{i_C}} \mathbf{e}_{i_C}$.⁸ If f incentivizes 0 responses, then the principal can choose $\Delta c_j = c_j - (P^T \mathbf{w})_j$ and set $\bar{c} = \alpha B$ in Algorithm 3 so that $\mathbf{a}_A^*(\mathbf{x}) = \alpha \mathbf{e}_j$.

Below, we also discuss the cases when $\mathbf{w} = \boldsymbol{\theta}$, e.g., the principal's optimal LS strategy $f_L^*(\mathbf{z}) = \boldsymbol{\theta}^T \mathbf{z}$.⁹

Lemma C.3. *If $\mathbf{w} = \boldsymbol{\theta}$ in f and f incentivizes 0 responses or improvement, then the optimal IC and IR discount mechanism is $G = 0$.*

This is straightforward since the principal cannot further lower the error from $err(\eta)$ and thus does not want to pay the agents.

If f incentivizes gaming, then the equilibrium individual error becomes, $\mathcal{E}(f, \mathbf{a}_C, \mathbf{x}) = [\boldsymbol{\theta}^T (\mathbf{x} + P \mathbf{a}_C^*) - \boldsymbol{\theta}^T \mathbf{x}]^2 + err(\eta) = (\boldsymbol{\theta}^T P \mathbf{a}_C^*)^2 + err(\eta)$, which is independent of the pre-response attribute \mathbf{x} .

Theorem C.4. *If $\mathbf{w} = \boldsymbol{\theta}$ in f , and f incentivizes gaming, then the optimal IC, IR, and BB $G_d \neq 0$ can be found as follows:*

Choose $i_A = \arg \max_{j \leq M_+} (P^T \boldsymbol{\theta})_j / c_j$ as the target dimension, and set $\Delta c_{i_A} = \Delta c_{i_A}^*$.

Then, derive the alternative form of individual subsidy surplus as $s(\alpha) = (2\alpha - \alpha^2)(\boldsymbol{\theta}^T P \mathbf{a}_C^*)^2 - \alpha B \Delta c_{i_A} (c_{i_A} - \Delta c_{i_A})^{-1}$ and get $\alpha^* = \arg \max_{\alpha \leq 1} s(\alpha) = 1 - \frac{B \Delta c_{i_A} (c_{i_A} - \Delta c_{i_A})^{-1}}{2(\boldsymbol{\theta}^T P \mathbf{a}_C^*)^2}$. Then find the optimal \bar{c} by $\bar{c} = \alpha^* B \Delta c_{i_A} (c_{i_A} - \Delta c_{i_A})^{-1}$.

Proof. In the special case, if improvement is incentivized by the mechanism, it is the dominant strategy to use the minimum effective discount amount, since a higher discount achieves the same error reduction but a higher subsidy cost.

For an AS best response $\mathbf{a}_A = \alpha \frac{B}{c_j - \Delta c_j^*} \mathbf{e}_j + (1 - \alpha) \frac{B}{c_{i_C}} \mathbf{e}_{i_C}$, where $\alpha < 1$, the alternative form of individual subsidy benefit is the reduction in the expected prediction error

$$(\boldsymbol{\theta}^T P \mathbf{a}_C^*)^2 - (1 - \alpha)^2 (\boldsymbol{\theta}^T P \mathbf{a}_C^*)^2,$$

the subsidy cost is $H(G) = \bar{c} = \frac{\alpha B \Delta c_j^*}{c_j - \Delta c_j^*}$, and thus we have the alternative individual subsidy surplus

$$s(\alpha) = (2\alpha - \alpha^2)(\boldsymbol{\theta}^T P \mathbf{a}_C^*)^2 - \alpha B \Delta c_{i_A} (c_{i_A} - \Delta c_{i_A})^{-1}.$$

□

An interesting observation is that the principal does not try to completely remove gaming with the discount mechanism. This is because when the error drops to a sufficiently low level, the marginal subsidy benefit becomes lower than the marginal subsidy cost, which is a constant.

