
FAIRNESS

Matt Kusner

ML IS AMAZING

Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification

Kaiming He

Xiangyu Zhang

Shaoqing Ren

Jian Sun

Microsoft Research

Human-level control through deep reinforcement learning

Volodymyr Mnih^{1*}, Koray Kavukcuoglu^{1*}, David Silver^{1*}, Andrei A. Rusu¹, Joel Veness¹, Marc G. Bellemare¹, Alex Graves¹, Martin Riedmiller¹, Andreas K. Fidjeland¹, Georg Ostrovski¹, Stig Petersen¹, Charles Beattie¹, Amir Sadik¹, Ioannis Antonoglou¹, Helen King¹, Dharshan Kumaran¹, Daan Wierstra¹, Shane Legg¹ & Demis Hassabis¹

Letter

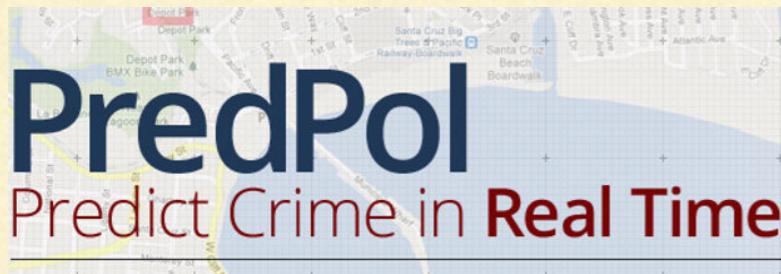
Dermatologist-level classification of skin cancer with deep neural networks

Andre Esteva ✉, Brett Kuprel ✉, Roberto A. Novoa ✉, Justin Ko, Susan M. Swetter, Helen M. Blau & Sebastian Thrun ✉

WHY NOT USE IT EVERYWHERE?

Policing

[Ensign et al., 2017]



Parole Sentencing

[Larson et al., 2016]

COMPAS

Advertising

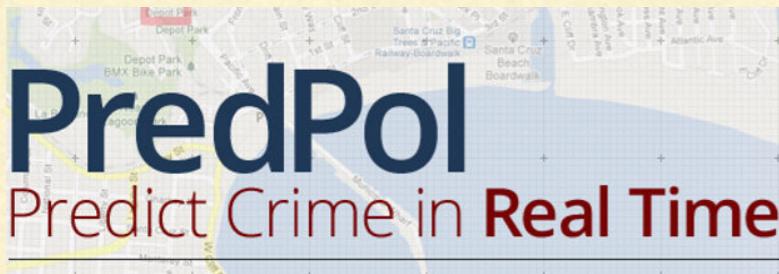
[Sweeney, 2013]

Google

WHY NOT USE IT EVERYWHERE?

Policing

[Ensign et al., 2017]



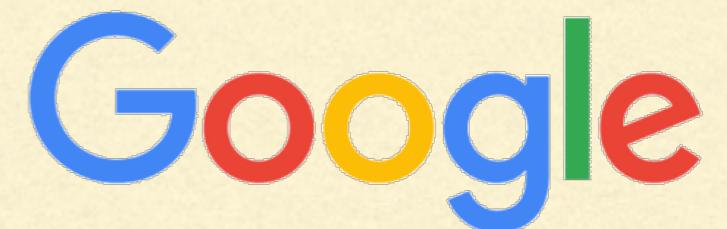
Parole Sentencing

[Larson et al., 2016]

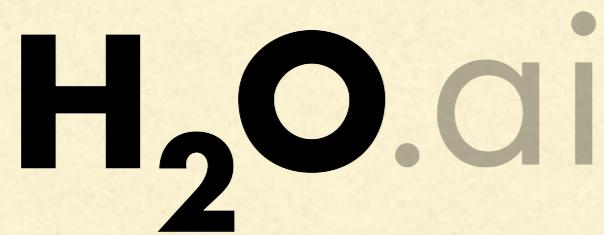


Advertising

[Sweeney, 2013]



Insurance



Lending



Hiring



ML CAN BE RACIST...

BBC [Sign in](#) [News](#) [Sport](#) [Weather](#) [iPlayer](#) [TV](#) [Radio](#)

NEWS LIVE BBC NEWS AT TEN

[News Front Page](#) Page last updated at 10:35 GMT, Thursday, 24 December 2009

[World](#) [E-mail this to a friend](#) [Printable version](#)

[UK](#)

[England](#)

[Northern Ireland](#)

[Scotland](#)

[Wales](#)

[Business](#)

[Politics](#)

[Health](#)

[Education](#)

[Science & Environment](#)

Technology

[Entertainment](#)

[Also in the news](#)

HP camera 'can't see' black faces

A YouTube video suggesting that face recognition cameras installed in HP laptops cannot detect black faces has had over one million views.

The short movie, uploaded earlier this month, features "Black Desi" and his colleague "White Wanda".

When Wanda, a white woman, is in front of the screen, the camera zooms to her face and moves as she moves.



"Black Desi" in the YouTube video

ML CAN BE SEXIST...

[Bolukbasi et al., 2016]



How does this happen?

I. ML MODEL MISUSE

I. ML MODEL MISUSE

Training Set



I. ML MODEL MISUSE

Training Set



ML model



I. ML MODEL MISUSE

Training Set



ML model



Real Life



I. ML MODEL MISUSE



Dataset Shift

[Quionero-Candela et al., 2009; Moreno-Torres et al., 2012;
Subbaswamy et al., 2019; Ovadia et al., 2019; Rabanser et al., 2019]

