

Group Fairness with Uncertainty in Sensitive Attributes

Abhin Shah, Maohao Shen, Jongha Jon Ryu, Subhro Das,
Prasanna Sattigeri, Yuheng Bu, Gregory W. Wornell



IBM **Research**



MIT-IBM
Watson
AI Lab

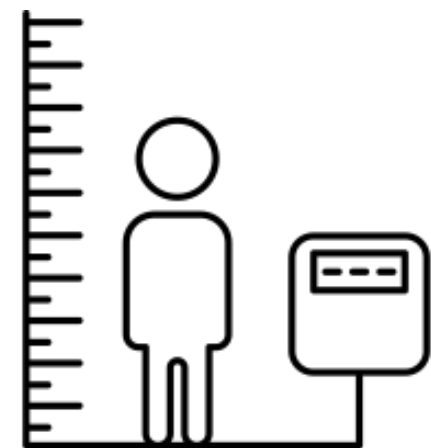
Group Fairness



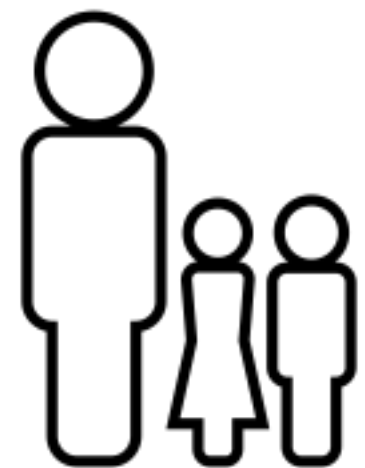
Age



Location



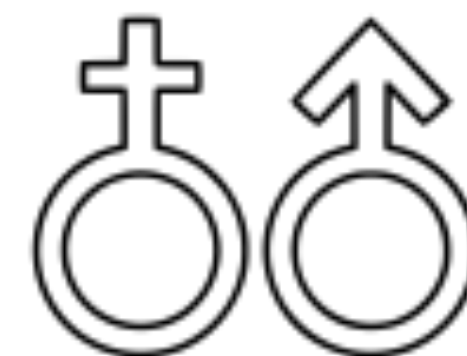
BMI



Number of Children



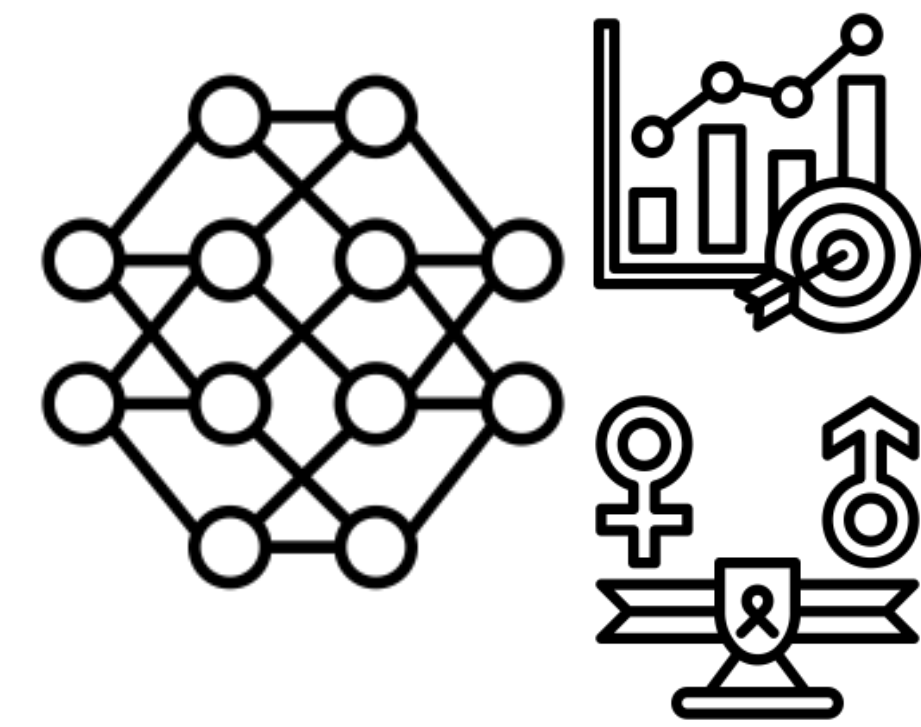
Smoker



Gender



Medical expenses



min Prediction Loss

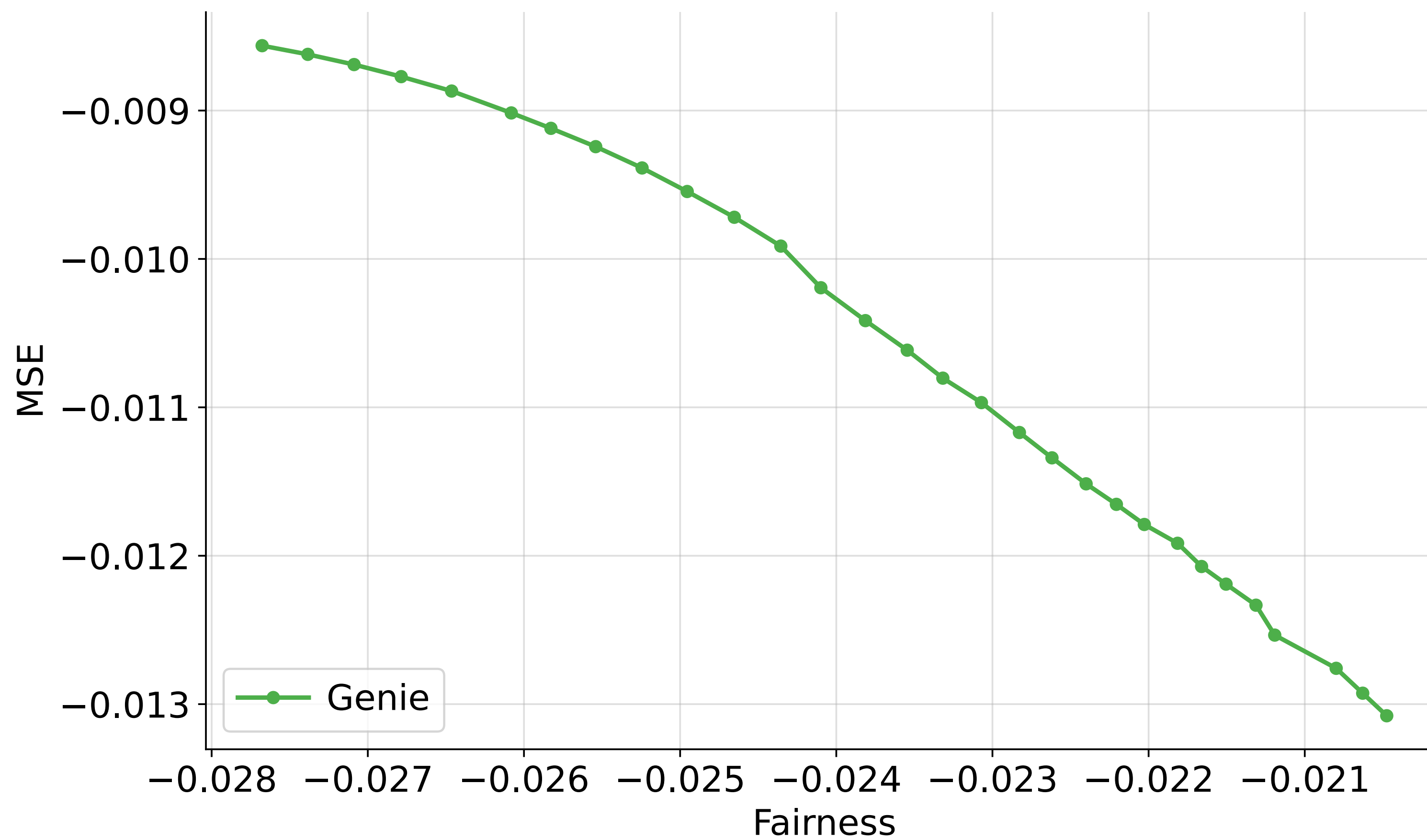
s.t. Fairness Loss $\leq \epsilon$

Independence :— Prediction \perp Gender

Separation :— Prediction \perp Gender | Medical expenses

Sufficiency :— Medical expenses \perp Gender | Prediction

Insurance Dataset



Lagrangian dual

$$\min \text{Prediction Loss} \quad \text{s.t.} \quad \text{Fairness Loss} \leq \epsilon$$



$$\min \max_{\lambda \geq 0} \text{Prediction Loss} + \lambda \underbrace{(\text{Fairness Loss} - \epsilon)}_{\text{Sensitive Attribute}}$$

Uncertainty in Sensitive Attribute

Limited



Additional annotation



Gender

Unreliable



Noisy



Response bias
in a survey



Legal regulations



GDPR
General Data
Protection Regulation



Limited sensitive attributes

Age	Location	BMI	Number of children	Smoker	Medical expenses	Gender	Limited Gender
19	Southwest	27.9	0	Yes	16884	Female	?
28	Southeast	33	3	No	4449	Male	Male
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
62	Southeast	26.29	0	Yes	27808	Female	?

$$\min \max_{\lambda \geq 0} \text{Prediction Loss} + \lambda \underbrace{(\text{Fairness Loss} - \epsilon)}_{\text{Sensitive Attribute}}$$