⁸Similar to the classification setting, we let the algorithm put discount on one action dimension. Any $\underline{c} \leq \bar{c}$ is equivalent to both the agents and the designer here since the agent will by default use the discount amount \bar{c} for the maximum improvement. The algorithm can return on condition $S > 0$ as well.

⁹The optimal CS strategy in regression does not guarantee $\mathbf{w} = \boldsymbol{\theta}$ when incentivizing improvement is impossible.

ALGORITHM 3: Grid Search an IC, IR and $S(f, G) > 0$ Mechanism for Regression

Choose $\epsilon > 0$;

$\mathbf{a}_C \leftarrow \frac{B}{c_{i_C}} \mathbf{e}_{i_C}; S_{max} \leftarrow 0$;

$ans \leftarrow (\mathbf{0}, [0, 0])$;

$E_C \leftarrow \int_{\mathcal{X}} \mathcal{E}(f, \mathbf{a}_C, \mathbf{x}) p(\mathbf{x}) d\mathbf{x}$;

for $j = 1 : M_+$ **do**

$\Delta \mathbf{c} \leftarrow \mathbf{0}; S \leftarrow 0; \alpha \leftarrow \epsilon$;

$\Delta c_j \leftarrow c_j - \frac{(P^T \mathbf{w})_j}{(P^T \mathbf{w})_{i_C}} c_{i_C}$;

while $S \geq 0$ **do**

$\alpha \leftarrow \alpha + \epsilon; \bar{c} = \frac{\alpha B \Delta c_j}{c_j - \Delta c_j}$;

$\mathbf{a}_A = \alpha \frac{B}{c_j - \Delta c_j} \mathbf{e}_j + (1 - \alpha) \frac{B}{c_{i_C}} \mathbf{e}_{i_C}$;

$E_A \leftarrow \int_{\mathcal{X}} \mathcal{E}(f, \mathbf{a}_A, \mathbf{x}) p(\mathbf{x}) d\mathbf{x}$;

$S \leftarrow E_C - E_A - \bar{c}$;

end

if $S > S_{max}$ **then**

$S_{max} \leftarrow S; ans \leftarrow (\Delta \mathbf{c}, [0, \bar{c}])$;

end

end

Return ans .

D SUPPLEMENTARY MATERIALS FOR SECTION 4

We provide a more detailed description of group differences as follows. We consider the following set of definitions; the first is new to the best of our knowledge and the other two were introduced in Milli et al. (2019).

Definition 1 (Group Disadvantages). *We say group 2 is*

1. *disadvantaged in attributes in classification if $F^{(2)}(l) > F^{(1)}(l)$ for $l \in (0, 1)$, where $F^{(d)}$ is the cumulative density function (cdf) of the conditional pre-response qualification status conditioned on $d \in \{1, 2\}$; the same in regression if $F^{(2)}(y) > F^{(1)}(y)$ for $y \in (0, \max_{\mathbf{x}} q(\mathbf{x}))$.*
2. *disadvantaged in positive individuals (in classification) if $F_+^{(2)}(l) > F_+^{(1)}(l)$, where $F_+^{(d)}$ is the cdf of conditional pre-response qualification status $(l(\mathbf{x})|Y = 1, D = d)$, $d \in \{1, 2\}$.*
3. *disadvantaged in action cost if $h^{(2)}(\mathbf{a}) > h^{(1)}(\mathbf{a})$, $\forall \mathbf{a} \neq \mathbf{0}$, where h^d denotes the action cost functions with sensitive attribute $d \in \{1, 2\}$. Moreover, the minimum effective discount values satisfy $(\Delta c^{(1)})_i^* \leq (\Delta c^{(2)})_i^*$, $\forall i$.*

D.1 FAIRNESS ISSUES IN THE CS/LS EQUILIBRIUM

We start with a number of fairness limitations of the CS equilibria in classification and regression; the same results apply to LS.

