



























· Suppose 1/2 = EF = EM











































even absent fairness



non-fair heuristics...

































(a) convergence plot: communities data set

(b) error spread plot: communities data set

Figure 2: (a) Error-unfairness trajectory plots illustrating the convergence of algorithm **AIF-Learn**. (b) Errorunfairness tradeoffs and individual errors for **AIF-Learn** vs. simple mixtures of the error-optimal model and random classification. Gray dots are shifted upwards slightly to avoid occlusions.



· Elicitation: "x x x' chauld receiv

-'x #x' should receive " same outcome"





constraints

5.2 Subjective Fairness Elicitation

In you	r view, a	as a matter of fa	airness, should the fol	lowing two individ	duals recieve the sa	me recidivism	prediction, or is	s it ok to	give them different p
sex	age	race	juv. felony count	juv. misdemea	nor count juv. o	ther count	priors count	sever	rity of charge
Male	25	Caucasian	0		1	0	6	Felon	у
vs.									
sex	age	race	juv. felony co	ount juv. miso	demeanor count	juv. other cou	unt priors o	ount	severity of charge
Male	29	African-Amer	rican	0	0		1	10	Felony
Should	bo trop		Ok to troat differently						

Figure 1: Screenshot of sample subjective fairness elicitation question posed to human subjects.





(b)

(a)

Other Strengthenings



















· Fair RL/Contol











The Real World

· Open-ended AI services





