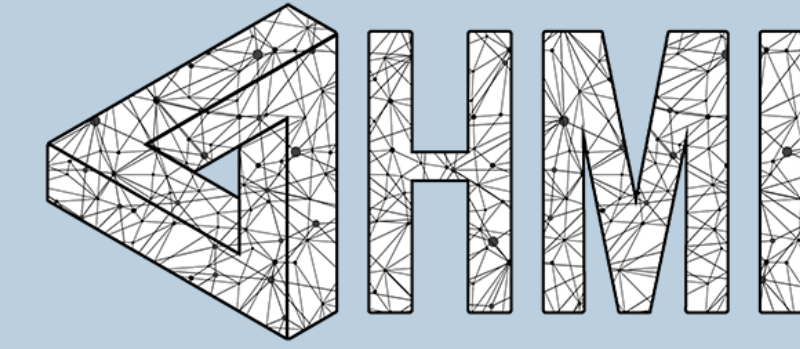


Fair Densities via Boosting the Sufficient Statistics of Exponential Families



Australian National University



amazon | science

Google Research

Summary

We propose a boosting algorithm for fair pre-processing of datasets with strong theoretical guarantees.

↔ two modes of fairness ↔ use in continuous domain ↔ boosting-style convergence | Others in paper: ↔ characterization of fair densities ↔ update interpretability

Fair Pre-Processing Setting

- ▶ \mathcal{X} domain of observations (can be continuous);
- ▶ \mathcal{Y} domain of labels (finite);
- ▶ \mathcal{S} domain of sensitive attributes (finite, separate from \mathcal{X}); such as age, race, etc.;
- ▶ P an input distribution which is unfair.

Goal: ① find a fair density Q ; but ② also close to P (e.g. KL).

Definitions of Data Fairness

Given a density $P \in \mathcal{D}(\mathcal{X} \times \mathcal{Y} \times \mathcal{S})$, we consider two forms of data fairness. The following properties hold, if $\forall s, s' \in \mathcal{S}$:

ρ -Representation Rate:

$$P(s)/P(s') \geq \rho.$$

Representation across \mathcal{S} .

τ -Statistical Rate:

$$P(y | s)/P(y | s') \geq \tau$$

Positive outcomes across \mathcal{S} .

Denote these rates of a density as $RR(P)$ and $SR(P)$, resp.

Fair Boosted Density Estimation

Start with an initially fair (but not accurate) density Q_0 . Let $RR_0 := RR(Q_0)$ and $SR_0 := SR(Q_0)$.