Theorem D.1. *In the equilibrium CS outcome of classification where two groups have the same action cost, then (i) if group 2 is disadvantaged in attributes, then there is a DP gap no matter if f incentivizes improvement or gaming; and (ii) if group 2 is disadvantaged in positive individuals, then there is an EO gap if f incentivizes gaming but not necessarily if f incentivizes improvement.*

Part (1) is a direct result of $1 - F^{(1)}(l) > 1 - F^{(2)}(l)$, and the two groups have the same implicit threshold, which is the lower side boundary of their manipulation margins (since every agent above it will manipulate to get a positive decision outcome), and $\mathcal{M}^{(1)}(f) = \mathcal{M}^{(2)}(f)$ since the two groups have the same action cost. For part (2), whether there is a quality gain gap entirely depends on whether f incentivizes improvement and the distribution of each group in its manipulation margin $\mathcal{M}^{(d)}(f)$. For example, we can have $Pr(\mathbf{x} \in \mathcal{M}^{(2)}(f)|D = 2) > Pr(\mathbf{x} \in \mathcal{M}^{(1)}(f)|D = 1)$ and thus group 2 have more agents to improve and may have an inverse quality gain gap.

Theorem D.2. *In the equilibrium CS outcome of classification and regression, if group 2 is disadvantaged in action cost but has the same pre-response attribute distribution as group 1 (for positive individuals as well), then there is (i) a quality gain gap only if f incentivizes improvement; (ii) an EO gap no matter if f incentivizes improvement or gaming; and (iii) a DP gap no matter if f incentivizes improvement or gaming.*

Proof. The DP gap is only related to $f(\mathbf{z})$ but not y or y' , when the two groups have the same action cost but group 2 is disadvantaged in attribute, the implicit threshold (the lower side boundary of $\mathcal{M}^d(f)$, $\hat{\tau}_L$) is the same for both groups and from the definition of attribute disadvantage, $PR^{(1)} = 1 - F^{(1)}(\hat{\tau}_L) > 1 - F^{(2)}(\hat{\tau}_L) = PR^{(2)}$, and we know that the DP gap exists.

When f incentivizes gaming, the reason of an EO gap is similar as above $TPR^{(1)} = 1 - F_+^{(1)}(\hat{\tau}_L) > 1 - F_+^{(2)}(\hat{\tau}_L) = TPR^{(2)}$. If f incentivizes improvement, then the EO gap depends on both the CS TPR in both groups, the CS PR in both groups, and the AS quality improvement in both groups. For example, if G only not incentivize agents in the manipulation margins, then

$$TPR_A = 1 - FNR_A = 1 - FNR_C \cdot \frac{PR_C}{PR_A} = 1 - \frac{(1 - TPR_C) \cdot PR_C}{\Delta Q_A + PR_C},$$

and we know that the AS EO gap depends on $\Delta Q_A^{(1)}, \Delta Q_A^{(2)}$ which is based on $p^{(1)}(\mathbf{x})$ and $p^{(2)}(\mathbf{x})$ and we can not easily conclude the EO gap changes. \square

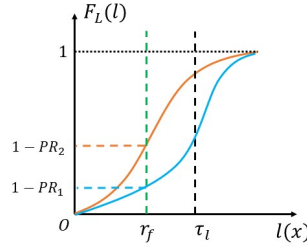


Figure 10: An illustration of the CS DP gap when group 2 is disadvantaged in attributes.

To understand the above result, we note that if group 2 is disadvantaged in cost, we have $\mathcal{M}^{(1)}(f) \supseteq \mathcal{M}^{(2)}(f)$, so even when group 2 has the same pre-response attribute distribution, a larger portion of group 1 are accepted in the equilibrium, causing the DP gap. This is similar to the reason of an EO gap when f incentivizes gaming. If f incentivizes improvement, then a larger portion of group 1 will improve and be accepted in the equilibrium, causing a quality gain gap and an EO gap simultaneously.

