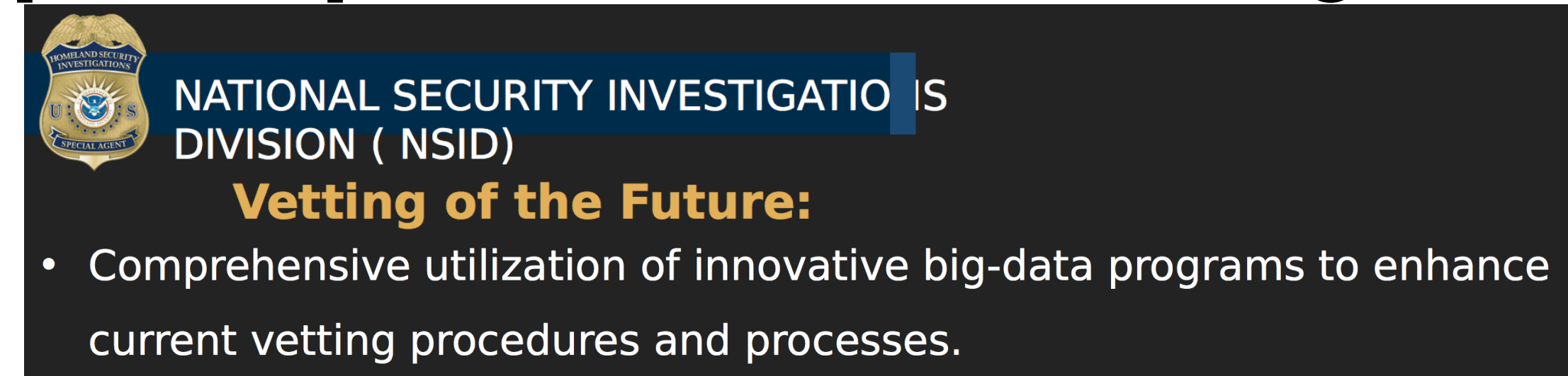


## ML is making Life-Changing Decisions Trump's Proposed Extreme Vetting Software



[Biddle, 2017]

### We Have Big Problems



### Our Solution

A fair classifier gives the same prediction had the person had a different race/sex.

We present a metric to check if any algorithm is fair

And an algorithm to learn fair classifiers

## Law School Admissions

**Data**      **Goal**

$\hat{Y}: \mathcal{X}, A \rightarrow Y$

sensitive attributes A	features X	label Y
male white	thumbs up	thumbs up
female black	thumbs up	thumbs down
male black	thumbs down	thumbs down

**Prior Fairness Methods**

- Fairness through Unawareness [Grgić-Hlača et al., 2016]
- Equality of Opportunity [Hardt et al., 2016]

## Prior Fairness Definitions Are Insufficient Unfair Influences

Fairness through Unawareness [Grgić-Hlača et al., 2016]