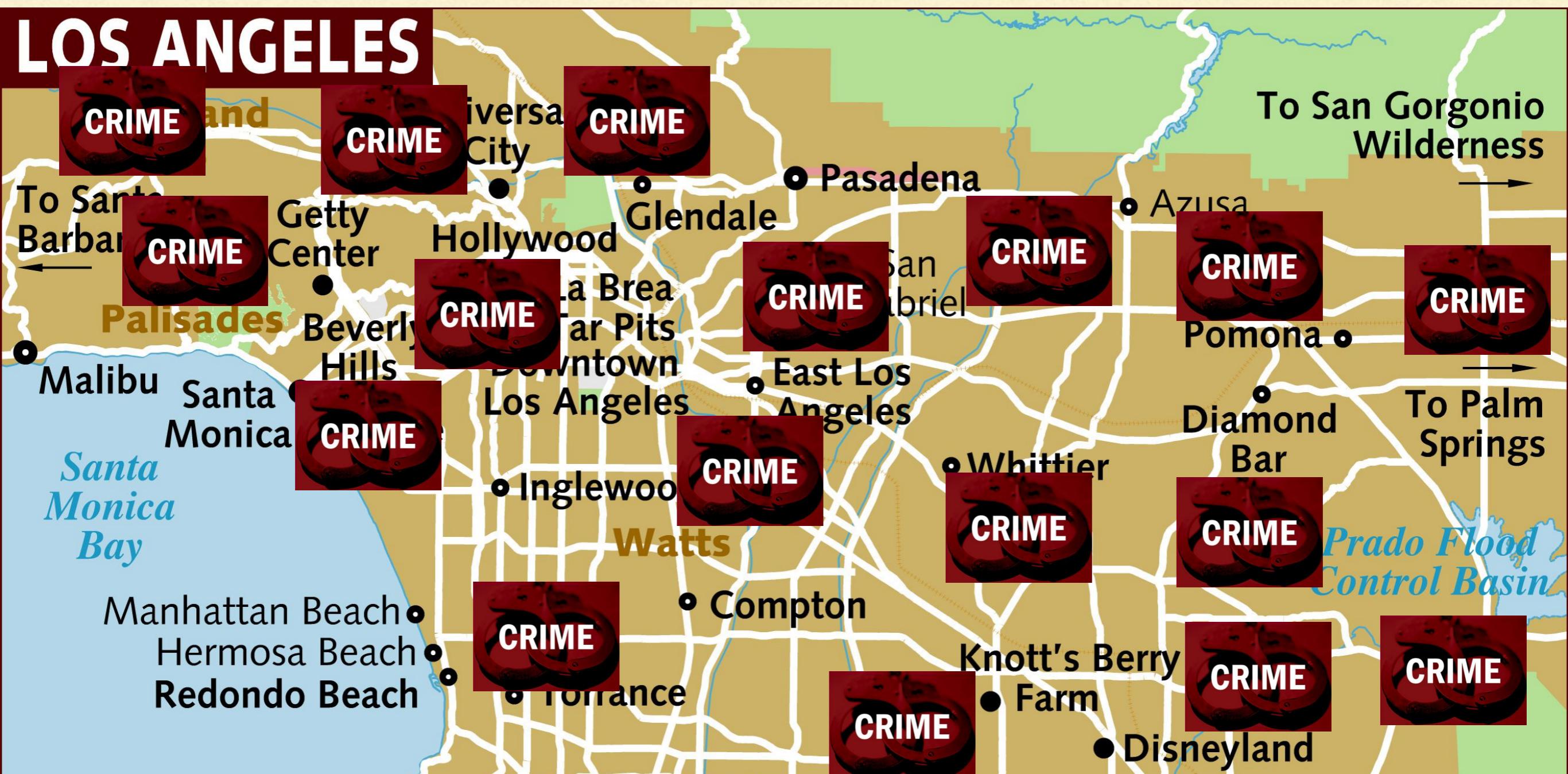
I. ML MODEL MISUSE

[Lum & Isaac, 2016]



I. ML MODEL MISUSE

[Lum & Isaac, 2016]



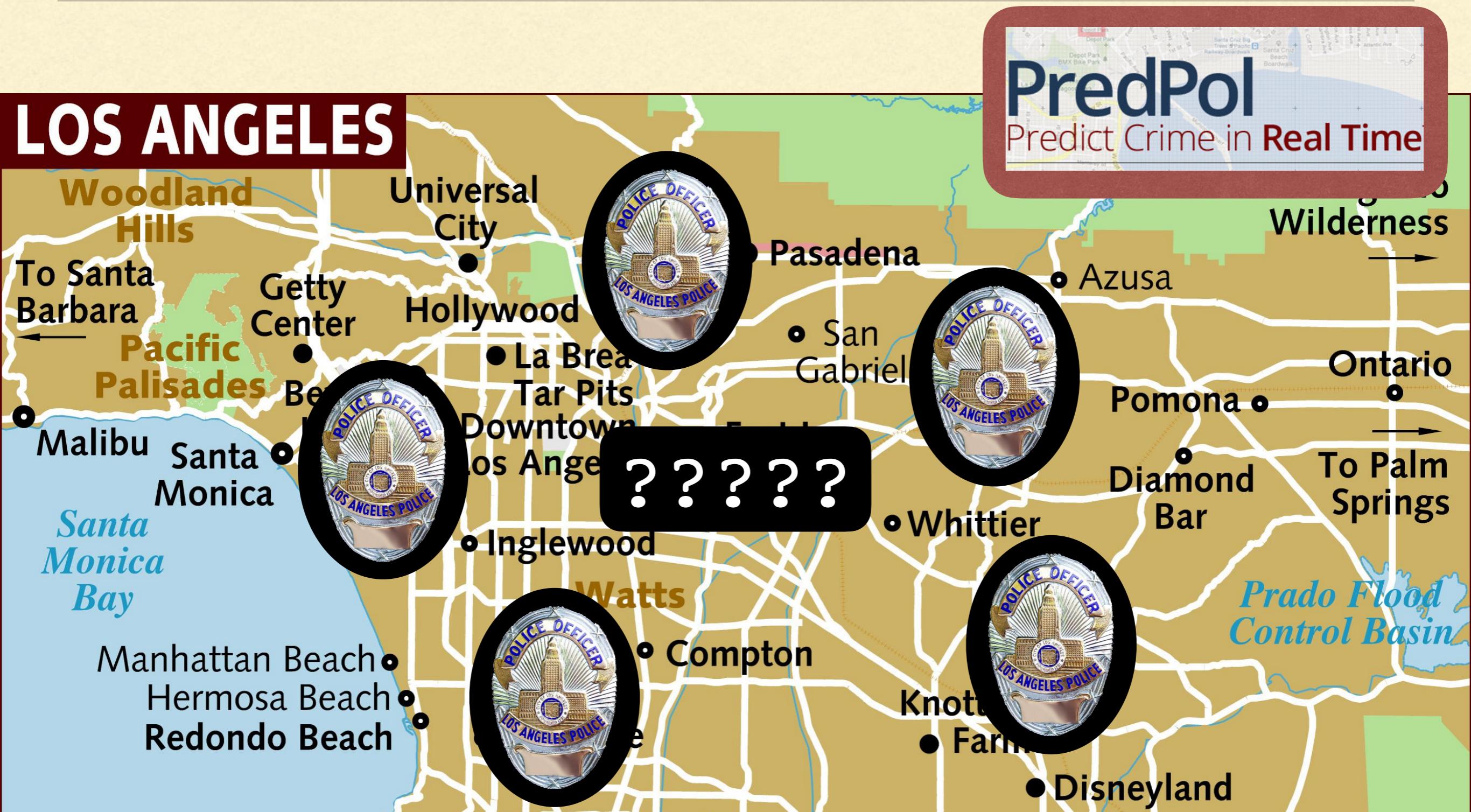
I. ML MODEL MISUSE

[Lum & Isaac, 2016]



