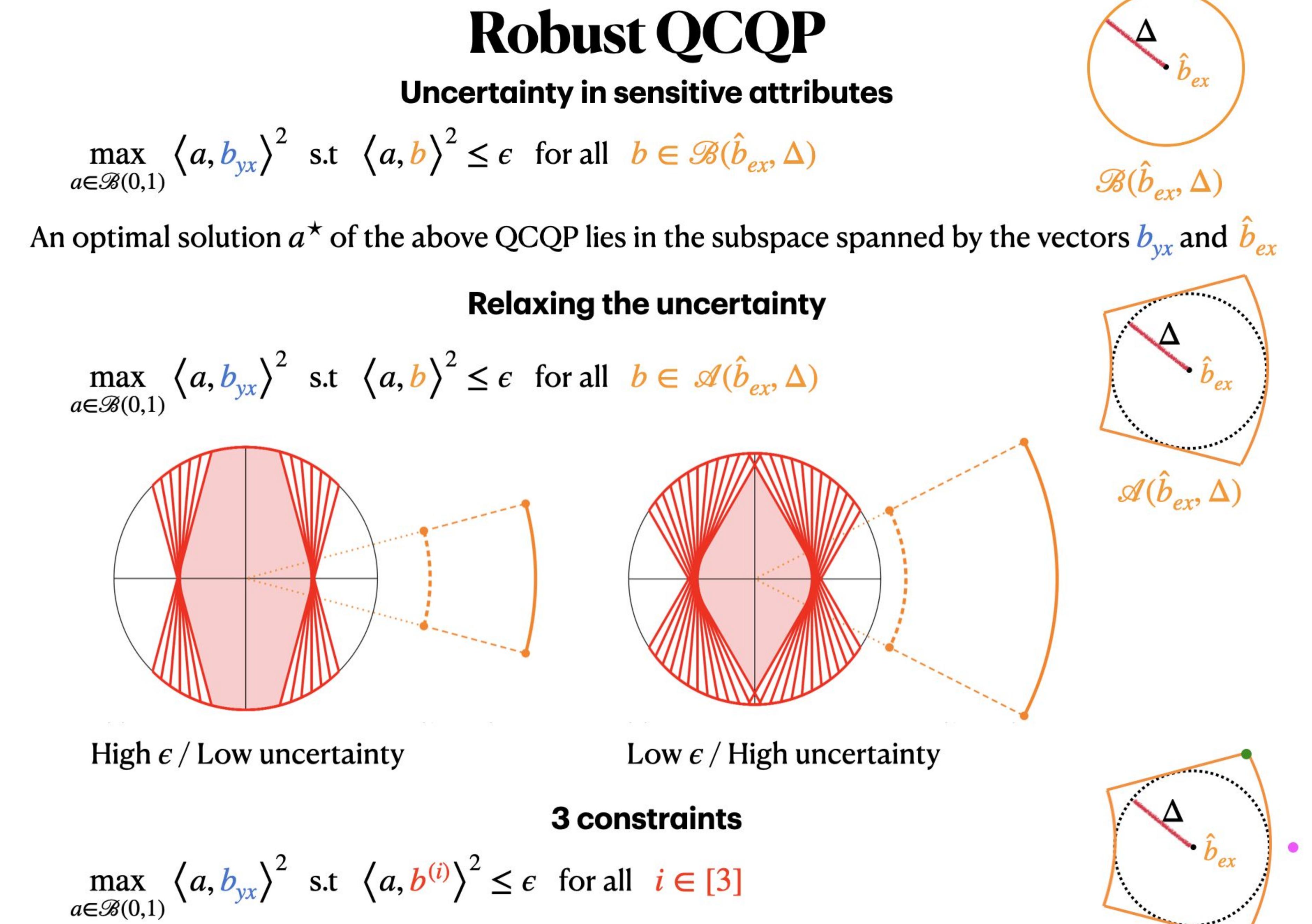