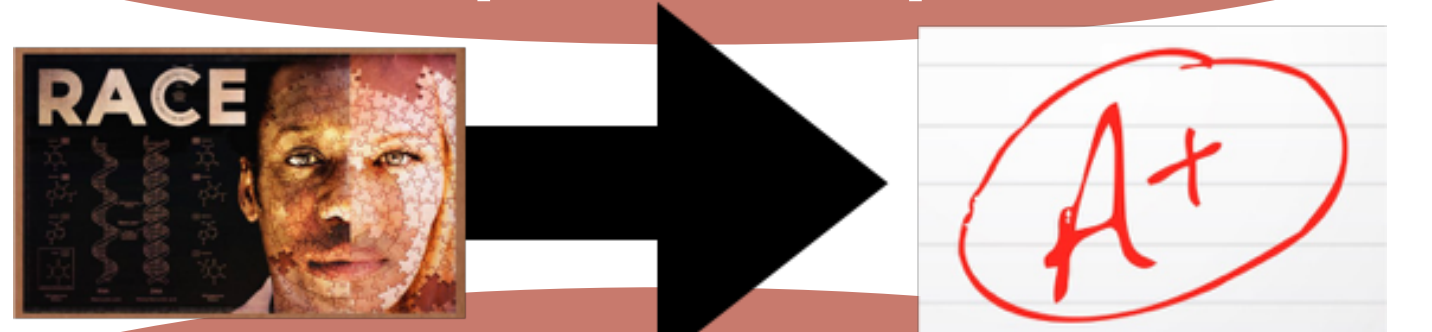
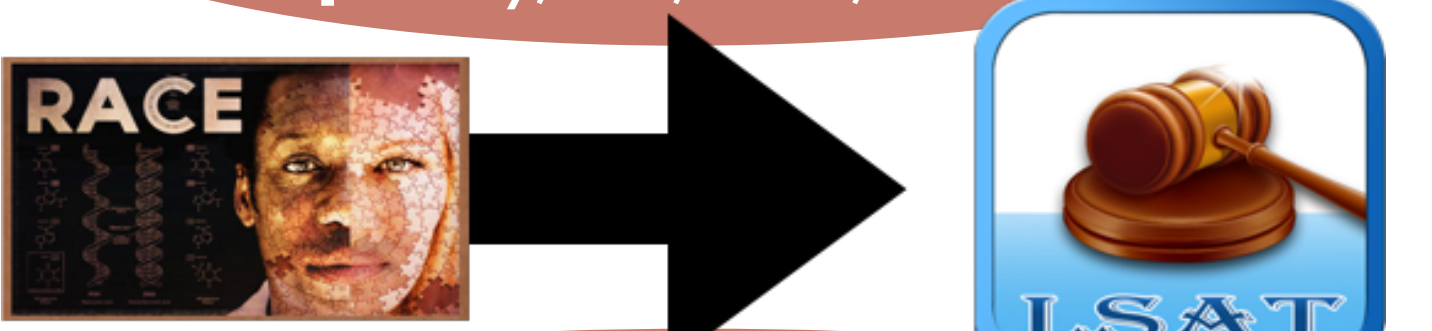
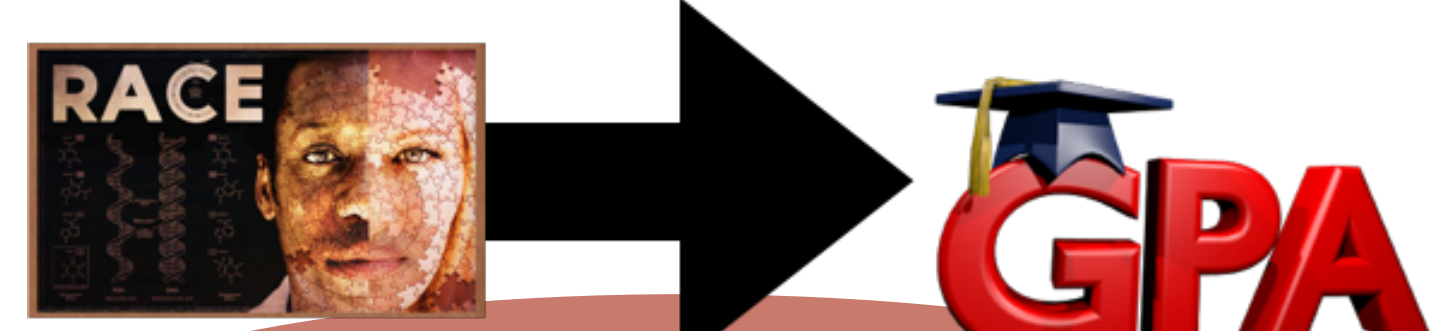
remove sensitive attributes: ~~gender~~ ~~race~~

$\hat{Y}: \mathcal{X} \rightarrow Y$

features X are biased

Equality of Opportunity [Hardt et al., 2016]

$\mathbb{P}(\hat{Y} = 1 | A = a, Y = 1) = \mathbb{P}(\hat{Y} = 1 | A = a', Y = 1)$