D.2 INFLUENCE OF THE DISCOUNT MECHANISM ON FAIRNESS

Here we analyze how the discount mechanism G alone may influence the fairness.

Theorem D.3. *If group 2 is disadvantaged in cost but has the same pre-response attribute distribution, then a rational principal will choose a G that widens the quality gain gap in both classification and regression.*

Proof. If group 2 is disadvantaged in cost, then it is cheaper to incentivize a group 1 agent than a group 2 agent to get the same qualification status improvement, and thus the principal subsidizes more group 1 agents and creates a quality gain gap. \square

Theorem D.3 means that a rational mechanism for the principal is always making the system more unfair when the quality gain gap is the metric. The rational mechanism influences the DP and EO gap but does not always make them worse.

D.3 ANALYTICAL RESULTS ON THIRD PARTY MECHANISM DESIGN

The mechanism designer can induce truthful revelation of the sensitive attribute by the agents as follows: (1) Let G consist of two group-specific mechanisms $G^{(1)}$ and $G^{(2)}$; agents who do not reveal their d participate in $G^{(1)}$; (2) Ensure that $\Delta c_i^{(1)} \leq \Delta c_i^{(2)}, \forall i$ and $(\Delta \mathbf{c}^{(1)})^T \mathbf{a} \in [\underline{c}^{(1)}, \bar{c}^{(1)}] \Rightarrow (\Delta \mathbf{c}^{(2)})^T \mathbf{a} \in [\underline{c}^{(2)}, \bar{c}^{(2)}]$. Then, group 1 agents are indifferent about revealing d while revealing d is the dominant strategy for group 2 agents. Figure 2 illustrates the three-party AS learning system.

Theorem D.4. *If there is a mechanism that is IC and IR and satisfies $S(f, G) > 0$, then a mechanism that satisfies IC, IR, and BB criteria exists and weakly improves the third party's social well-being objective (either efficiency or fairness oriented) compared to the original AS equilibrium.*

Proof. We still need G to be IR for the principal, where the maximum tax a rational principal accepts is the subsidy benefit $\mathcal{T}(G) \leq S(f, G) + H(G)$, and the BB condition requires $S(f, G) + H(G) \geq \mathcal{T}(G) \geq H(G)$. So, as long as $S(f, G) \geq 0$, there is an IC, IR, and BB third party mechanism. Therefore, finding the optimal IC, IR, and BB third party mechanism is the same as

$$\text{maximize}_G W(f, G), \quad \text{subject to } S(f, G) \geq 0,$$

and if $S(f, G) > 0$ the mechanism can further improve its objective by setting the surplus at 0. \square

We also discuss how the objective of the mechanism designer and the corresponding incentive mechanisms influence the equilibrium efficiency and fairness oriented social well-being metrics. We compare the different AS, CS, and LS equilibrium outcomes where they have the same decision rule f and focus on how the incentive mechanisms for different objectives affect the outcome.

Definition 2. We say a mechanism $G^d \neq 0$ is an ideal mechanism if it is IC and IR for group d agents and achieves $S(f, G^d) > 0$ on group d , $\forall d \in \{1, 2\}$.

Theorem D.5. If group 2 is disadvantaged in action cost but has the same pre-response attribute distribution as group 1 (for positive individuals as well), then in the equilibrium,

1. the DP gap in weak ascending order is: AS-fair, CS(LS), AS-dm, AS-eff;
2. the EO gap (or quality gain gap) in weak ascending order is: AS-fair, CS(LS), AS-dm, AS-eff;
3. The social quality improvement in weak descending order is: AS-eff, AS-dm, CS(LS).

If there is an ideal mechanism for group 1, then AS-fair is strictly the lowest in DP gap; the orders in EO gap (or quality gain gap) and quality improvement becomes strict for CS(LS), AS-dm and AS-eff. Moreover, if there is an ideal mechanism for group 2, AS-fair is strictly the lowest in EO gap (or quality gain gap).