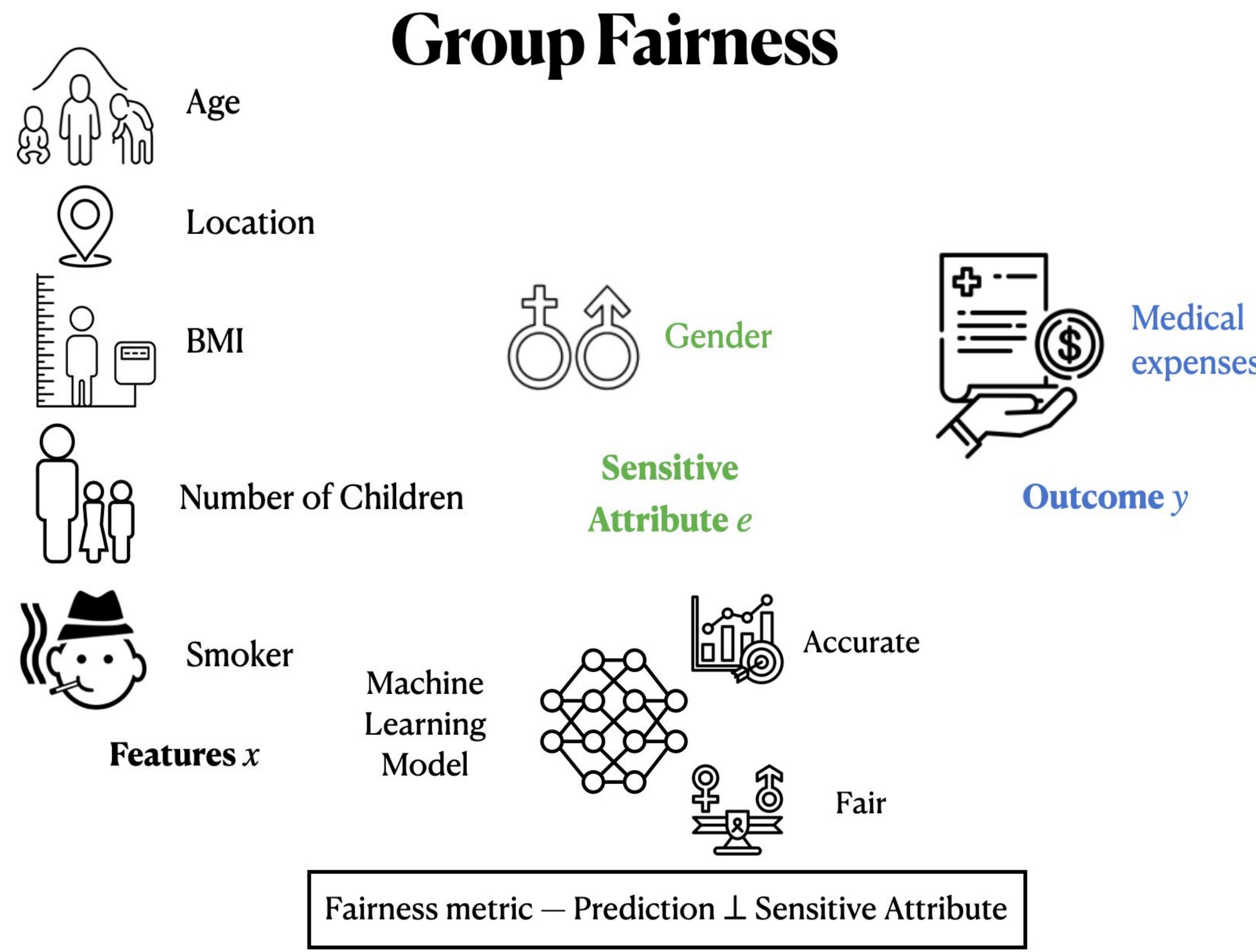