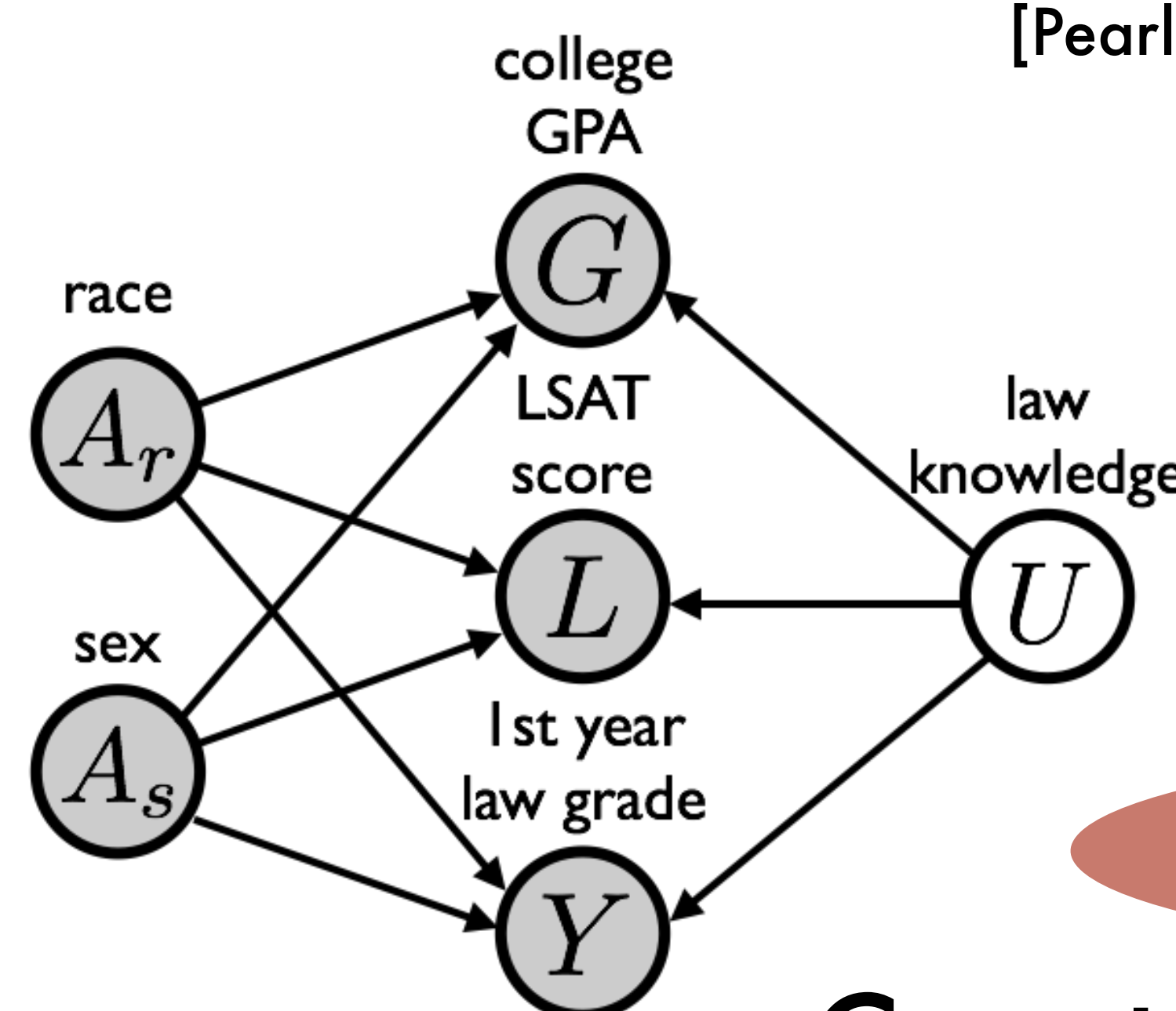


35% of black students  
50% of white students

label Y is biased

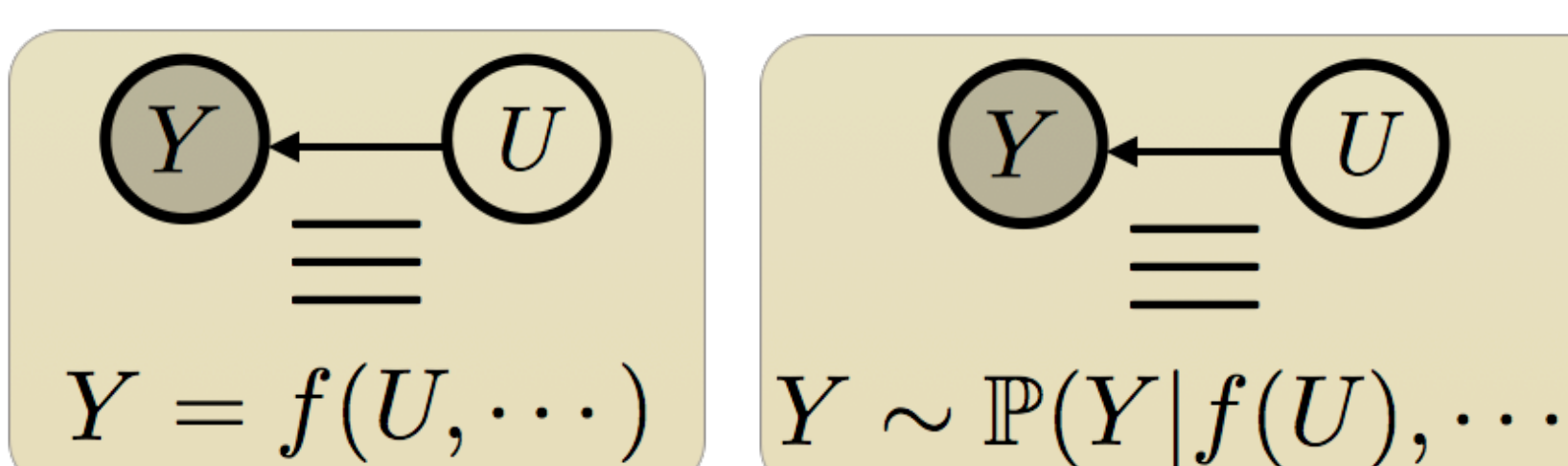
## Causal Inference

[Pearl et al., 2016]



### structural causal models

deterministic      non-deterministic



how can we make fair predictions?

## Counterfactuals

Given an individual:					
$A_s$	$A_r$	G	L	Y	
male	black	thumbs down	thumbs up	thumbs down	
1. Change race attribute					
$A_s$	$a'$	G	L	Y	
male	white	thumbs down	thumbs up	thumbs down	
2. Compute unobserved variables in causal model					
$A_s$	$a'$	G	L	Y	U
male	white	thumbs down	thumbs up	thumbs down	thumbs up
3. Recompute observed variables in causal model					
$A_s$	$a'$	$G_{A_r \leftarrow a'}$	$L_{A_r \leftarrow a'}$	$Y_{A_r \leftarrow a'}$	
male	white	thumbs up	thumbs up	thumbs up	

## Counterfactual Fairness

Definition. A predictor  $\hat{Y}$  is counterfactually fair if given observations  $\mathcal{X} = \mathbf{x}$  and  $A = a$  we have that,

$$\mathbb{P}(\hat{Y}_{A \leftarrow a} = y | \mathcal{X} = \mathbf{x}, A = a) = \mathbb{P}(\hat{Y}_{A \leftarrow a'} = y | \mathcal{X} = \mathbf{x}, A = a)$$