Proof. For Part (1), the fairness oriented third party can implement the ideal mechanism on group 2 and even further subsidize other group 2 agents to reduce the gap while avoiding subsidizing more group 1 agents to enlarge the fairness gaps.

For Part (2), any ideal mechanism makes sure the efficiency oriented third party has “remaining budget” to incentivize more agents to improve compared to AS-dm outcome and thus has the strictly highest equilibrium social quality improvement. \square

Below we provide some explanations about the statements in Theorem 4.1. For an efficiency oriented third party, the set of agents it incentivizes is a superset of the agents incentivized by the principal, making AS-eff the best in part (3). This is because subsidizing the agents with a positive individual subsidy surplus not only helps the third party improve the objective but also raises the budget to subsidize agents with a negative individual subsidy surplus (individual subsidy deficit). Moreover, the efficiency oriented third party tries to incentivize more agents from group 1 since they are “cheaper” to incentivize and thus exacerbates the fairness issues in part (1) and (2).

For a fairness oriented third party, it can also incentivize a superset of agents incentivized by the principal, but that means incentivizing some group 1 agents, which results in two contradicting effects: it helps the third party gather more “funding” to subsidize group 2 agents, but also makes the fairness issue worse simultaneously. As a result, the social quality improvement in AS-fair is better than CS(LS) and worse than AS-eff, but how it compares to AS-dm depends on the specific game parameters and thus is not discussed in part (3). When there is an ideal mechanism for group 2, the third party can ignore the dilemma of subsidizing group 1 agents and focus on subsidizing only group 2 agents to improve fairness in part (1) and (2).

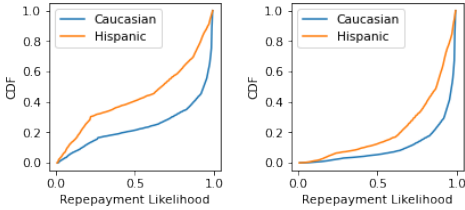
The ideal mechanisms in Theorem 4.1 makes the comparison strict. The existence of an ideal $G^{(2)}$ is a sufficient condition to the existence of an ideal $G^{(1)}$ when group 2 is disadvantaged in cost but has the same distribution. This is because $G^{(2)}$ itself is ideal for group 1.

Theorem D.6. In both classification and regression problems, if group 2 is disadvantaged in attributes (resp. positive individuals) but has the same action cost as group 1 then

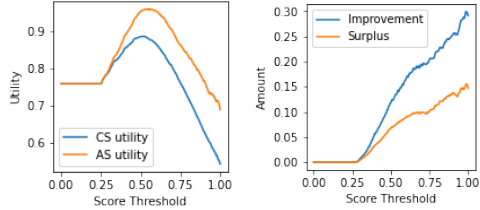
1. the DP (resp. EO) gap in AS-fair outcome is weakly the lowest, and is strictly the lowest if there is an ideal mechanism for group 2;
2. the social quality improvement in AS-eff outcome is weakly the highest, and is strictly the highest if there is an ideal mechanism for either group.

When group 2 has the same cost, then an ideal $G^{(2)}$ is no longer sufficient or necessary for an ideal $G^{(1)}$ to exist for general classification problems, and that’s why the condition in part (2) looks different from Theorem D.6. But the existence of an ideal $G^{(2)}$ is sufficient and necessary to the existence of an ideal $G^{(1)}$ in regression, as well as in a special class of classification problems where $\mathbf{w} = \boldsymbol{\theta}$ in f and l is convex on $[0, \tau]$.¹⁰ From Theorem D.1, we know that DP and EO gap always

¹⁰We are excluding extreme distributions in the “iff” claim, e.g., $Pr(\mathbf{x} \in \mathcal{M}(f)) = 0$.