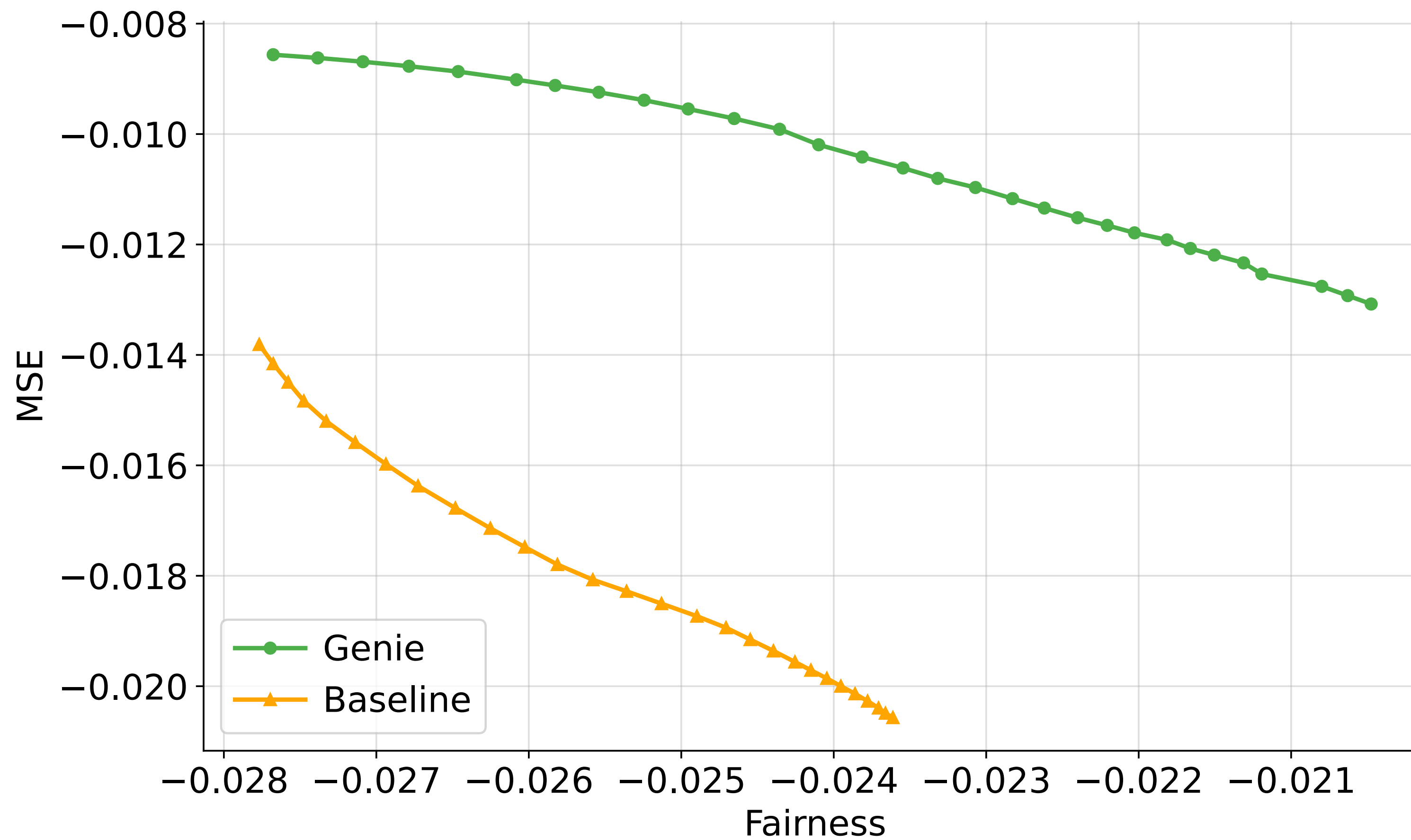
Unreliable sensitive attributes

Age	Location	BMI	Number of children	Smoker	Medical expenses	Gender	Unreliable Gender
19	Southwest	27.9	0	Yes	16884	Female	Male
28	Southeast	33	3	No	4449	Male	Male
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
62	Southeast	26.29	0	Yes	27808	Female	Female

$$\min \max_{\lambda \geq 0} \text{Prediction Loss} + \lambda \underbrace{(\text{Fairness Loss} - \epsilon)}_{\text{Sensitive Attribute}}$$