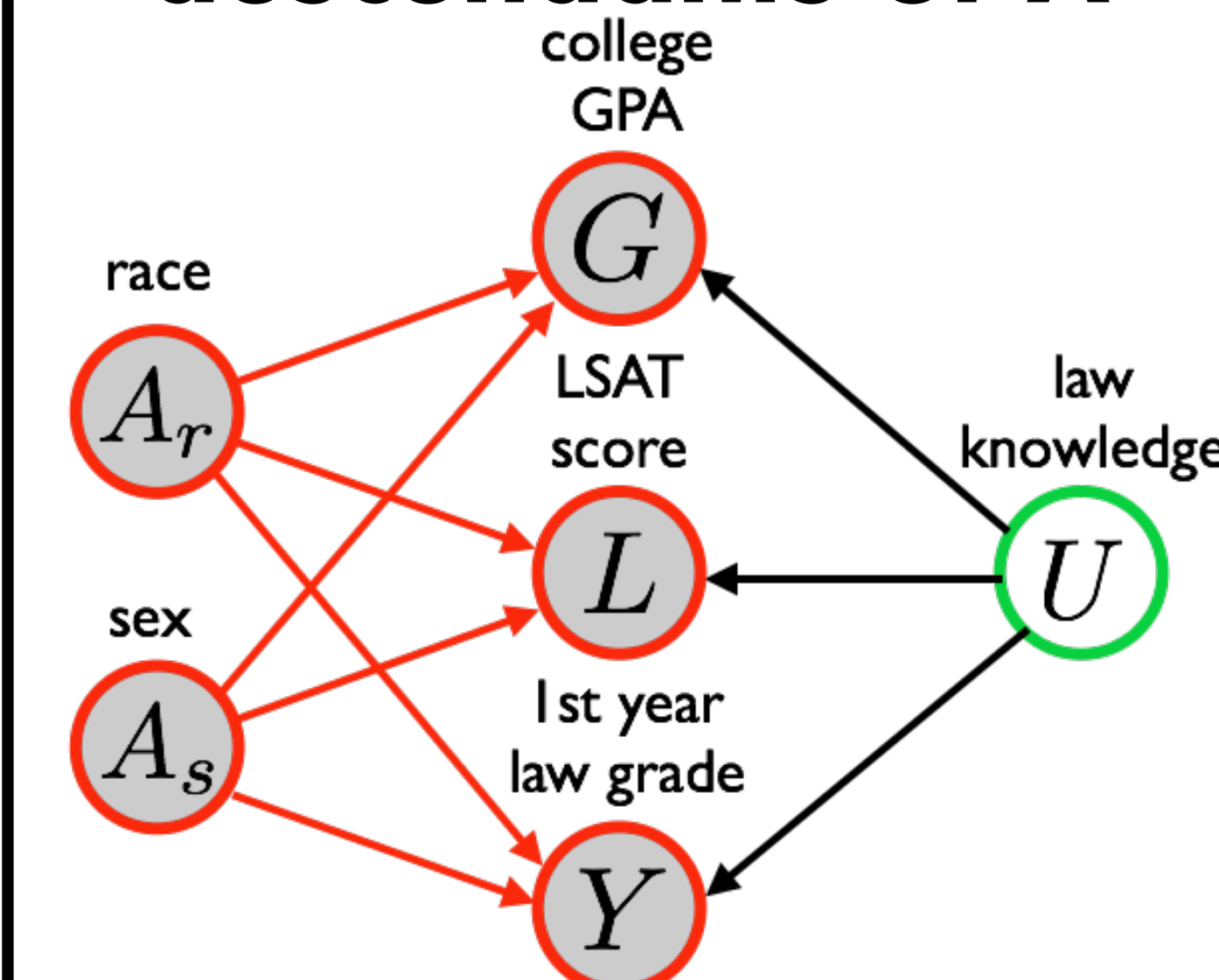
for all  $y$  and  $a' \neq a$ .

Compares the same individual with a different version of him/herself

## Learning Fair Predictors

Counterfactuals alter descendants of A

Algorithm (non-deterministic models)



Given:  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)}, a^{(i)})\}_{i=1}^d$

- For each data point  $i \in \mathcal{D}$ , sample  $u_1^{(i)}, \dots, u_m^{(i)} \sim \mathbb{P}(U|\mathbf{x}^{(i)}, a^{(i)})$
- $\hat{\theta} \leftarrow \arg \min_{\theta} \sum_{i \in \mathcal{D}} \sum_{j=1}^m \ell(y^{(i)}, \hat{Y}_{\theta}(u_j^{(i)}, \mathbf{x}_{\setminus A}^{(i)}))$
- Return:  $\hat{Y}_{\hat{\theta}}$

