# Demographic Parity Constrained Minimax Optimal Regression under Linear Model

# Summary

- We investigate the minimax optimality of regression with the constraint of demographic parity.
- Our model poses the following additional challenges compared to the existing results of Chzhen et al. (2022):
- (Direct discrimination) Mitigating outcome's variance disparity. • (Indirect discrimination) Addressing indirect discrimination.
- We reveal the minimax optimal error rate as  $\sigma_{\xi}^{2}B^{2}dM/n$ .

# Setup

Consider X as non-sensitive features in  $\mathbb{R}^d$  and S as a sensitive feature within [M]. Given noise  $\xi \sim N(0, \sigma_{\varepsilon}^2)$ , the outcome Y is:

$$Y = f^*(X, S) + \xi.$$

# Fairness

### Definition: demographic parity (Pedreshi et al. 2008)

A regressor f satisfies (strong) demographic parity if for all  $s, s' \in [M]$ , and for all  $E \in \sigma(f(X,S))$ ,

$$\mathbb{P}\{f(X,S) \in E | S = s\} = \mathbb{P}\{f(X,S) \in E | S = s'\}$$

- Fairness consistency requires the learned regressor to approach an (exactly) fair regressor as n tends to infinity.
- We use the Wasserstain distance-based unfairness score to define "approaching".

$$U(f) = \max_{s,s' \in [M]} W_2(\nu_{f|s}, \nu_{f|s'})$$

### Definition: $(\alpha, \delta)$ -fairness consistency

A learning algorithm is  $(\alpha, \delta)$ -consistently fair for an unfairness score U if there exists constants  $n_0 \geq 0$  and C > 0 independent of n such that  $\mathbb{P}\{U(\hat{f}_n) > Cn^{-\alpha}\} \leq \delta \text{ for all } n \geq n_0.$ 

# Accuracy

- Goal: to obtain a fair version of  $f^*$ , defined as  $f_{\rm DP}^* = \arg\min_{f \in \mathcal{F}_{\rm DP}(\mu)} \mathbf{E}[(f(X, S) - f^*(X, S))^2].$
- Inaccuracy of f is measured by the mean squared deviation from  $f_{\text{DP}}^*$ :

$$\mathcal{E}(f;\beta_{\cdot}^*,\mu_{\cdot}) = \mathbf{E}\left[\left(f(X,S) - f_{\mathrm{DP}}^*(X,S)\right)^2\right].$$

### Definition: minimax optimal error

Given  $\alpha > 0$  and  $\delta \in (0, 1)$ , the minimax optimal error is defined as  $\mathcal{E}_n(\alpha, \delta) = \inf_{\hat{f}_n:(\alpha, \delta) \text{-consistently fair } \beta^* \in \mathcal{B}, \mu \in \mathcal{M}} \sup \mathbf{E} \Big[ \mathcal{E}(\hat{f}_n; \beta^*, \mu) \Big],$ 

Kazuto Fukuchi<sup>1,3</sup> and Jun Sakuma<sup>2,3</sup>

<sup>1</sup> University of Tsukuba  $^{-2}$  Tokyo Institute of Technology <sup>3</sup> RIKEN AIP

# **Sources of Unfairness (Direct v.s. Indirect)**

- (Direct discrimination) Sensitive attribute directly affects the outcome, regardless of non-sensitive features.
- (Indirect discrimination) Sensitive attribute indirectly affects the outcome through its correlation with non-sensitive features.