Insurance

Uncertainty — limited sensitive attribute

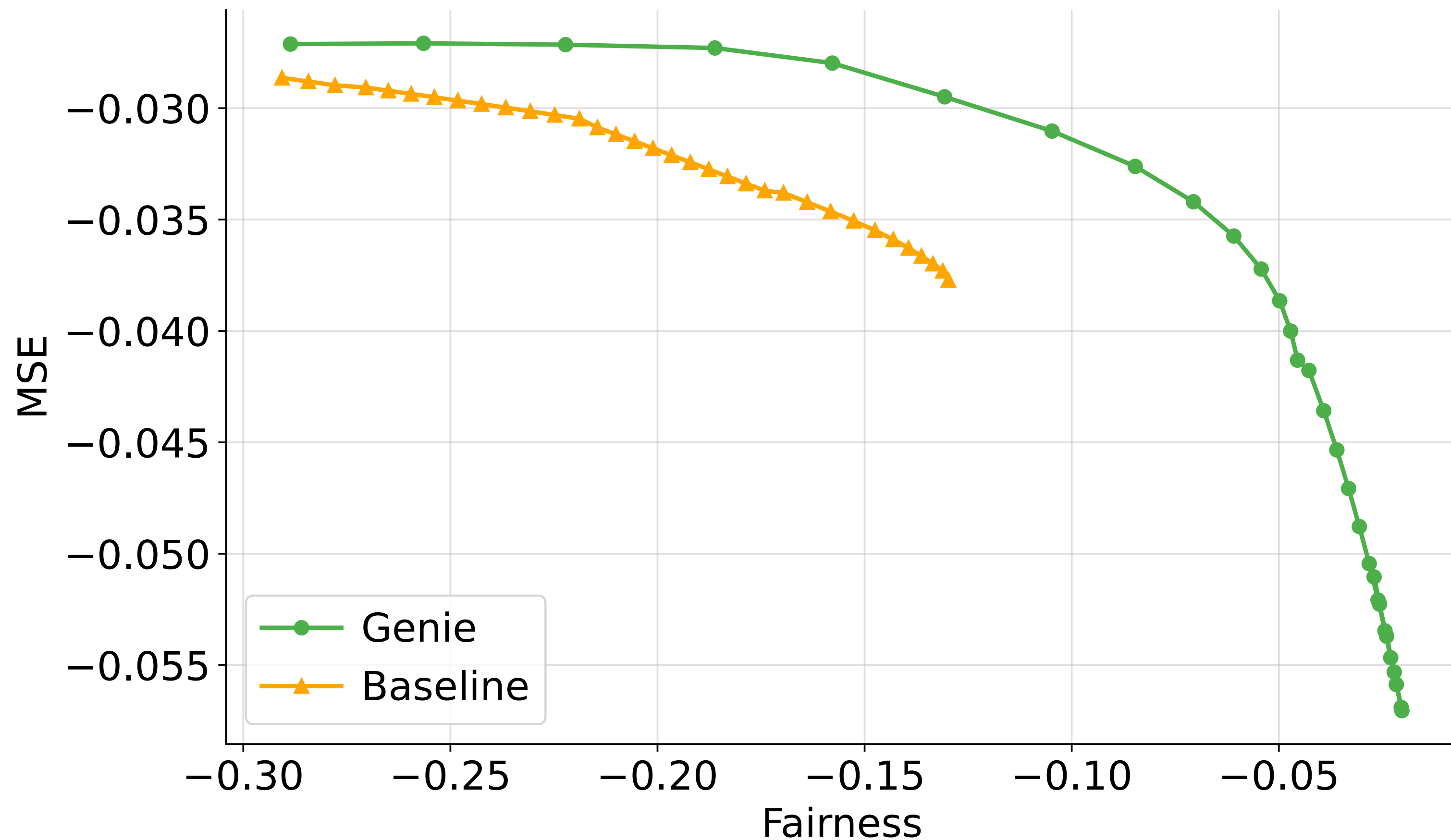


Communities and Crime

Task— predict density of violent crimes (Regression)

Sensitive attribute — Race (Continuous)

Uncertainty — unreliable sensitive attribute



Goal

Learn a fair model despite uncertain sensitive attribute data.