features that are not descendants of A

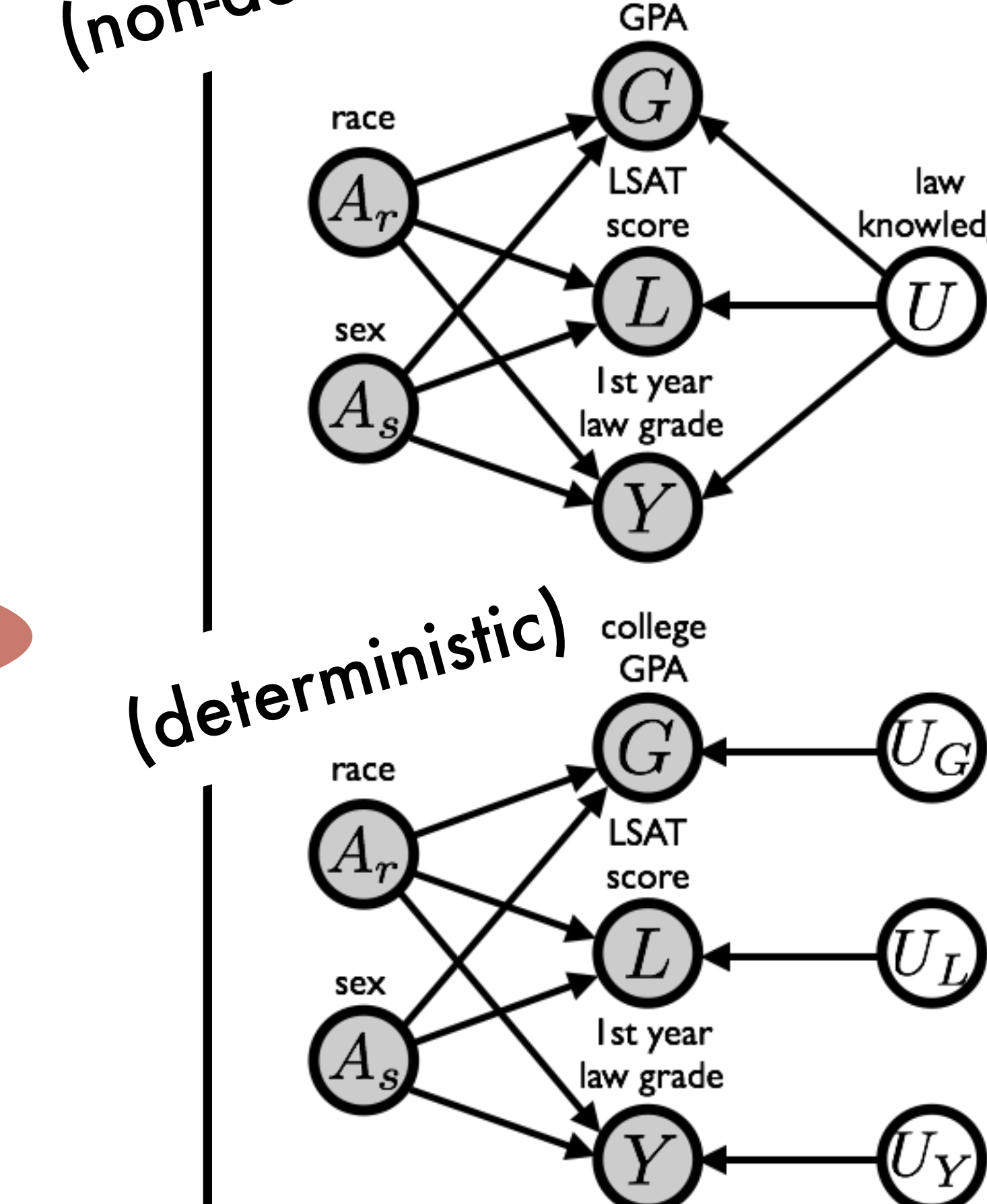
[Kilbertus et al., 2017]

## Related Work

$$\mathbb{P}(\hat{Y} = y | do(A = a), \mathcal{X} = \mathbf{x}) = \mathbb{P}(\hat{Y} = y | do(A = a'), \mathcal{X} = \mathbf{x})$$

