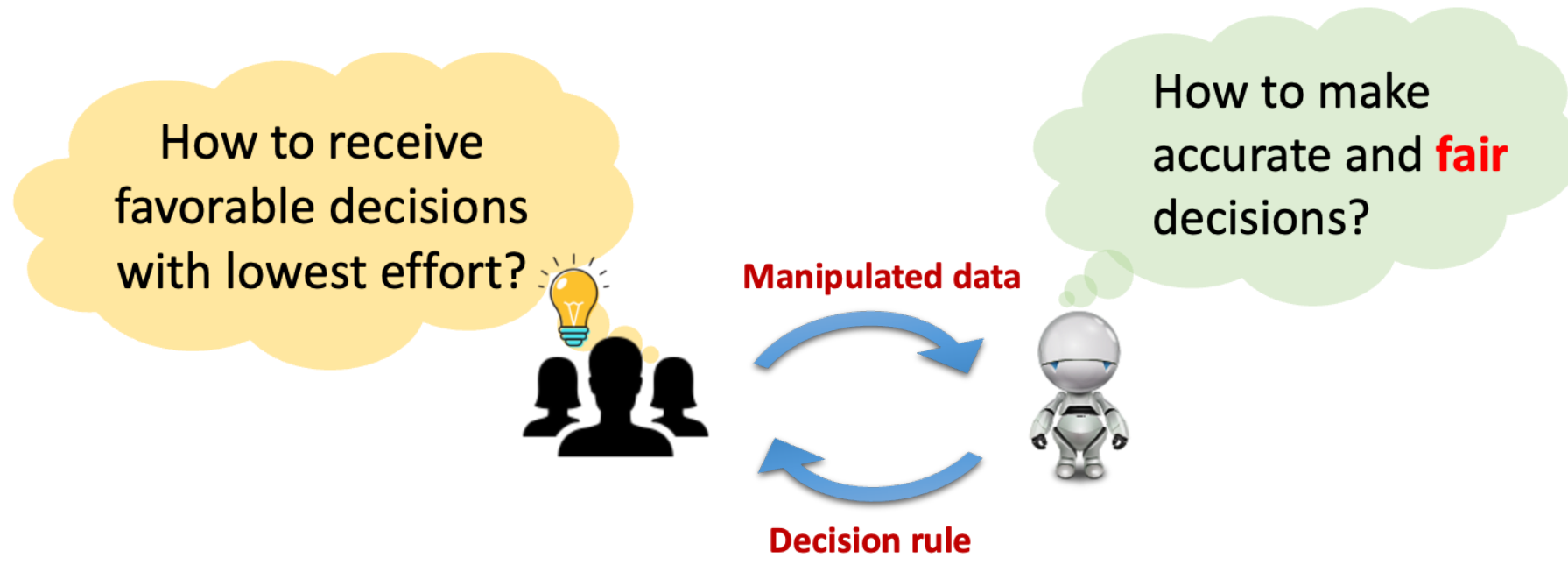


Fairness Interventions as (Dis)Incentives for Strategic Manipulation

Xueru Zhang¹, Mohammad Mahdi Khalili², Kun Jin³, Parinaz Naghizadeh¹, Mingyan Liu³

BACKGROUND

- ML has been increasingly used to help make decisions about people
 - lending, hiring, college admission, ...
- Two challenges:
 - ML is vulnerable to **strategic manipulation**
 - ML can be **biased** against certain social groups