Limited



Additional annotation



Gender

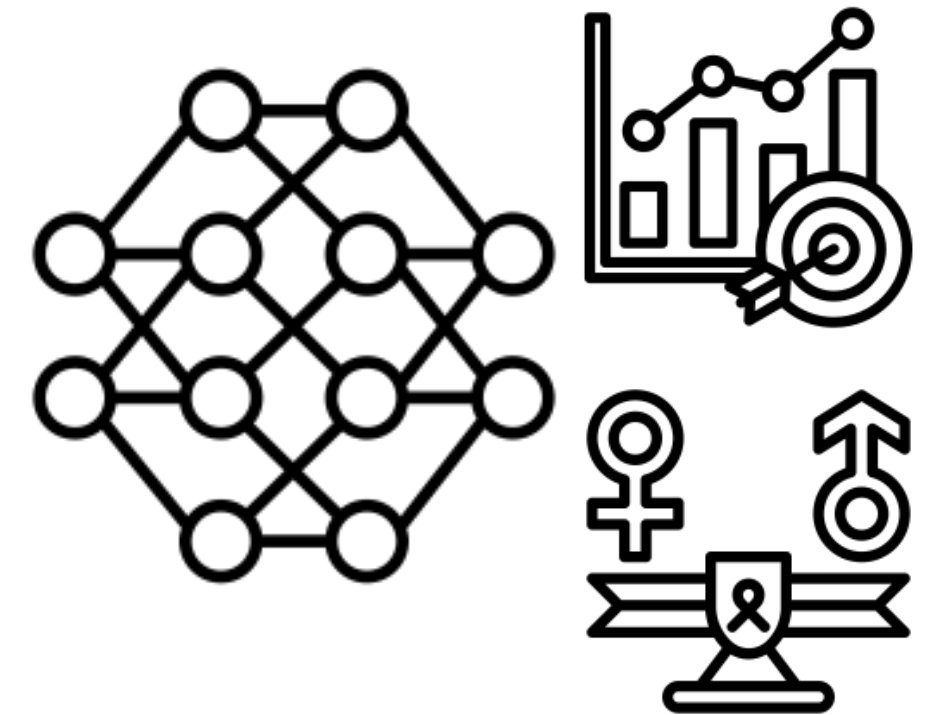


Noisy

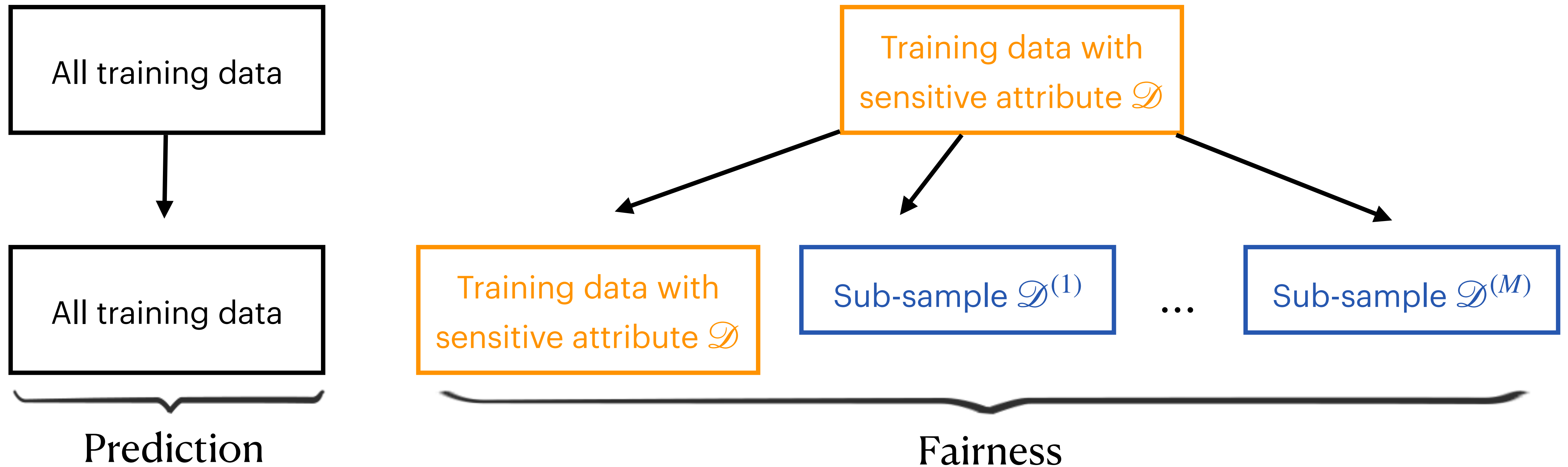
Unreliable



Legal regulations



A General Purpose Algorithm



For some $k \in [n]$ and $M \geq 1$, uniformly draw $\mathcal{D}^{(1)}, \dots, \mathcal{D}^{(M)}$ each of size k from \mathcal{D} with replacement.

$$\min_{\lambda \geq 0} \max \text{Prediction Loss} + \lambda (\text{Fairness Loss}(\mathcal{D}) - \epsilon) +$$