Indirect discrimination

# Models

- Chzhen et al. (2022) is the sole study demonstrating minimax optimality in fair regression.
- Contrast with Chzhen et al. (2022): our model accounts for a broader source of discrimination.

		partial coefficients	intercept <sup>r</sup>	non-sensitive features	
	Chzhen et al. (2022) ours	$\checkmark$	$\checkmark$	$\checkmark$	
Chzhen	et al. (2022)'s model				
Non-s	sensitive features' mo	del: for a po $X \sim N(0,$	sitive semi- $\Sigma$ ).	definite matrix $\Sigma$ ,	, (6)
No dependency on $S$ . No indirect discrimination.					
<ul> <li>Outco</li> </ul>	ome's model:	$r / Q^* \mathbf{V}$	h + C		(7)
	Intercepts ( $b_s$ ) may cause direct discrimination.				
Our mo	del				
Non-s	sensitive features' mo Means depend on	del: for $\sigma_X^2 > X \sim N(\mu_S, q)$ on <i>S</i> , leading	> 0, $\sigma_X^2 I$ ). g to indirect	t discrimination.	(8)
<ul> <li>Outco</li> </ul>	ome's model:	$Y = \langle \beta^* X$	$\rangle + \xi$		(9)

No dependency on 
$$S$$
. No indir

$$Y = \langle \beta^*, X \rangle + \frac{b_s}{2} +$$

$$Y = \underbrace{\langle \beta_s^*, X \rangle}_{} + \xi$$

Both partial coefficients and intercept may cause direct discrimination.



(3)



(4)





# (Direct discrimination)

- disparities.

# (Indirect discrimination)

- maintaining consistent output across varying  $\mu_S$ .

. There is a finite universal constant B > 0 such that

Greater variation in  $\|\beta_s^*\|$  increases B, characterizing dispersion of outcome variances.

Main theorem

Given  $\alpha \in (0, 1/2]$  and  $\delta \in (0, 1)$  $4\ln(M/\delta))/\min_{s\in[M]}p_s.$ 

$$c\frac{\sigma_{\xi}^2 B^2 dM}{n} - o\left(\frac{1}{n}\right) \leq \mathcal{E}_n(\alpha,\delta)$$

The upper bound is achieved by a carefully designed plugin estimator (See our paper for detail).

- on U and  $\sigma_X^2$ .
- The term  $\sigma_{\xi}^2 dM/n$  aligns with standard non-fair regression.
- outcome's variances.
- cost-free if X's dependence of S is solely on its mean.

References Chzhen, Evgenii and Nicolas Schreuder (Aug. 2022). "A minimax framework for quantifying risk-fairness trade-off in regression". In: The Annals of Statistics 50.4, pp. 2416–2442. ISSN: 0090-5364. DOI: 10.1214/22-A0S2198. Pedreshi, Dino, Salvatore Ruggieri, and Franco Turini (2008). "Discrimination-aware data mining". In: Proceeding of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining - KDD 08, pp. 560–568. ISBN: 9781605581934. DOI: 10.1145/1401890.1401959.

# Challenges

 Variability in partial coefficients leads to varied outcome variances against S, while diverse intercepts only change the outcome's mean. Our model poses a challenge of addressing both variance and mean

• Our model introduces indirect discrimination via  $\mu_S$  changes relative to S. • Counteracting this requires estimating  $\mu_S$  to fine-tune the regressor,

# Main result

 $\|\beta_s^*\| \le B \text{ and } \frac{(\sum_{s \in [M]} p_s \|\beta_s^*\|)^2}{M} \sum_{[M]} \|\beta_s^*\|^{-2} \le B^2$ (10)

2. There exists a finite universal constant U > 0 such that  $\|\mu_s\| \leq U$ .

), suppose 
$$M(d-1) > 16$$
 and  $n \ge 12(3d \lor d)$ 

$$\leq C \frac{\sigma_{\xi}^2 B^2 dM \vee \sigma_X^2 B^2 M \vee B^2 U^2}{n} + o\left(\frac{1}{n}\right). \tag{11}$$

# Implications

• The constructed estimator is minimax optimal up to constant depending

• (Direct discrimination) The greater the dispersion of outcome variances, the more difficult it becomes to mitigate direct discrimination due to the

• (Indirect discrimination) Indirect discrimination can be mitigated