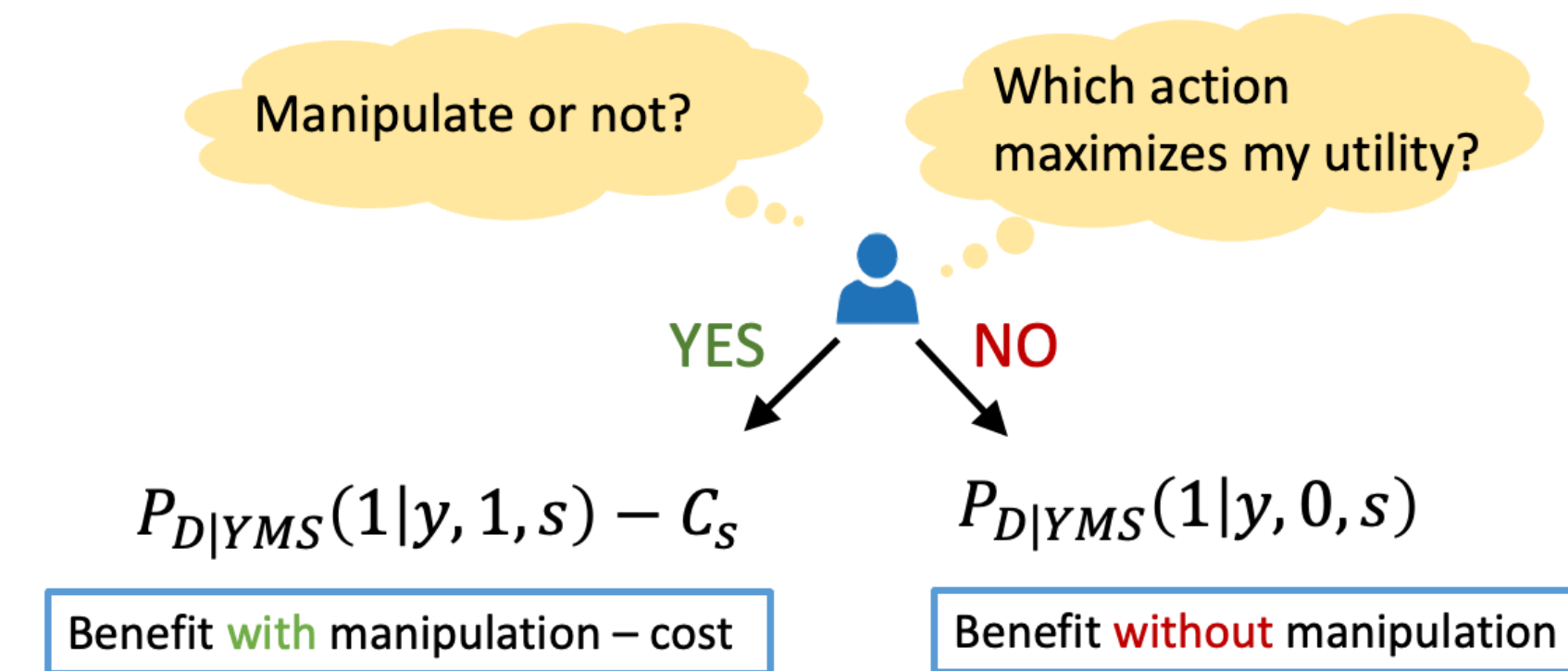


MODEL: STRATEGIC INTERACTION

Two demographic groups $\mathcal{G}_a, \mathcal{G}_b$

- Sensitive attribute $S \in \{a, b\}$
- Features $X \in \mathbb{R}^d$
 - feature generation** $P_{X|YS}(x|y, s)$
- Qualification state $Y \in \{0, 1\}$
 - qualification rate** $\alpha_s = P_{Y|S}(1|s)$
- Decision $D \in \{0, 1\}$
 - policy** $\pi_s(x) = P_{D|XS}(1|x, s)$
- Manipulation action $M \in \{0, 1\}$
 - Manipulation doesn't affect Y but results in better feature distribution
 - Manipulation cost $C_s \geq 0$