(a) Entire Group (b) Positive Individuals
Figure 12: The Likelihood CDF



(a) CS/AS Utilities (b) ΔQ and $S(f, G)$
Figure 13: Single Group (Caucasian) Results

exist in the CS(LS) problem, but if there is an ideal $G^{(2)}$, the fairness oriented third party can further incentivize group 2 agents to reduce the gap in part (1) (those not in $\mathcal{M}^{(2)}(f)$ to reduce the DP gap).

Theorem D.7. *Suppose group 2 is disadvantaged in cost but has the same pre-response distribution (for positive individuals as well) Denote $p^{(d)} := \Pr(D = d)$, then an IC, IR, and BB mechanism $G \neq 0$ that satisfies $\gamma_A^Q = \gamma_A^{EO} = \gamma_A^{DP} = 0$ exists if $S(f, G^{(1)}) + (1 - p^{(1)})H(G^{(1)}) \geq p^{(2)}H(G^{(2)})$, s.t. $h_A^{(1)}(\mathbf{a}) = h_A^{(2)}(\mathbf{a}), \forall \mathbf{a}$.*

Proof. If $h_A^{(1)}(\mathbf{a}) = h_A^{(2)}(\mathbf{a}), \forall \mathbf{a}$, then the equilibrium feature and attribute distribution are the same for both groups, and thus there is no fairness gap. Meanwhile, the subsidy benefit are the same in both groups, so the overall benefit is $S(f, G^{(1)}) + H(G^{(1)})$, and the overall subsidy cost is $p^{(1)}H(G^{(1)}) + p^{(2)}H(G^{(2)})$. \square

In general, this condition can hold if $p^{(1)}$ is much larger than $p^{(2)}$, i.e., the disadvantaged group is also the minority group in the population or $S(f, G^{(1)})$ is very high.

Remark 2. *Our results generalize to multiple groups when the definitions of group disadvantages and fairness metrics are consistent.*

E FULL NUMERICAL EXPERIMENT AND RESULTS

This section presents numerical results obtained using the FICO score Reserve (2007) dataset preprocessed in Hardt et al. (2016b). The credit card holders are considered as agents and they have repayment rates that can map to the likelihood function l in our model. The principal uses binary classification to predict whether the agents will default. We assume that $\theta = 1, P = [1, 1]$, and the agent can either choose a_1 to improve or a_2 to game the classifier $f(z) = \mathbf{1}(z \geq \tau)$, i.e., x is the pre-response normalized FICO score as well as the attribute, $x' = x + a_1$ is the post-response attribute, and $z = x + a_1 + a_2$ is the post-response normalized FICO score. Figure 11 shows how the repayment rate $l(x)$ changes with x ; it has an S-shape, with $l(x) = 0.5$ approximately corresponding to $x = 0.3$ and $l(x)$ (nearly) convex on $[0, 0.3]$. We assume that the principal chooses $w = 1$, which aligns with the LS and CS optimal solution from Section 3 when $c_2 < c_1$.

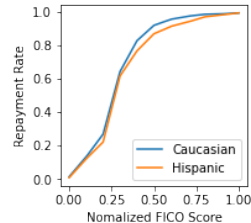


Figure 11: Repay Rate $l(x)$

We start with the properties of the discount mechanism and show how the principal’s CS and AS utility changes with different choices of threshold τ . We then show the impact the incentive mechanisms have on social well-being metrics.

Throughout this section, we use a quadratic outcome likelihood cost function and assume that $c_1^{(1)} = c_1 = 8$ and $c_2^{(1)} = c_2 = 4$ (for the advantaged group if there are action cost differences). For the multiple group case, we make the following two sets of comparisons. (1) Groups with different distributions: the Hispanic group is disadvantaged in features and in positive individuals compared to the Caucasian group (see Figure 12). (2) Groups with different costs: we will assume there are two subgroups (A and B) in the Caucasian group, and group 2 has higher action costs $c_1^{(2)} = 10$ and $c_2^{(2)} = 5$. We set $p^{(1)} = 0.8, p^{(2)} = 0.2$ as the population proportions.