$$\sum_{i \in [M]} \lambda_i (\text{Fairness Loss}(\mathcal{D}^{(i)}) - \epsilon)$$

A General Purpose Algorithm

Bootstrap-M

$$\min_{\lambda \geq 0} \max \text{Prediction Loss} + \lambda (\text{Fairness Loss}(\mathcal{D}) - \epsilon) +$$

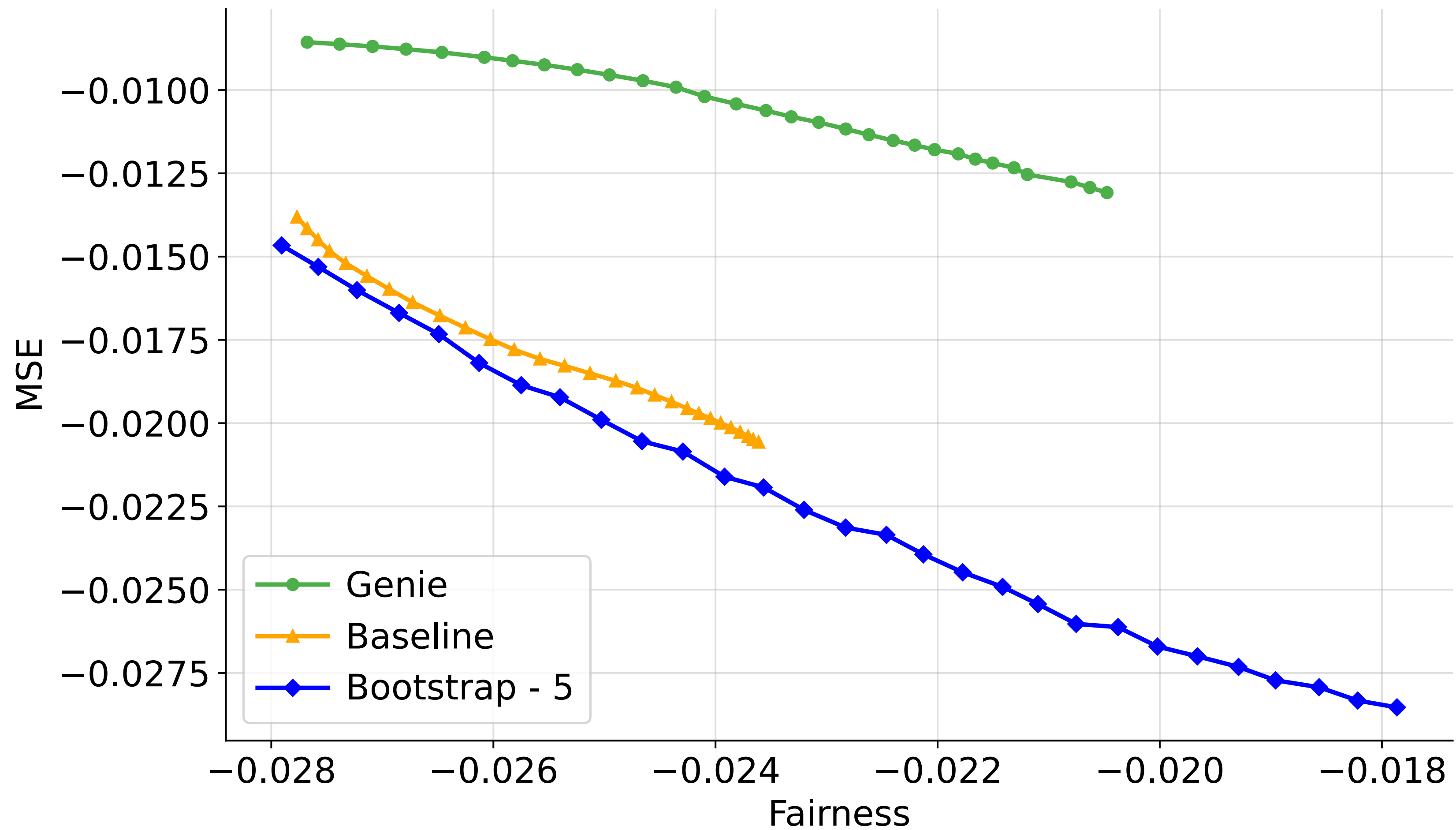
$$\sum_{i \in [M]} \lambda_i (\text{Fairness Loss}(\mathcal{D}^{(i)}) - \epsilon)$$