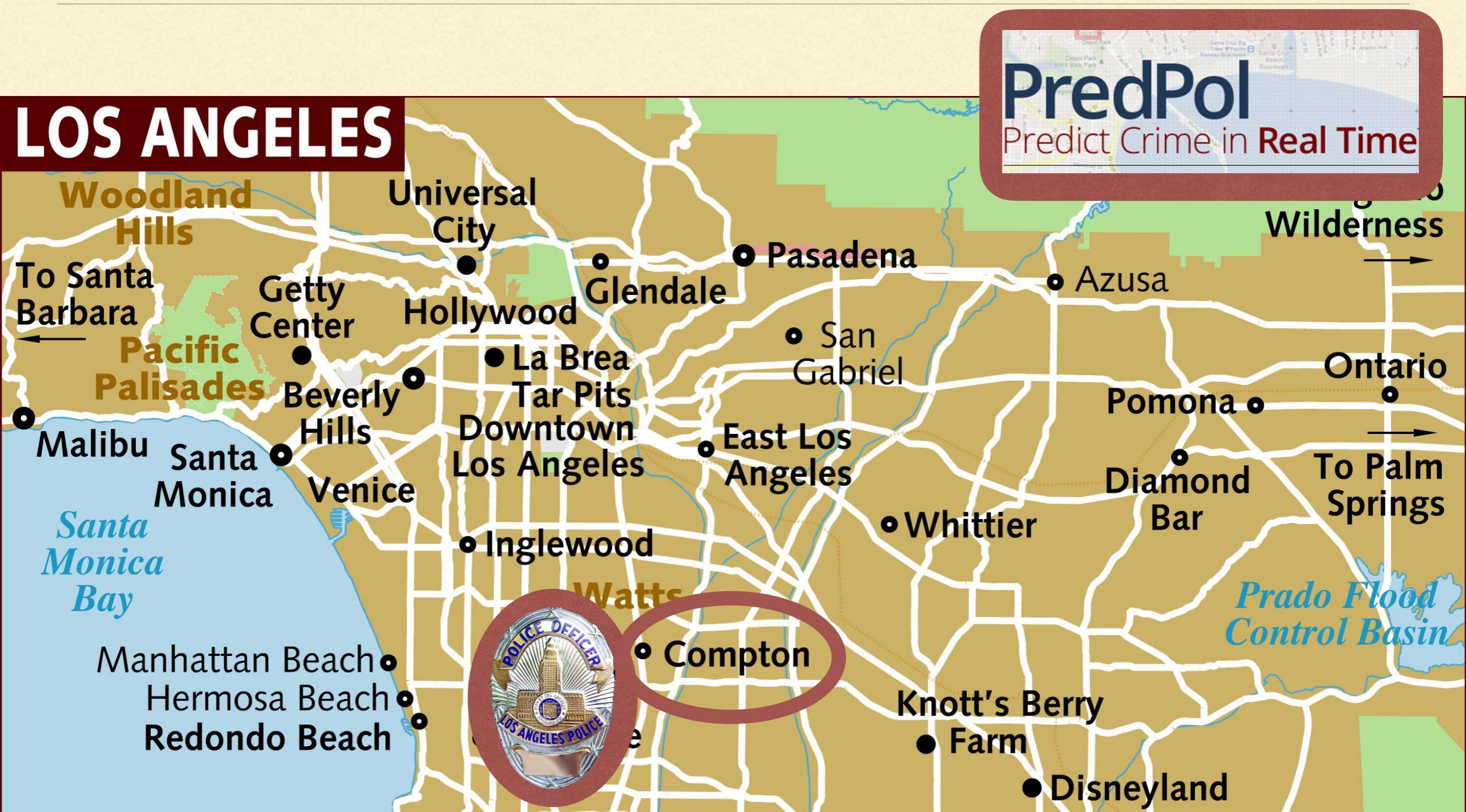
I. ML MODEL MISUSE

[Lum & Isaac, 2016]



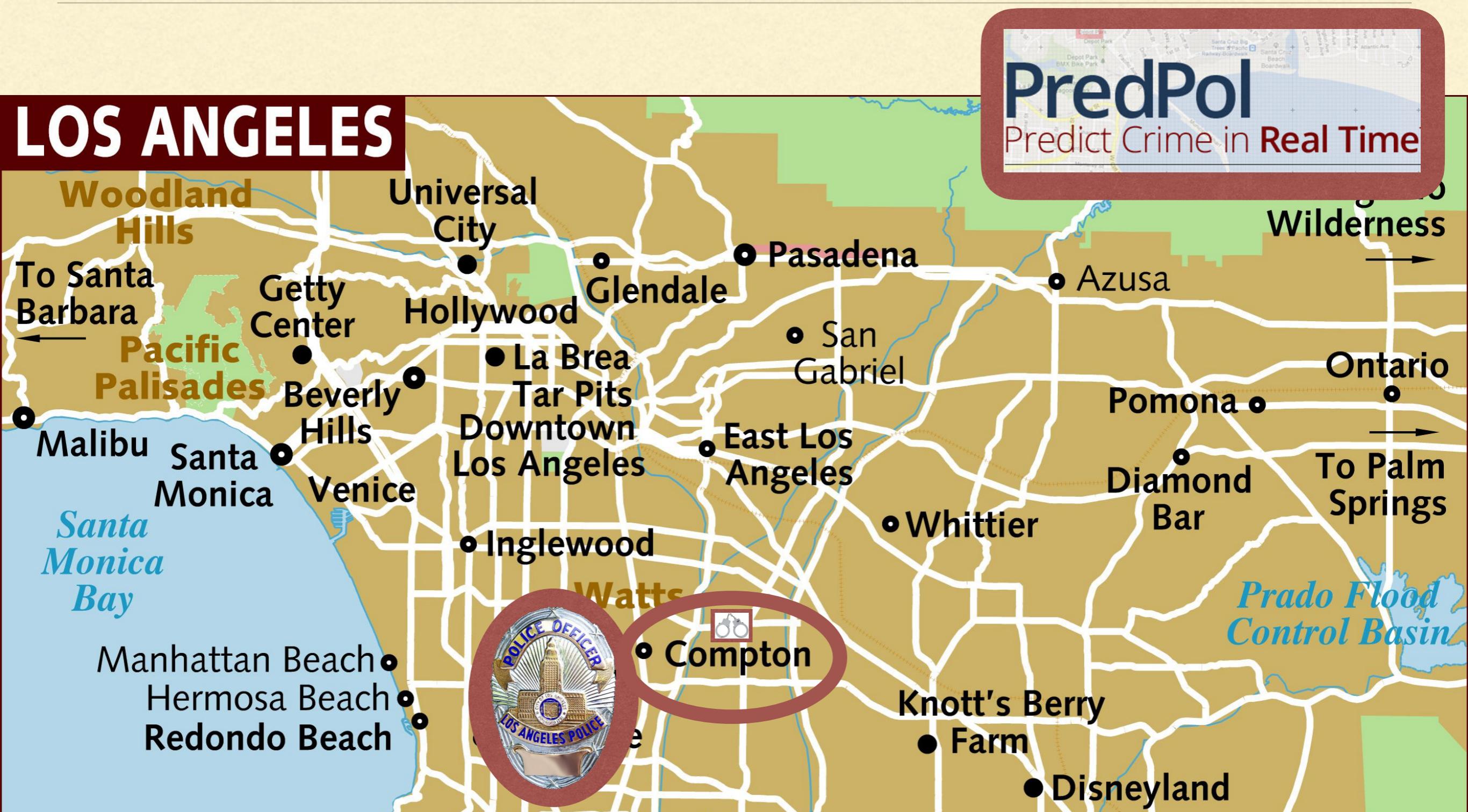
I. ML MODEL MISUSE

[Lum & Isaac, 2016]