Compares different individuals with the same observed features

## Results: US law schools causal models



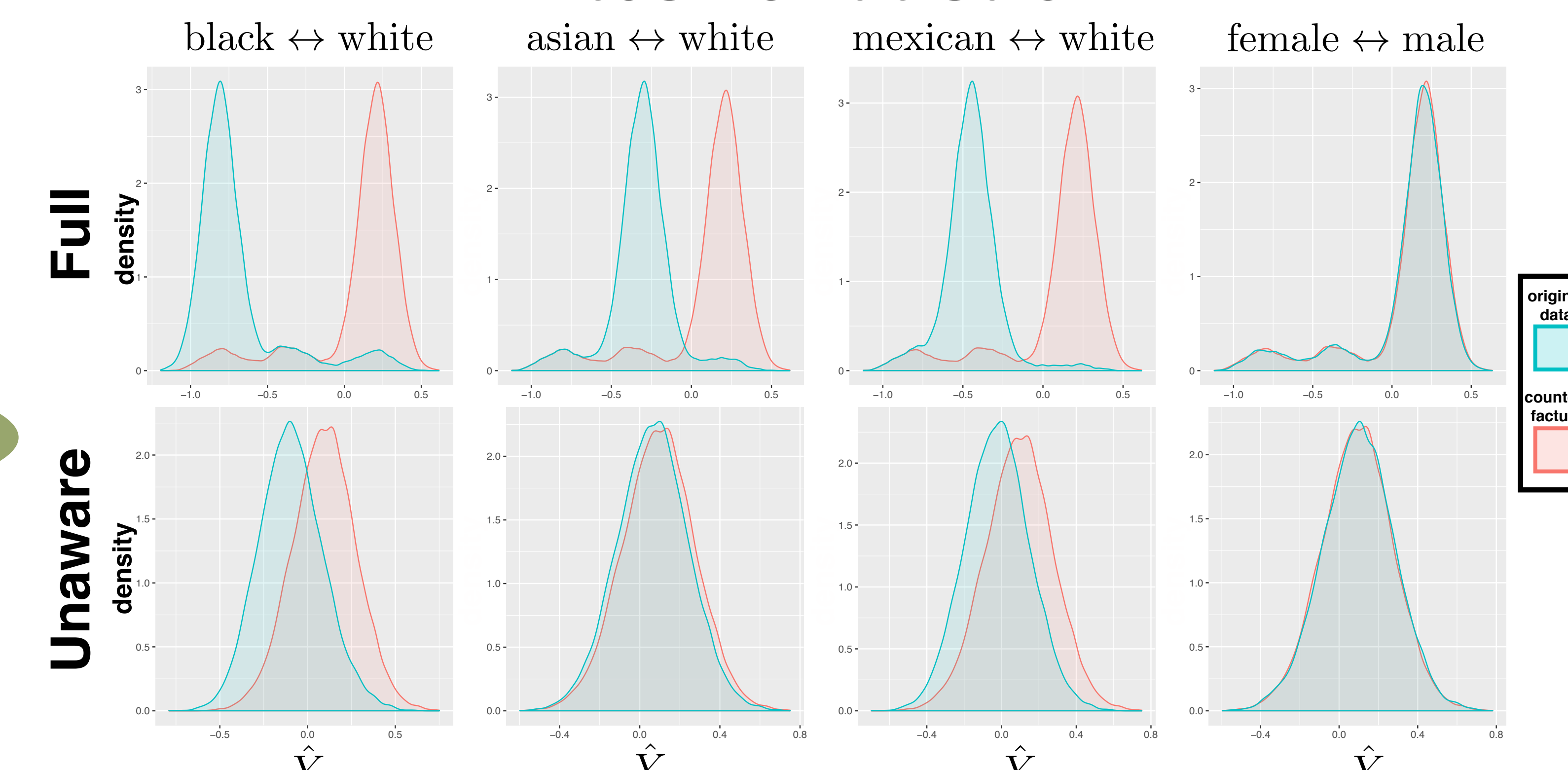
$G \sim \text{Normal}(b_G + [U, A_r, A_s]w_G, \sigma_G)$   
 $L \sim \text{Poisson}(\exp(b_L + [U, A_r, A_s]w_L))$   
 $Y \sim \text{Normal}([U, A_r, A_s]w_Y, 1)$   
 $U \sim \text{Normal}(0, 1)$

$G = b_G + [A_r, A_s]w_G + U_G$   
 $L = b_L + [A_r, A_s]w_L + U_L$   
 $Y = b_Y + [A_r, A_s]w_Y + U_Y$   
 $U_G \sim \mathbb{P}(U_G)$   
 $U_L \sim \mathbb{P}(U_L)$   
 $U_Y \sim \mathbb{P}(U_Y)$

### predictive error

	Full	Unaware	C-Fair (Non-Det.)	C-Fair (Det.)
RMSE	0.873	0.894	0.929	0.918

### counterfactuals



## Results: NYC stop-and-frisk