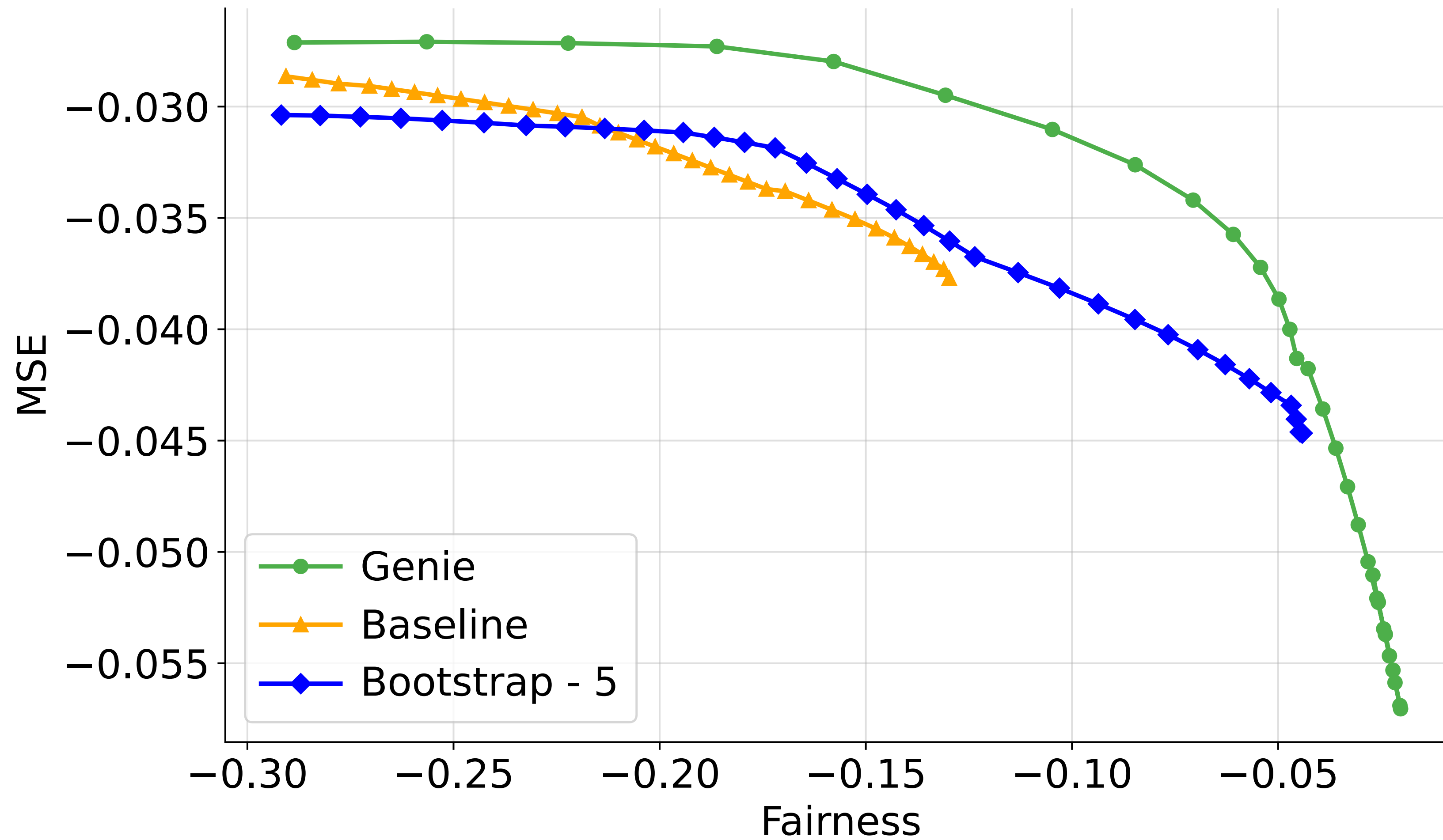
$$\min \text{Prediction Loss} \quad \text{s.t.} \quad \text{Fairness Loss}(\mathcal{D}) \leq \epsilon$$

$$\text{s.t.} \quad \text{Fairness Loss}(\mathcal{D}^{(i)}) \leq \epsilon \text{ for all } i \in [M]$$