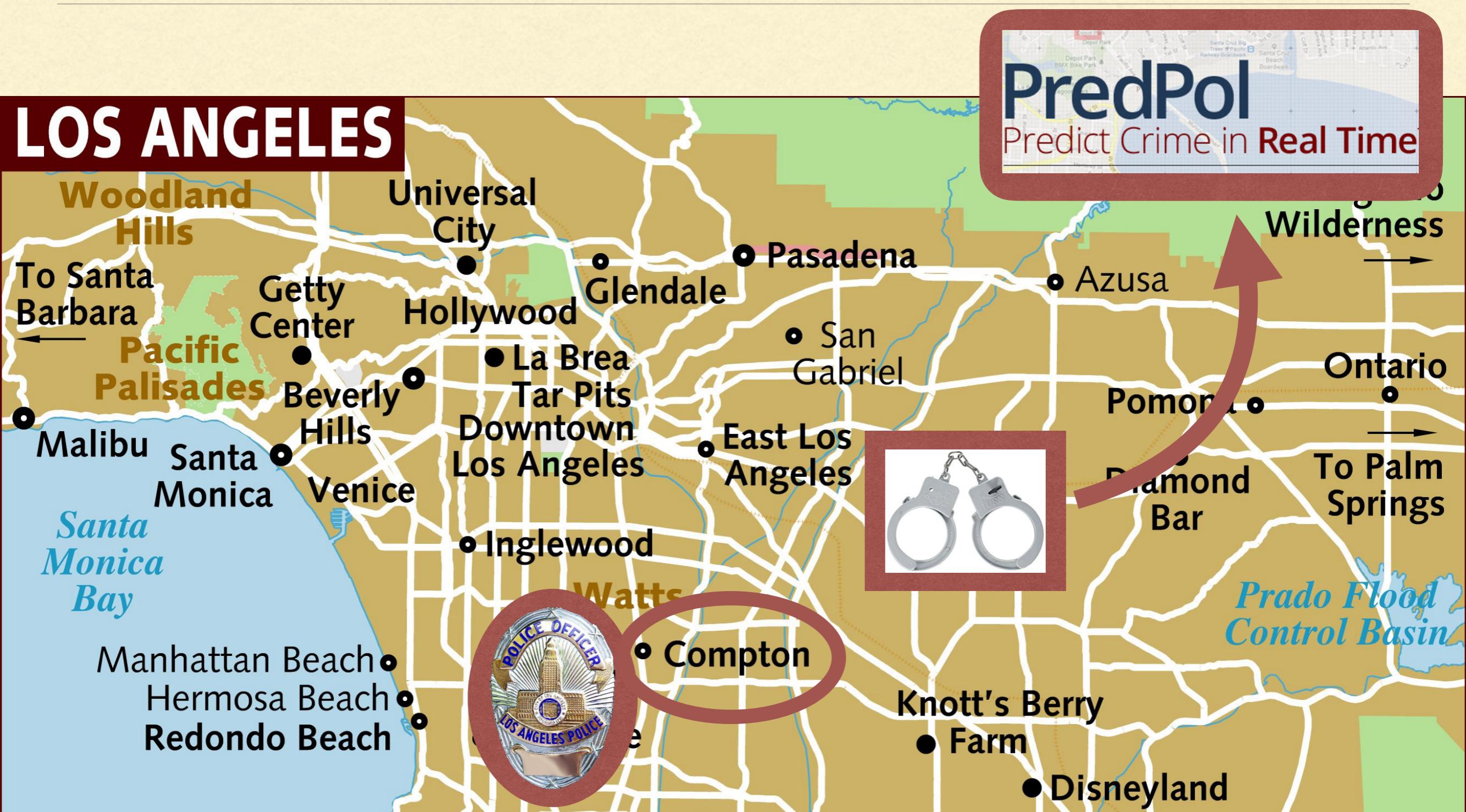
I. ML MODEL MISUSE

[Lum & Isaac, 2016]



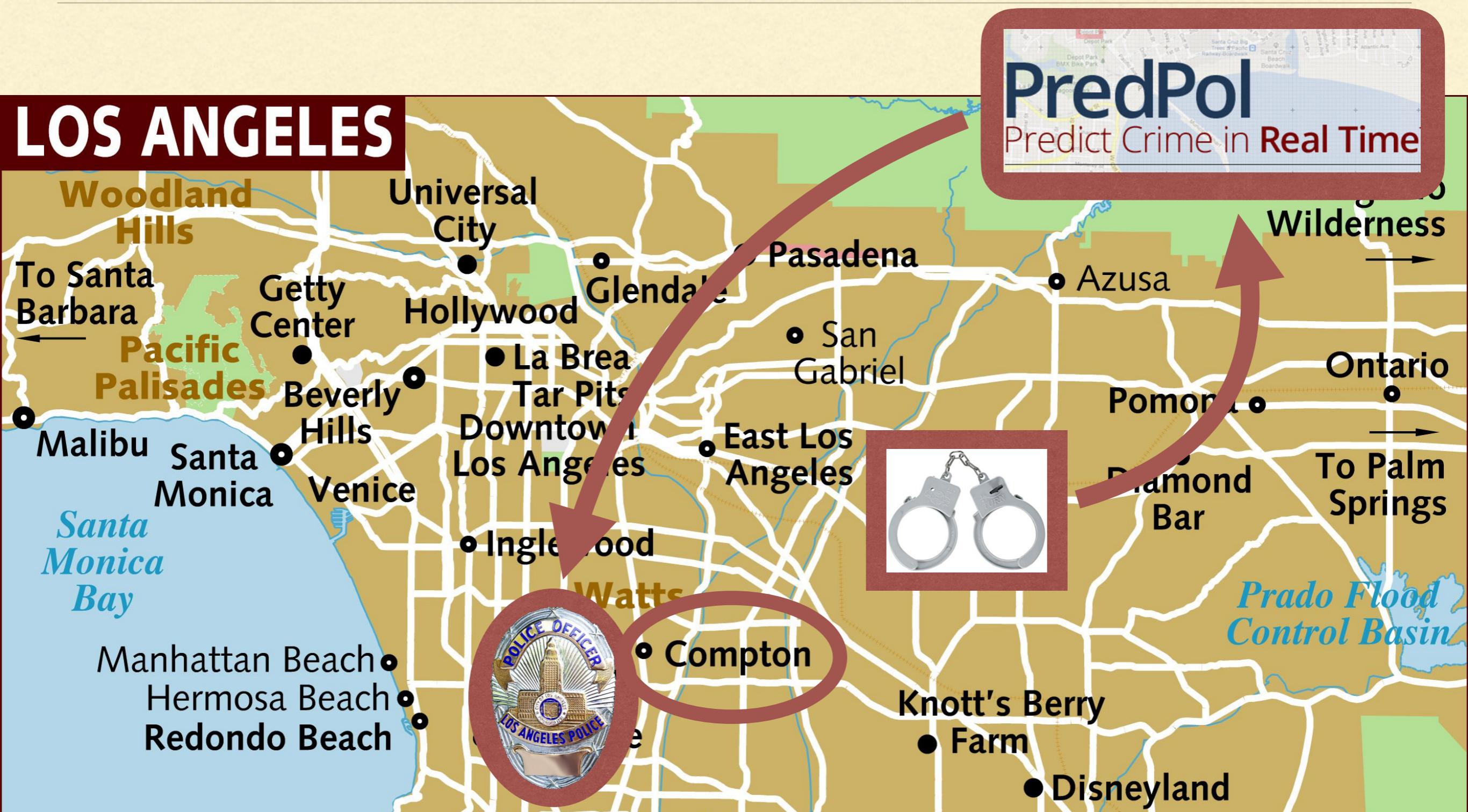
I. ML MODEL MISUSE

[Lum & Isaac, 2016]