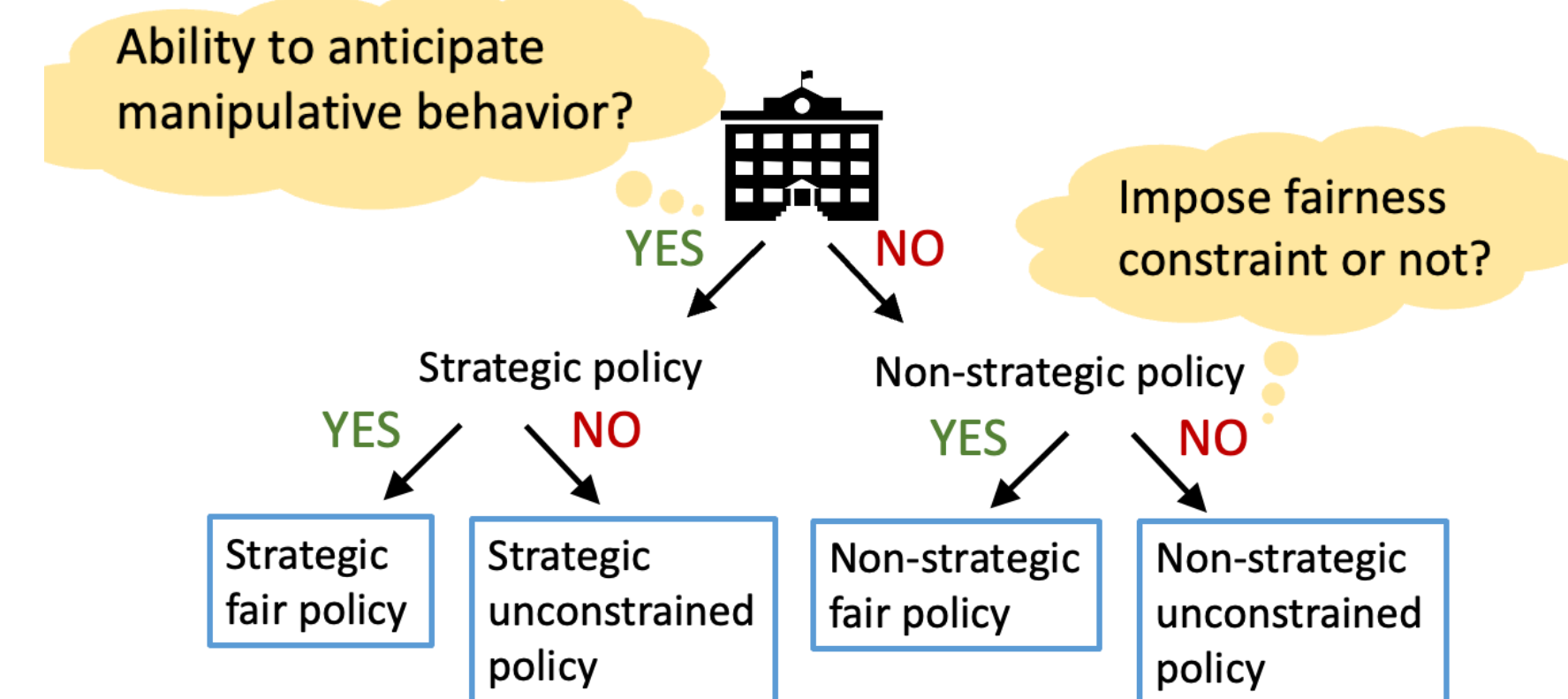
Individual best response



- Individual chooses to manipulate at cost if manipulation brings the higher utility
- Manipulation probability:

$$\Pr(C_s \leq P_{D|YMS}(1|y, 1, s) - P_{D|YMS}(1|y, 0, s))$$

Decision-maker's optimal (fair) policies



$$\begin{aligned} \max_{\pi_a, \pi_b} \quad & \mathbb{E}[R(D, Y)] \\ \text{s.t.} \quad & \text{fairness constraint} \end{aligned}$$