Insurance Dataset



Crime Dataset



Setup

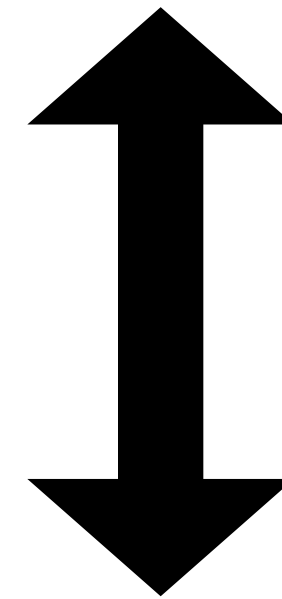
- x :— d -dimensional input
- y :— 1-dimensional target
- e :— 1-dimensional sensitive attribute
- $p_{x,y}$ and p_e :— known

Gaussian Data

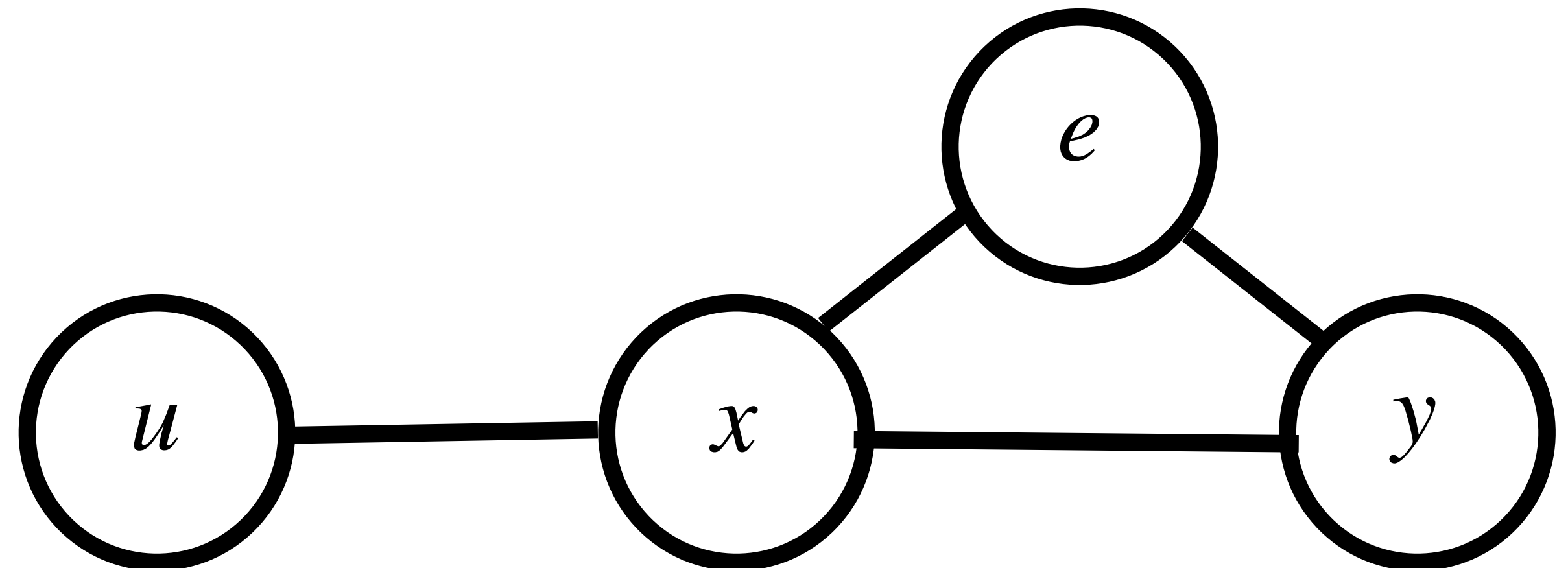
Model the distribution of (x, y, e, u) as Gaussian

min Prediction Loss s.t. Fairness Loss $\leq \epsilon$

Quadratically
Constrained
Quadratic
Program
(QCQP)



$\max_{a \in \mathcal{B}(0,1)} \langle a, b_{yx} \rangle^2$ s.t. $\langle a, b_{ex} \rangle^2 \leq \epsilon$ where $a = b_{ux}$ and $b_{vw} \triangleq \Sigma_v^{-1/2} \Sigma_{vw} \Sigma_w^{-1/2}$



Gaussian Data

Quadratically
Constrained
Quadratic
Program
(QCQP)

Model the distribution of (x, y, e, u) as Gaussian

$$\max_{a \in \mathcal{B}(0,1)} \langle a, b_{yx} \rangle^2 \quad \text{s.t.} \quad \langle a, b_{ex} \rangle^2 \leq \epsilon \quad \text{where} \quad a = b_{ux} \quad \text{and} \quad b_{vw} \triangleq \Sigma_v^{-1/2} \Sigma_{vw} \Sigma_w^{-1/2}$$