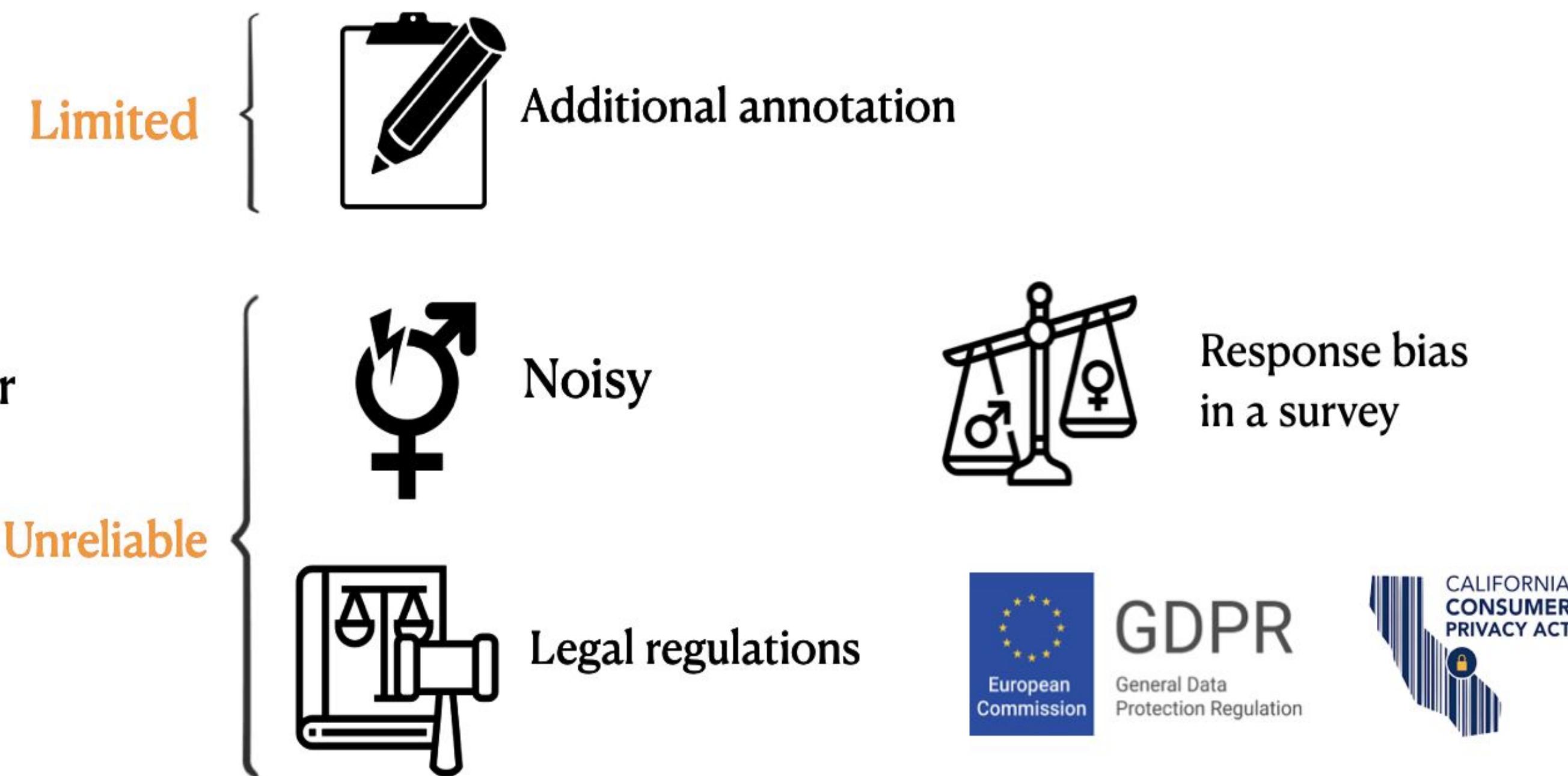