EXISTING WORK

Fair machine learning



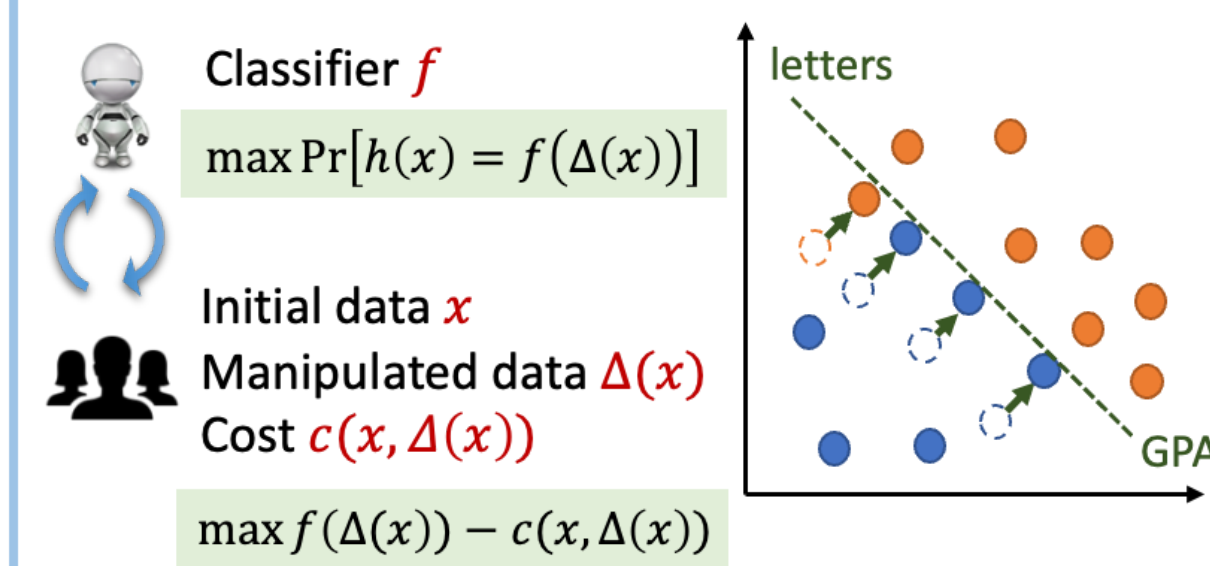
$$\begin{aligned} \min \quad & \text{Loss} \\ \text{s.t.} \quad & \phi(\text{blue}) \approx \phi(\text{red}) \\ & \text{Fairness constraint} \end{aligned}$$

- Demographic parity:** equal positive rate
- Equal Opportunity:** equal true positive rate

Strategic classification

Stackelberg game formulation

Hardt et al., 2016a; Dong et al. 2018; Milli et al., 2019; Hu et al., 2019; Braverman & Garg, 2020



- Most works studied these two problems separately
- Existing Stackelberg game formulation assumes:
 - Manipulation outcome is **deterministic & known**
 - Manipulation cost is a **deterministic function** of features before & after manipulation

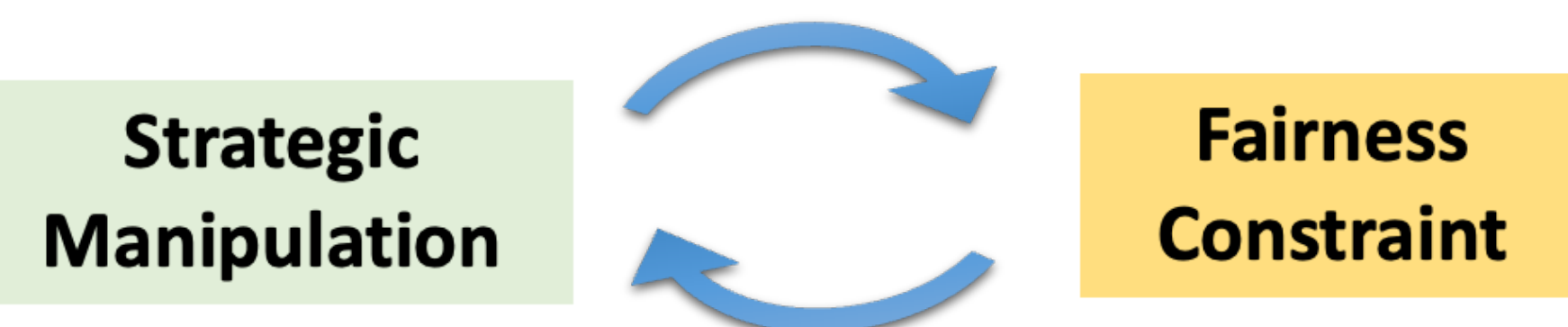
THEORETICAL RESULTS

- Characterize the equilibrium strategies of individuals & decision-maker (four types of policies)
- Impact of decision-maker's anticipation of strategic manipulation
 - Strategic policy **over(under)** accepts **majority-qualified(majority-unqualified)** group
 - Anticipation of manipulation can **worsen** the fairness of a strategic policy when one group is **majority-qualified** while the other is **majority-unqualified**
 - When both groups are **majority-unqualified**, a strategic policy may **mitigate** unfairness and even **flip the disadvantaged group**
- Impact of fairness interventions on policies and individuals' manipulation
 - Conditions under which non-strategic decision maker may **benefit from fairness constraints**
 - Conditions under which fairness constraints serve as **(dis)incentives** for strategic manipulation