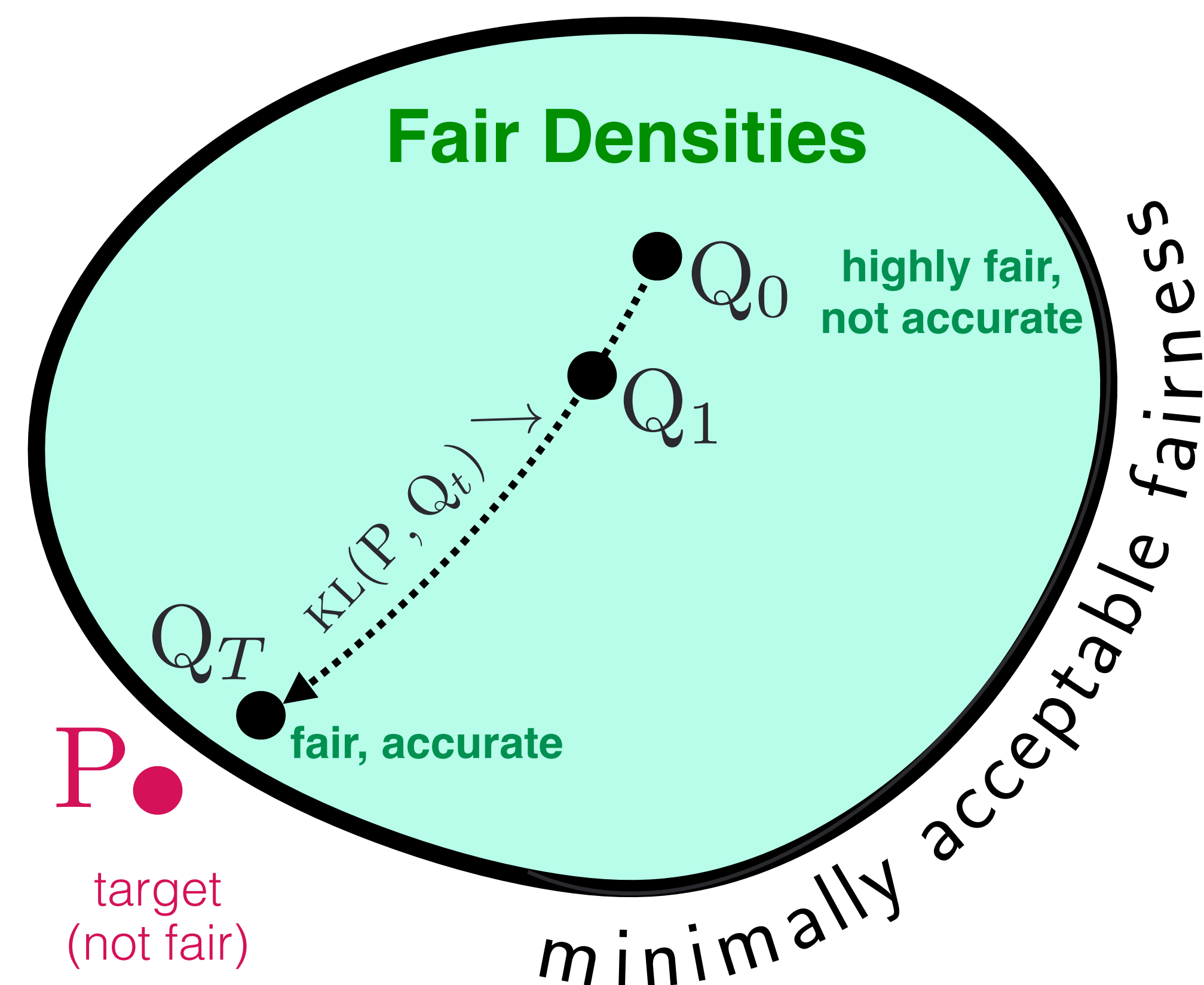
$$Q_t(x, y, s) \stackrel{\text{Update}}{\propto} Q_{t-1}(x, y, s) \cdot \exp(\vartheta_t c_t(x, y, s))$$

Weak Learner $c_t(x, y, s)$:

Classifies the 'realness' of samples of P vs Q_{t-1} .

Leverage Coefficient ϑ_t :

Controls the 'step size' of each update.



Fairness Guarantees [Abridged]

Given that Q_0 has high SR_0 , we present two modes of fairness by changing the leverage ϑ_t . Let $\tau \in (0, 1]$ be a 'target fairness'. Then,

Exact Fairness:

$$\vartheta_t^E := -(C2^{t+1})^{-1} \log(\tau/SR_0)$$

Theorem 3.2:

Suppose $\vartheta_t = \vartheta_t^E$ hold, then

$$SR(Q_T) > \tau.$$

Holds when we forever-boost.

Relative Fairness:

$$\vartheta_t^R := -(4Ct)^{-1} \log(\tau/SR_0)$$

Theorem 3.3:

Suppose $\vartheta_t = \vartheta_t^R$ hold, then

$$SR(Q_T) > \tau \cdot \tau^{\log T}.$$

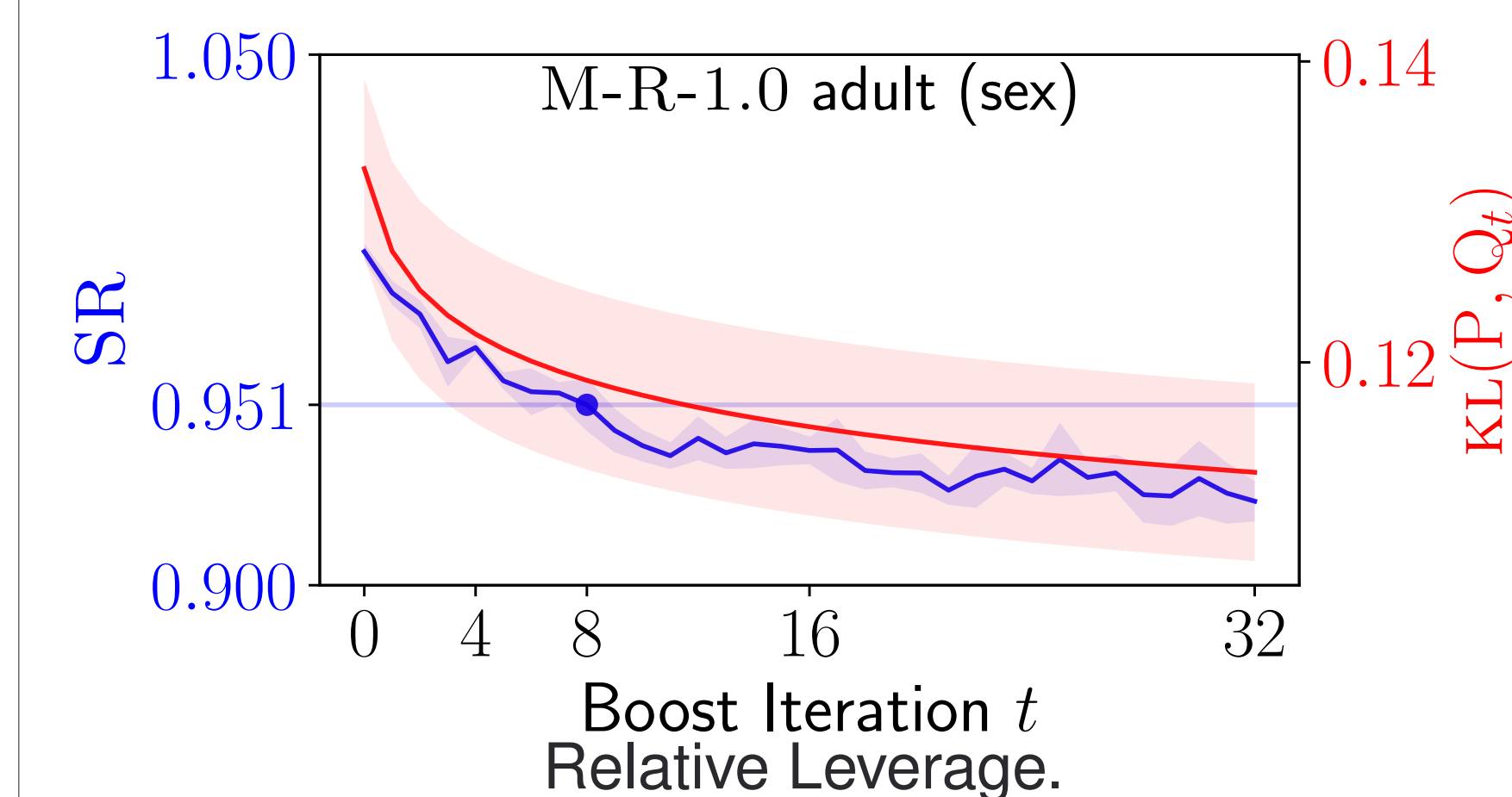
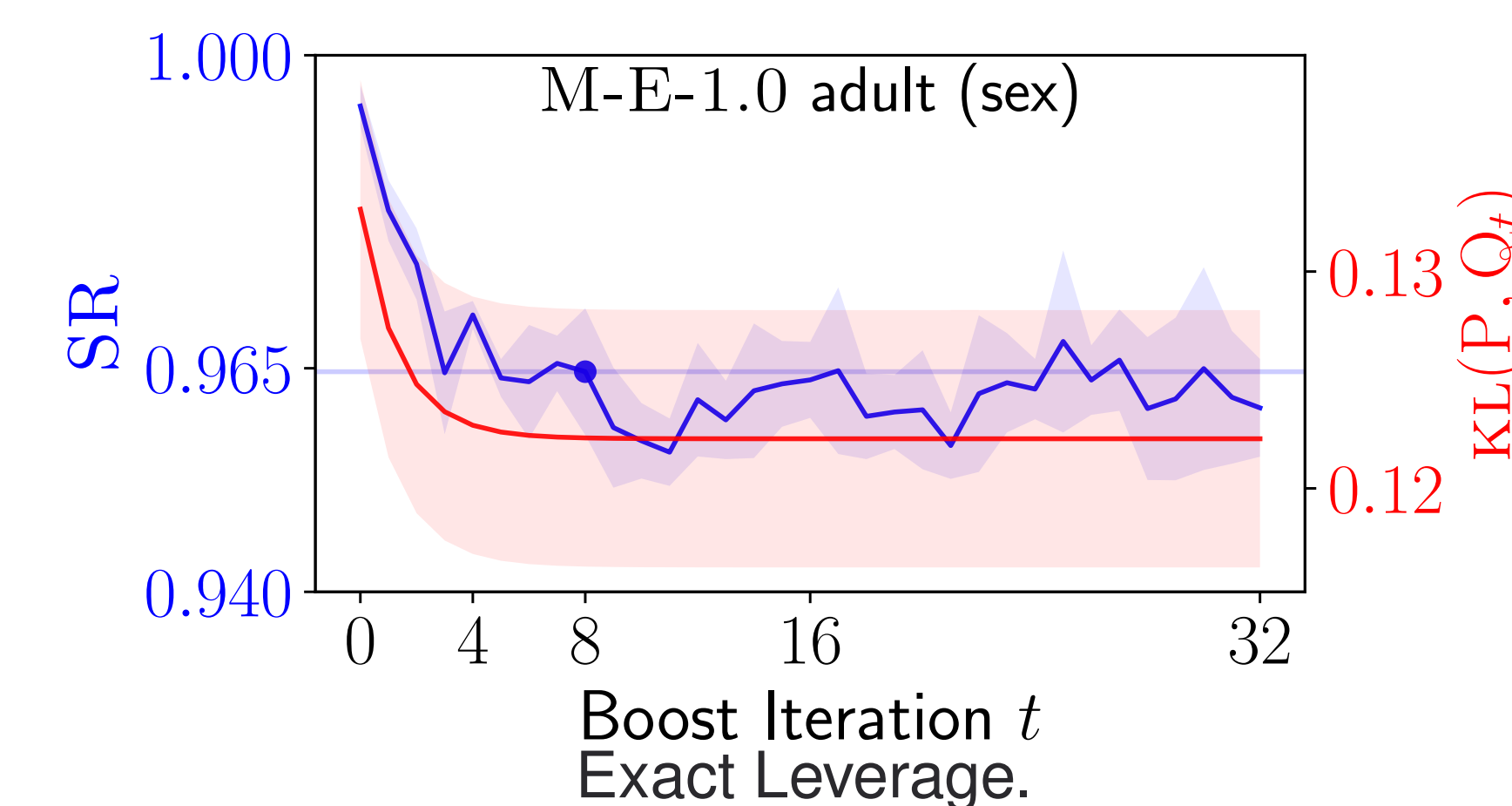
Decays with boosting iterations.