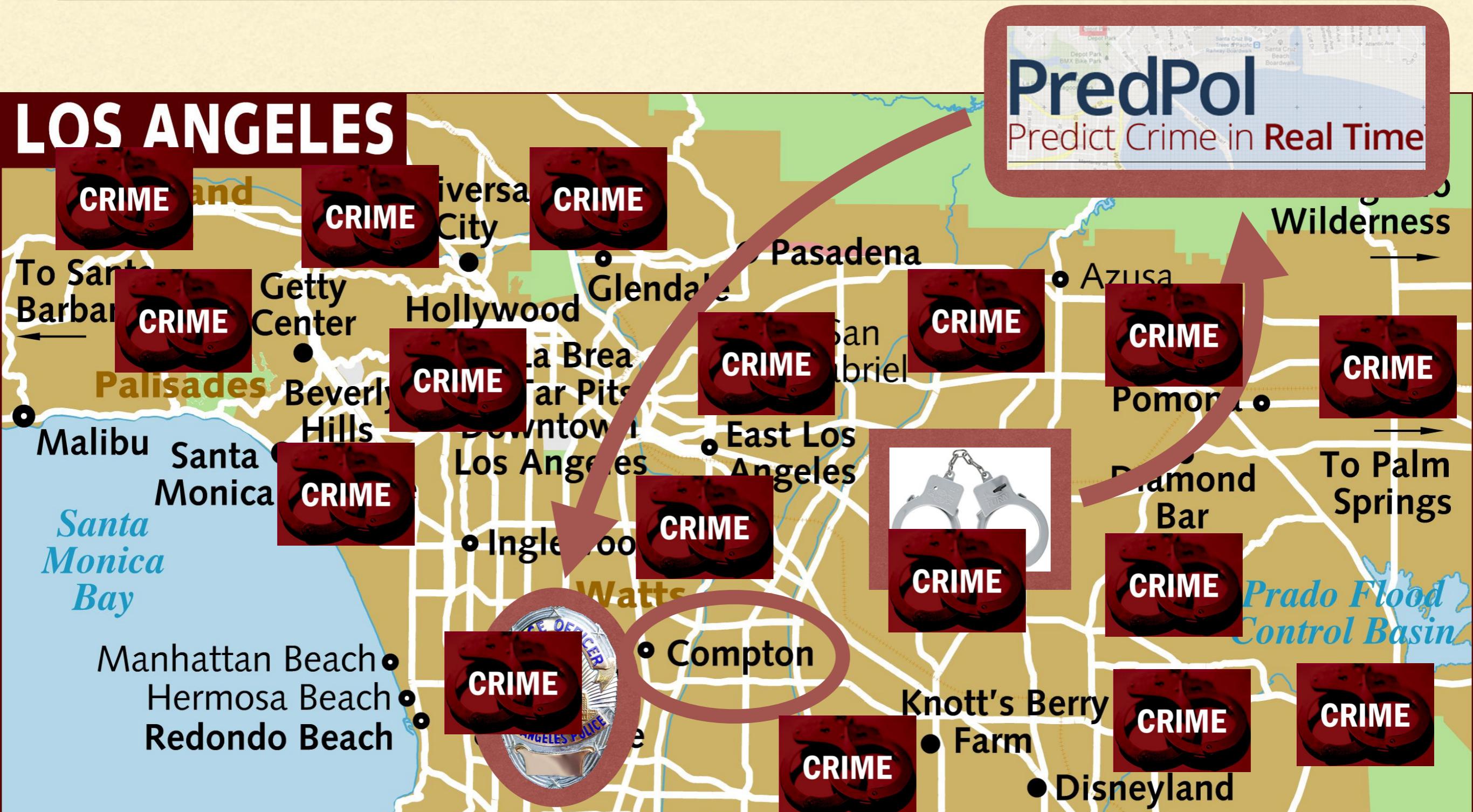
I. ML MODEL MISUSE

[Lum & Isaac, 2016]



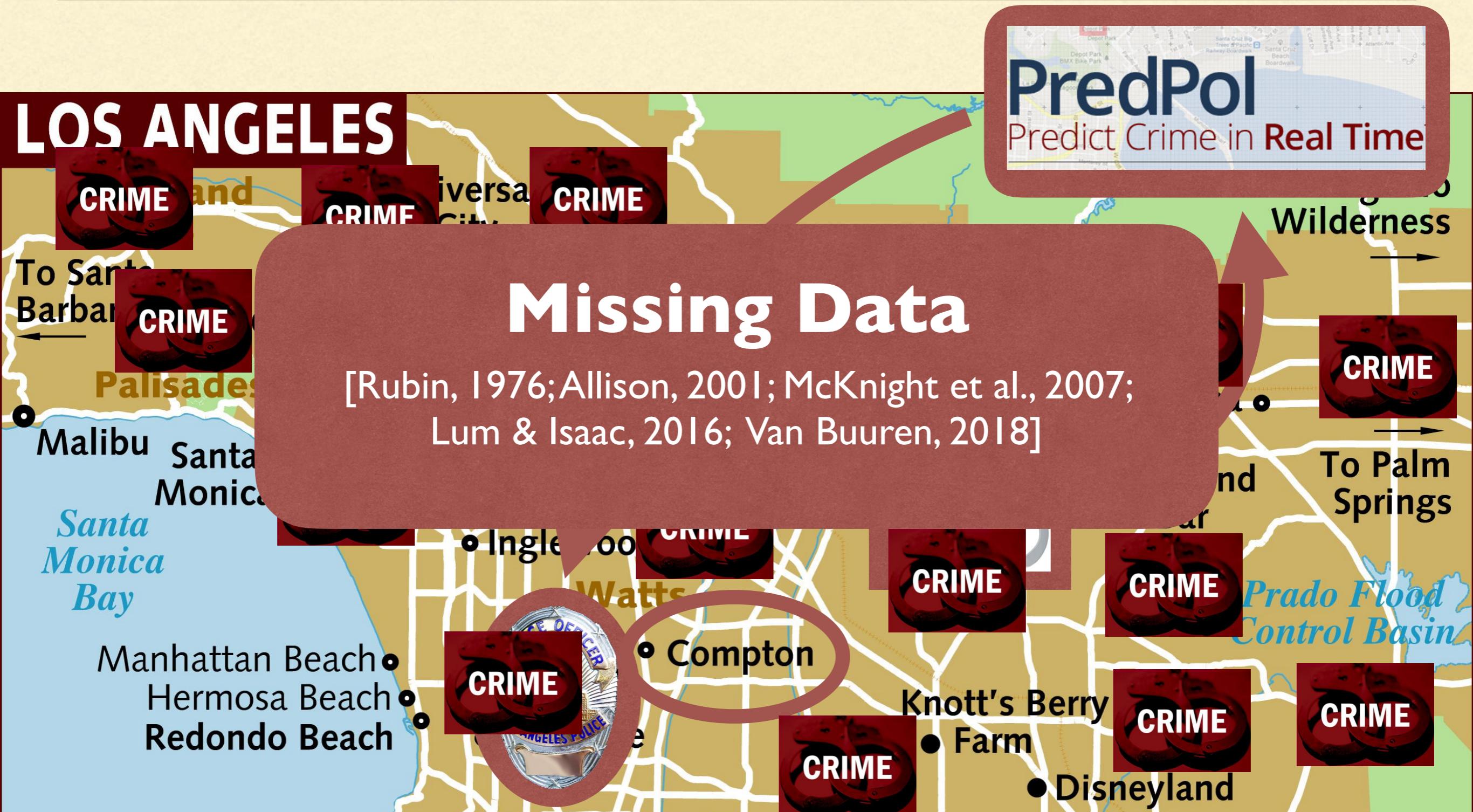
I. ML MODEL MISUSE

[Lum & Isaac, 2016]



I. ML MODEL MISUSE

[Lum & Isaac, 2016]