THIS WORK



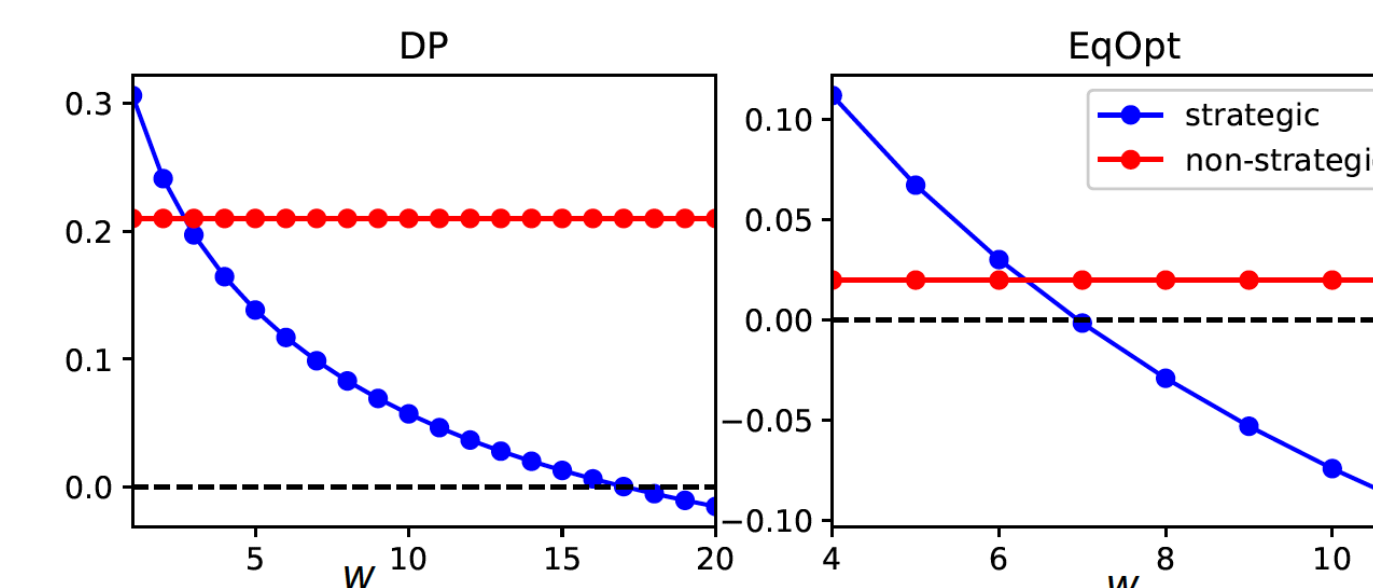
- A new Stackelberg game formulation:
 - Uncertain** manipulation outcomes
 - Manipulation cost is determined before observing manipulation outcomes
- Understand the impacts **strategic manipulation** and **fairness intervention** have on each other



Motivating Example: students cheat on exams to get admitted to college

EXPERIMENTS: FICO CREDIT SCORE

- Hispanic & Black:** strategic policy mitigates unfairness



- Black & \mathcal{G}_a :** strategic policy worsens unfairness

	\mathcal{G}_a	strategic		non-strategic
		$C_a = C_b$	$C_a \neq C_b$	
EqOpt	Caucasian	0.355	0.556	0.136
	Hispanic	0.292	0.493	0.034
	Asian	0.333	0.533	0.123
DP	Caucasian	0.611	0.680	0.449
	Hispanic	0.421	0.490	0.242
	Asian	0.634	0.703	0.522

- White & Asian:** non-strategic fair policy has higher utilities

δ_u	C_a	$U_a(\hat{\theta}_a^{JN})$	$U_a(\hat{\theta}_a^C)$	$U_b(\hat{\theta}_b^{JN})$	$U_b(\hat{\theta}_b^C)$
0.8	Beta(10, 10)	-0.190	-0.189	0.024	0.034
0.756	Beta(10, 1)	0.396	0.397	0.181	0.201

- More results:**