Baseline

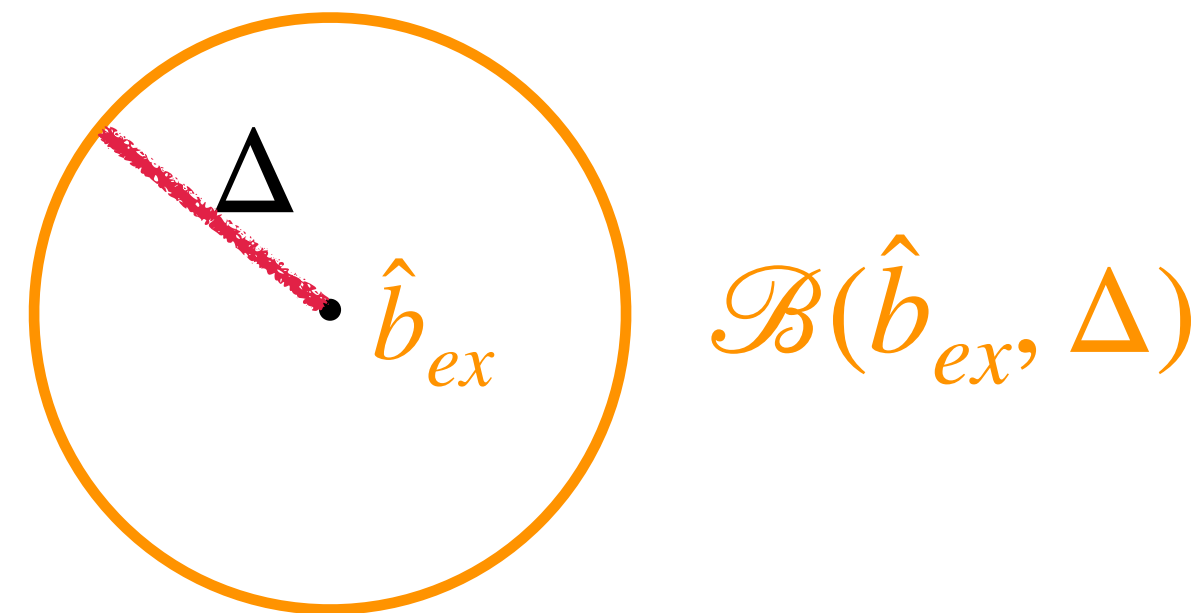
$$\max_{a \in \mathcal{B}(0,1)} \langle a, b_{yx} \rangle^2 \quad \text{s.t.} \quad \langle a, \hat{b}_{ex} \rangle^2 \leq \epsilon$$

This does not guarantee fairness

$$b_{vw} \triangleq \Sigma_v^{-1/2} \Sigma_{vw} \Sigma_w^{-1/2}$$

Robust QCQP

Uncertainty in sensitive attributes

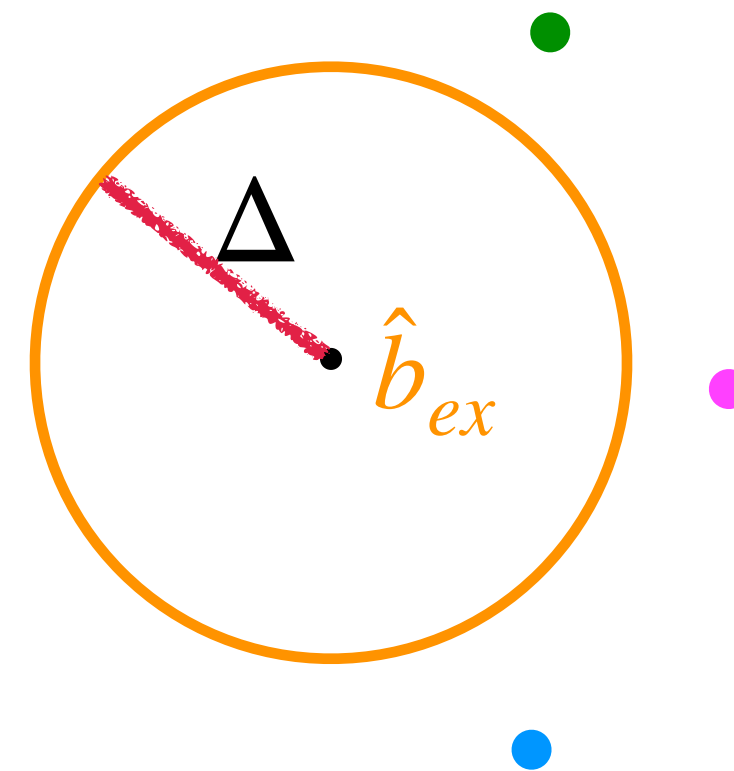


$$\max_{a \in \mathcal{B}(0,1)} \langle a, b_{yx} \rangle^2 \quad \text{s.t.} \quad \langle a, b \rangle^2 \leq \epsilon \quad \text{for all } b \in \mathcal{B}(\hat{b}_{ex}, \Delta)$$

$$b_{vw} \triangleq \Sigma_v^{-1/2} \Sigma_{vw} \Sigma_w^{-1/2}$$

Robust QCQP

Relaxing the uncertainty



$$\max_{a \in \mathcal{B}(0,1)} \langle a, b_{yx} \rangle^2 \quad \text{s.t.} \quad \langle a, b^{(i)} \rangle^2 \leq \epsilon \quad \text{for all } i \in [3]$$