2. DATA ITSELF

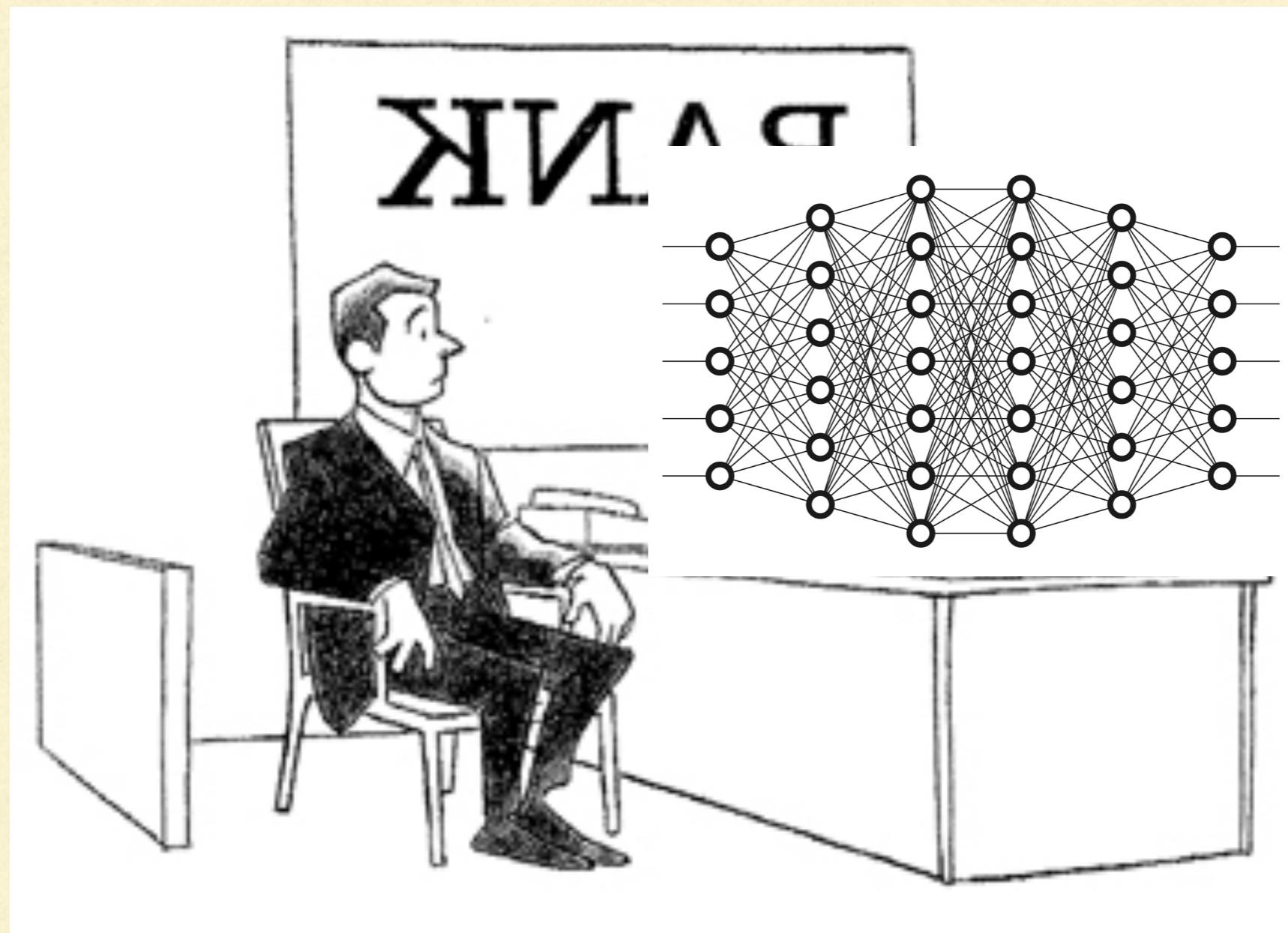
2. DATA ITSELF



2. DATA ITSELF

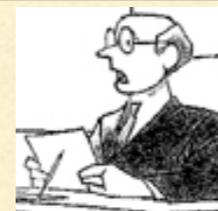


2. DATA ITSELF



2. DATA ITSELF

	income	credit score	enough in savings?	
	\$\$			 success
	\$			 success
	\$\$			 failure
	\$\$\$			 failure



2. DATA ITSELF

	income	credit score	enough in savings?		
	\$\$				
	\$				
	\$\$				
	\$\$\$				

this
banker is a
racist!

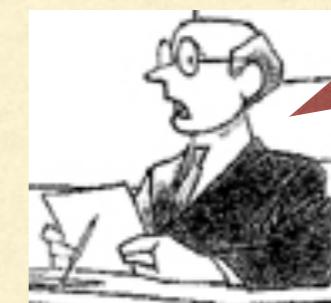
2. DATA ITSELF



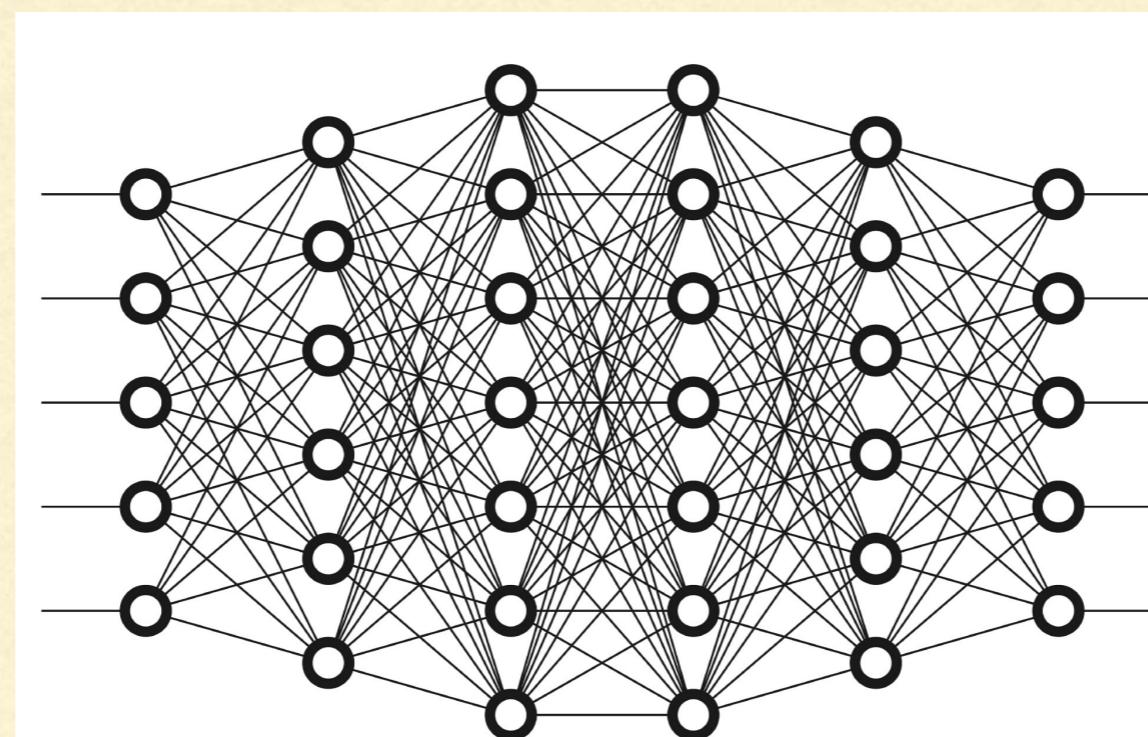
income credit enough in

score savings?

\$\$\$\$



Loan
application
failed



2. DATA ITSELF



income credit enough in
score savings?