Convergence [Informal]

Theorem 3.5: If the weak learners c_t satisfies a *Weak Learning Assumption* (in predicting samples of P against Q_{t-1}), then

$$KL(P, Q_{t-1}) - KL(P, Q_t) \geq \vartheta_t \cdot \Lambda_t, \quad \text{RHS} \geq 0 \text{ with light assumptions.}$$

Using $RR_0 = 1$, $SR_0 = 1$, and $\tau = 0.8$:



	DATA	M-E-1.0	M-R-1.0	MAXENT	TABFAIR	
DATA	RR	.496 ± .002	.958 ± .003	.919 ± .004	.995 ± .002	.516 ± .023
	SR	.360 ± .005	.961 ± .005	.924 ± .006	.979 ± .004	.862 ± .101
	KL	—	.122 ± .006	.113 ± .006	.182 ± .005	1.68 ± .538
PRED	SR _c	.360 ± .003	.818 ± .010	.793 ± .011	.919 ± .011	.823 ± .118
	EO	.471 ± .008	.959 ± .016	.935 ± .016	.981 ± .010	.867 ± .105
	ACC	.803 ± .003	.785 ± .002	.787 ± .002	.773 ± .005	.781 ± .006

Table: Data and prediction results (Adult Dataset)

MAXENT: "Data preprocessing to mitigate bias: A maximum entropy based approach. 2020."
TABFAIR: "TabFairGAN: Fair tabular data generation with generative adversarial networks. 2022."