As a result, we show the AS-fair equilibrium outcome is the best well-rounded system design for the augmented strategic learning problems.

The principal’s AS and CS utility. Using only the Caucasian data, the set of results in Figure 13 show how the AS/CS principal utilities, subsidy surplus and qualification status improvement change with the threshold τ .

We can see that the AS utility is always higher than the CS utility (Fig. 13). This is because their difference is the subsidy surplus, which is non-negative for a rational principal. We note that the CS utility should always be single-peaked but the AS utility may have multiple local maxima since the value of subsidy surplus is not monotone in τ and depends on $p(x)$. For other choices of c_1, c_2 values, we find that the larger the difference $c_1 - c_2$, the smaller the utility difference and the closer the optimal thresholds are ($|\tau_{AS}^* - \tau_{CS}^*|$ lower). Both the subsidy surplus and the qualification status improvement are positive, indicating the principal’s selfish strategy is also benefiting the efficiency oriented social well-being. The improvement and subsidy surplus are also highly positively correlated with a correlation coefficient of 0.92.

Social well-being of the strategic incentive mechanism. Figure 14 (resp. Figure 15) shows the quality improvement, PR and TPR, (and thus we can see the DP, and EO gap from the curve differences) when the Hispanic group (resp. Caucasian subgroup 2) is disadvantaged in features and positive individuals (resp. costs) compared to the Caucasian group (resp. Caucasian subgroup 1) in the CS(LS) and AS-dm equilibrium. The principal does not incentivize agents outside of the manipulation margin and thus the CS and AS PR curves are the same.

We can see from Figure 14a that when τ is in the lower score ranges, the Hispanic group has a slightly higher qualification status improvement compared to the Caucasian group, whereas if τ is in the higher score ranges, the Caucasian group has a much higher improvement. Intuitively, this is because the Hispanic (resp. Caucasian) group has a higher probability mass in the lower (resp. higher) score ranges and a low (resp. high) τ incentivizes a higher proportion of agents to improve in the Hispanic (resp. Caucasian) group. Figure 14b shows that the PR is 1 when $\tau < 0.25$; this is because all agents can manipulate to get $f(z) = 1$. When $\tau > 0.25$, the PR is strictly decreasing in τ for both groups and the Caucasian group always has a higher PR, i.e., the Hispanic group will suffer from a DP gap in both CS and AS-dm equilibrium. This is because the lower side boundary of the manipulation margin becomes an implicit threshold, where all agents above the implicit threshold can manipulate (no matter improvement or gaming) to get accepted. The implicit threshold is the same for both groups since they have the same action cost, and the DP gap is caused by the disadvantage in pre-response attribute distribution (Theorem D.1 part (1)). For similar reasons, Figure 14c shows that the CS and AS TPR is 1 when $\tau < 0.3$. Therefore, we can see that the AS TPR is always higher than the CS TPR for either group, because now some agents improved their qualification status and get accepted at the same time, making the numerator and denominator of the TPR formula increase by the same amount and thus increase the TPR. On the other hand, the Hispanic group suffers from an EO gap in both the CS and the AS-dm equilibrium, as previously discussed in Theorem D.1 part (2).