\$\$\$\$

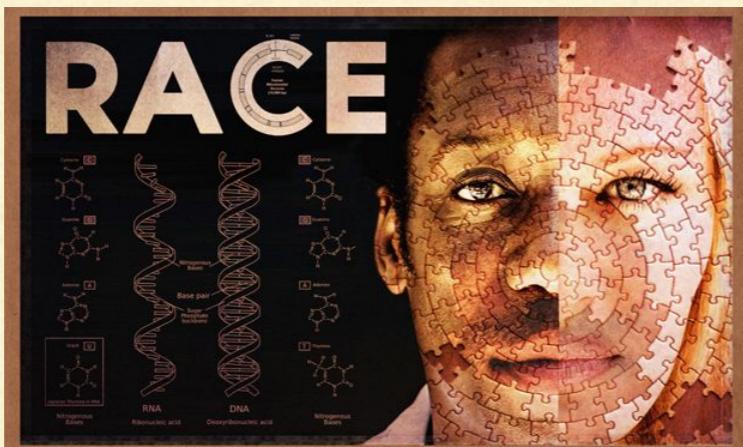


Loan
application
failed

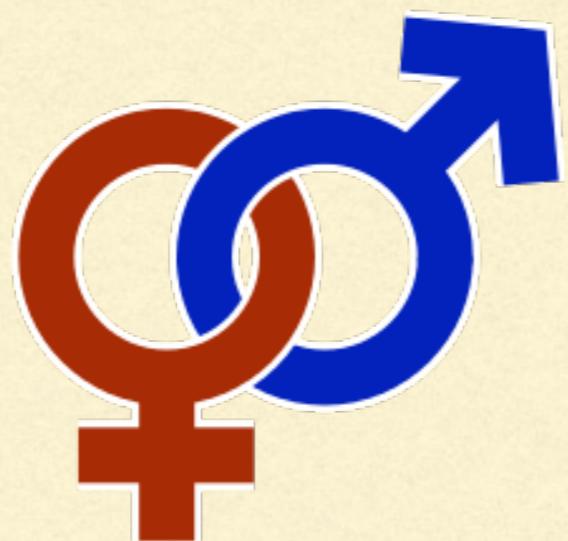
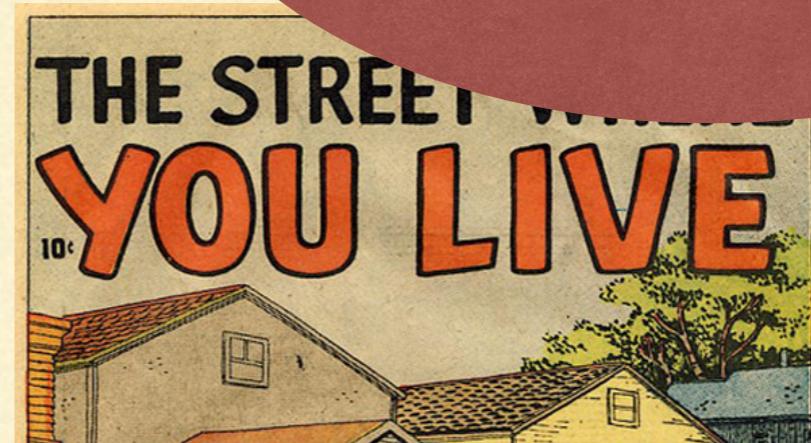


2. DATA ITSELF

Proxies



?



?



red-lining
excluded black people
from housing and services
[Sagawa and Segal, 2000]

taller applicants
are more often
selected for jobs
[Hensley & Cooper, 1987]

ILLEGAL ALGORITHMS?

Regulated domains in the US

- Credit (Equal Credit Opportunity Act)
- Education (Civil Rights Act of 1964)
- Employment (Civil Rights Act of 1964)
- Housing (Fair Housing Act)
- ‘Public Accommodation’ (Civil Rights Act of 1964)

ILLEGAL ALGORITHMS?

Legally protected sensitive attributes

- Race (Civil Rights Act of 1964)
- Color (Civil Rights Act of 1964)
- Sex (Equal Pay Act of 1963; Civil Rights Act of 1964)
- Religion (Civil Rights Act of 1964)
- National origin (Civil Rights Act of 1964)
- Citizenship (Immigration Reform and Control Act)
- Age (Age Discrimination in Employment Act of 1967)
- Pregnancy (Pregnancy Discrimination Act)
- Familial status (Civil Rights Act of 1968)
- Disability status (Rehabilitation Act of 1973)
- Veteran status (Vietnam Era Veterans' Readjustment Assistance Act)
- Genetic information (Genetic Information Nondiscrimination Act)

ILLEGAL ALGORITHMS?

It isn't easy to subject algorithms to these laws!



[TOPICS ▾](#) [SERIES ▾](#) [NEWS APPS](#) [GET INVOLVED](#) [IMPACT](#) [ABOUT](#) [⚲](#)



MACHINE BIAS



Facebook (Still) Letting Housing Advertisers Exclude Users by Race



After ProPublica revealed last year that Facebook advertisers could target housing ads to whites only, the company announced it had built a system to spot and reject discriminatory ads. We retested and found major omissions.

by Julia Angwin, Ariana Tobin and Madeleine Varner, Nov. 21, 2017, 1:23 p.m. EST

A CHALLENGE



A 2016 US report urged data scientists to analyze “how technologies can deliberately or inadvertently **perpetuate, exacerbate, or mask discrimination**”

A CHALLENGE



A 2016 US report urged data scientists to analyze “how technologies can deliberately or inadvertently **perpetuate, exacerbate, or mask discrimination**”



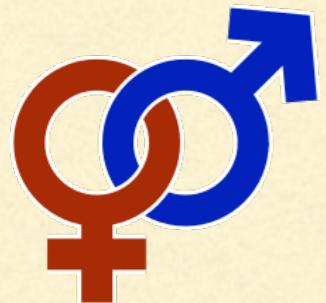
In 2016 the EU issued the General Data Protection Regulation (GDPR) which makes it a requirement to “prevent, inter alia, **discriminatory effects on natural persons**”

Ideas?

APPLYING TO LAW SCHOOL



WHO SHOULD BE ADMITTED?

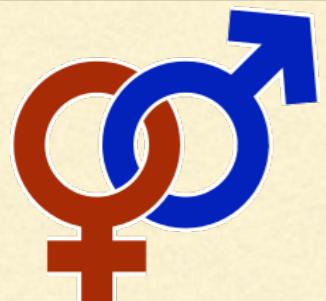


male white

female black

male black

WHO SHOULD BE ADMITTED?



male

white



female

black

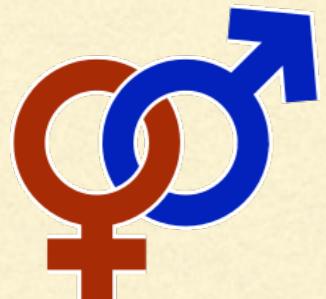


male

black



WHO SHOULD BE ADMITTED?



male

white



female

black

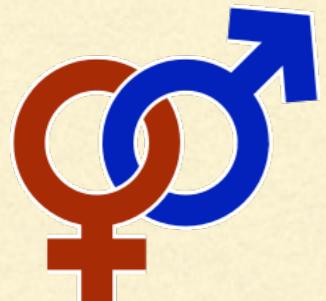


male

black



WHO SHOULD BE ADMITTED?



male

white



female

black

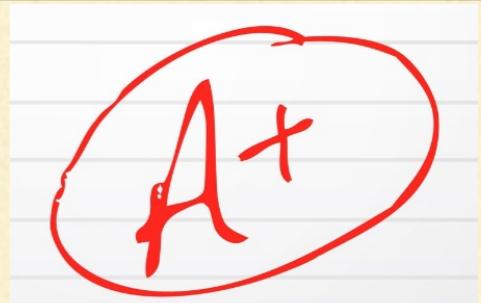
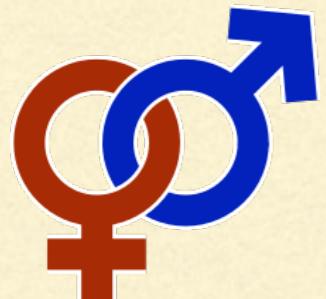


male

black



WHO SHOULD BE ADMITTED?



male

white



female

black



male

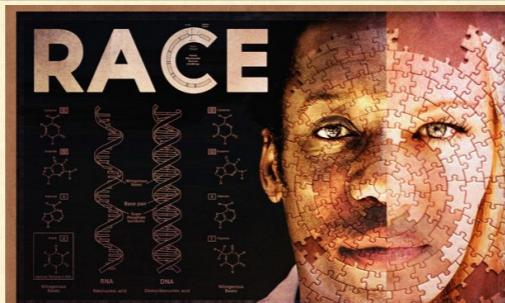
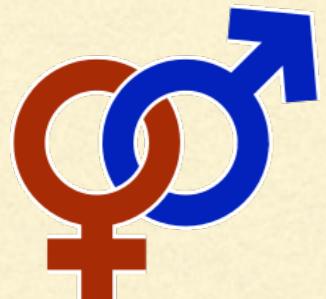
black



features χ

label Y

WHO SHOULD BE ADMITTED?