Group Fairness with Uncertain Sensitive Attributes



Uncertainty in Sensitive Attribute



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

Goal — Learn a fair model despite uncertain sensitive attribute data

Limitations of Existing Work

- A. Proxy variables — effectiveness depends on the degree of correlation between the sensitive attribute and the proxy variables
- B. Perturbed sensitive attribute — focus on specific perturbation models

Problem Formulation

Design a randomized mapping $p_{u|x}$ such that the representation u is maximally informative about y while limiting how informative it is about e

$$\min_{p_{u|x}} \min_{\hat{y}(u)} MSE(y, \hat{y}(u)) \text{ s.t. Dependence}(e; u) \leq \epsilon$$

Gaussian Data

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

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

An optimal solution a^* of the above QCQP lies in the subspace spanned by the vectors b_{yx} and b_{ex}

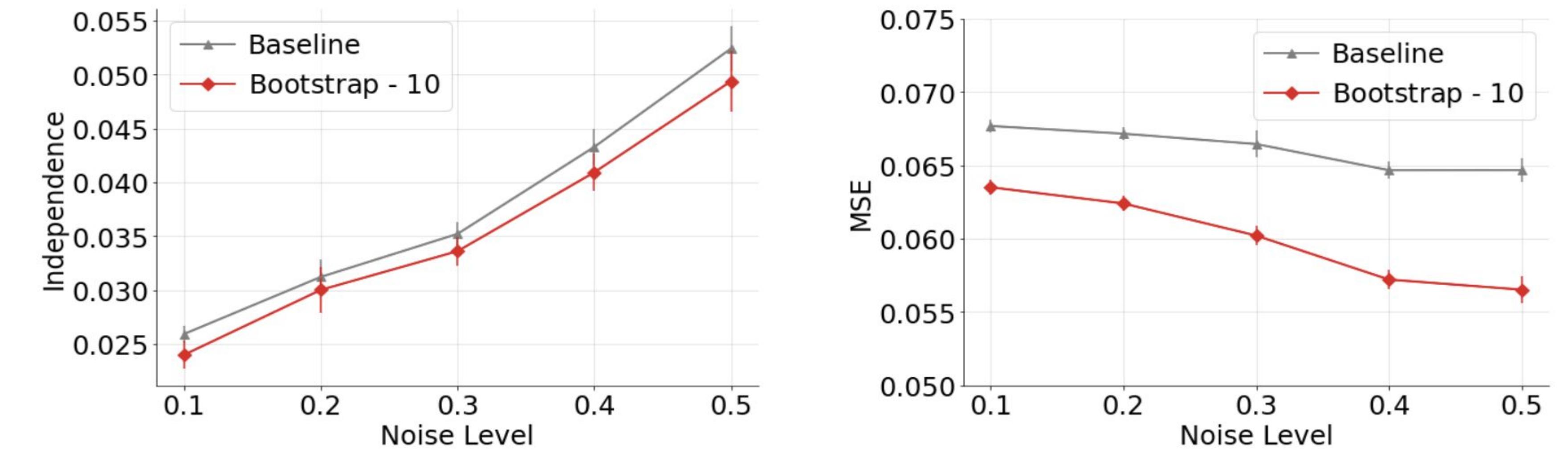
Baseline

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

This does not guarantee fairness

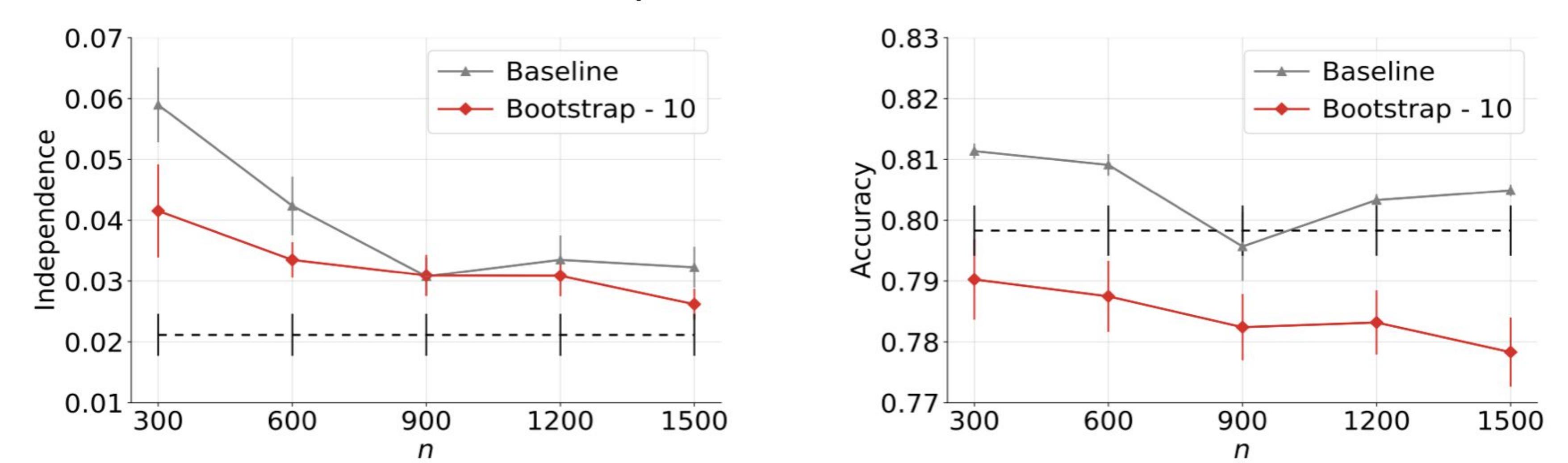
Quadratically Constrained Quadratic Program (QCQP)

Task — predict density of violent crimes (Regression)
Sensitive attribute — Race (Continuous)
Uncertainty — unreliable sensitive attribute



Adult Dataset

Task — predict whether an individual's income > \$50k (Classification)
Sensitive attribute — Gender (Discrete)
Uncertainty — limited sensitive attribute



Crime Dataset

Task — predict density of violent crimes (Regression)
Sensitive attribute — Race (Continuous)
Uncertainty — unreliable sensitive attribute