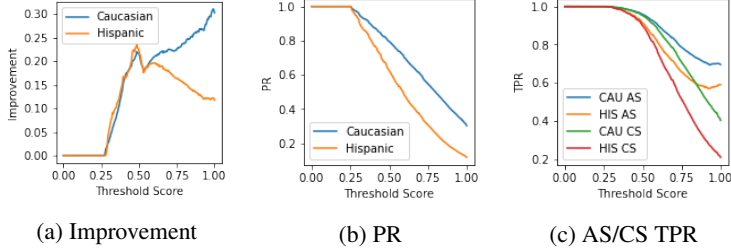


Figure 14: Disadvantaged in features

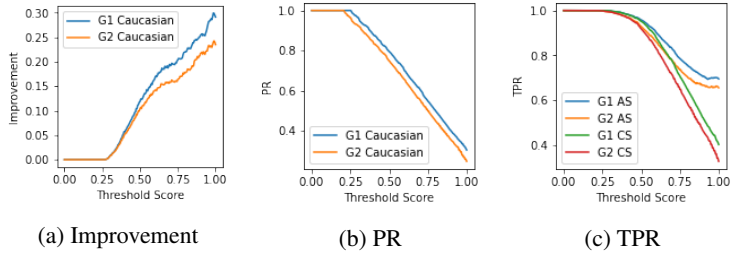


Figure 15: Disadvantaged in costs

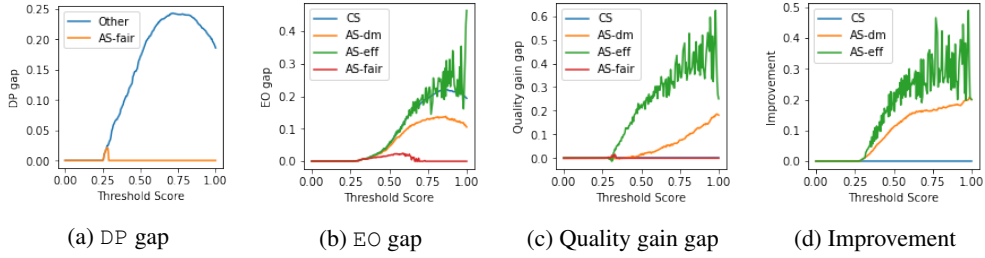


Figure 16: Third Party Outcomes with Attribute Distribution Differences

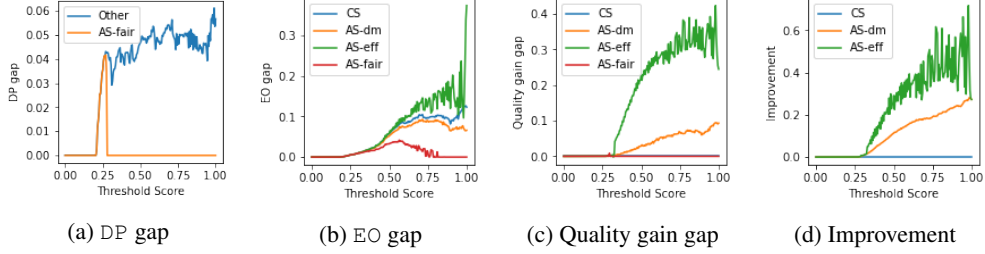


Figure 17: Third Party Outcomes with Cost Differences

Figure 15a and 15c support our claims in Theorem D.3 part (3), where the incentive mechanism widens the quality gain gap and the EO gap. Figure 15b shows PR curves and the DP gap between the two subgroups, which is determined by the pre-response attribute probability mass within $[\tau - 1/c_2^{(1)}, \tau - 1/c_2^{(2)}]$ (the difference between the manipulation margins in the two groups). Figure 15c shows the CS and AS TPR curves and the EO gaps; the implicit threshold creates the CS EO gap, and the fact that group 1 agents are cheaper to incentivize jointly creates the AS EO gap.

Social Well-being metrics with the third party incentive. Social well-being results under the third party model are shown in Figure 16 where groups have attribute distribution differences (Caucasian and Hispanic group), and in Figure 3 where groups have cost differences (Caucasian subgroups).

We can see in both sets of results that the AS-fair equilibrium outcome significantly reduces and even removes the fairness issues in the system, whereas the AS-eff equilibrium outcome has the worst fairness metrics. On the other hand, the AS-eff equilibrium achieves the highest social qualification status improvement. We note that the chosen AS-fair outcomes used mechanisms that incentivized a superset of agents compared to those that are incentivized by the principal, and thus it achieves a higher social qualification status improvement than AS-dm as well.