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

male

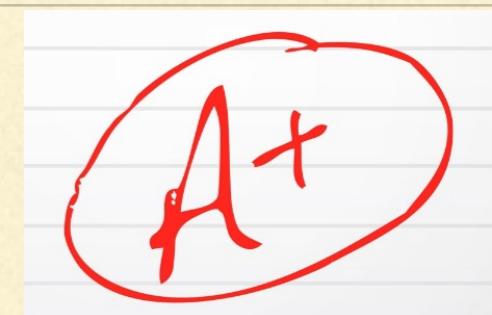
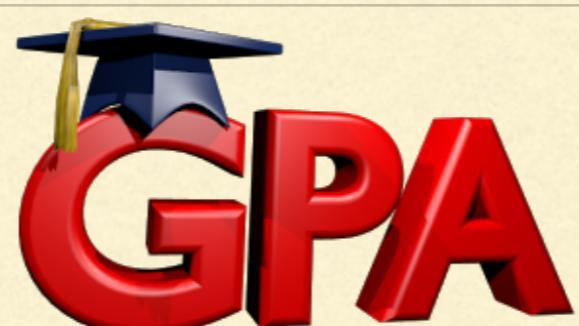
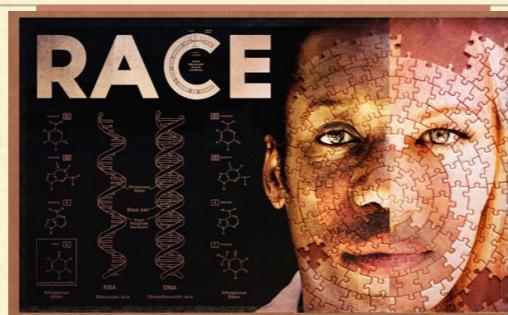
black



features \mathcal{X}

label Y

WHO SHOULD BE ADMITTED?



male	white
female	black
male	black

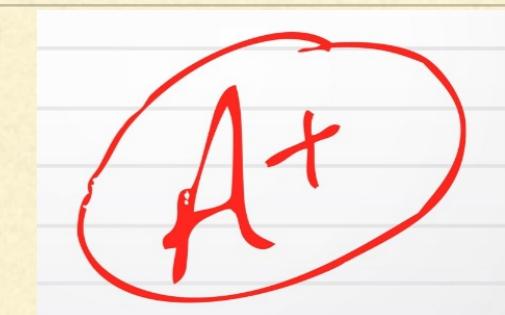
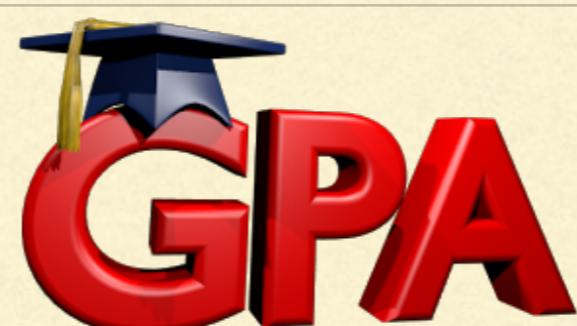
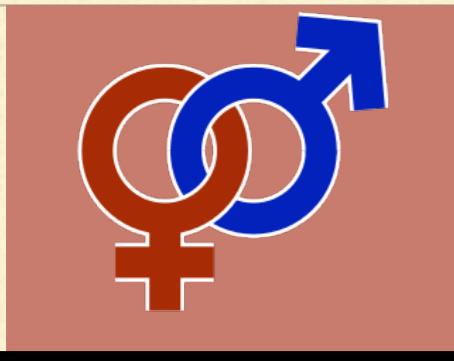
white	red
black	green
black	red



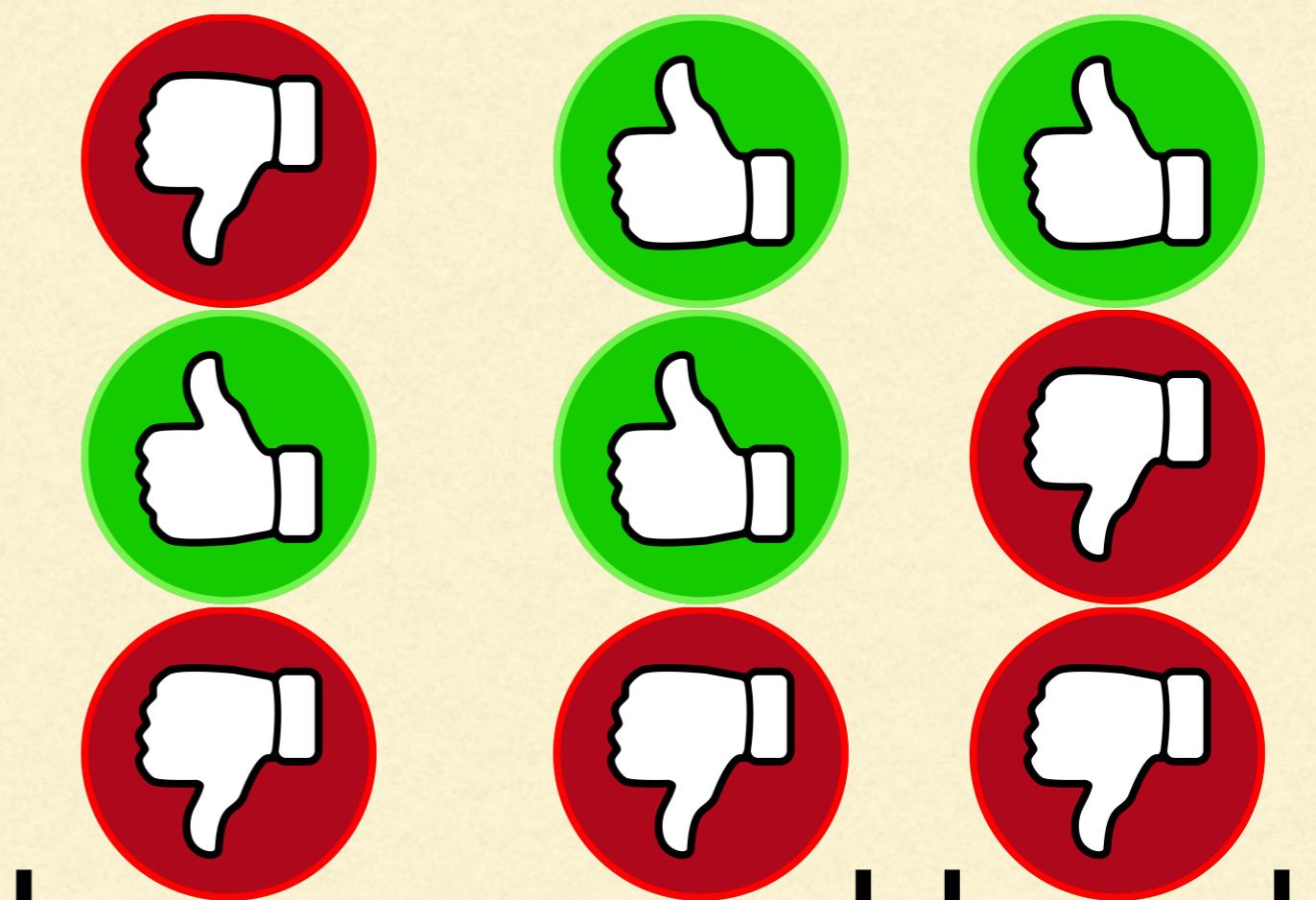
features χ

label Y

WHO SHOULD BE ADMITTED?



male	white
female	black
male	black

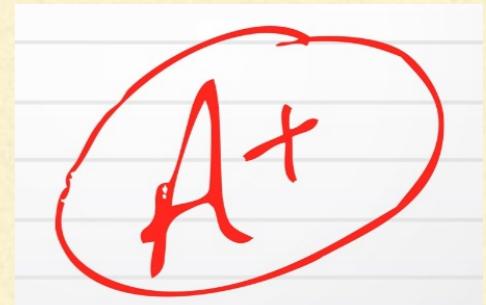
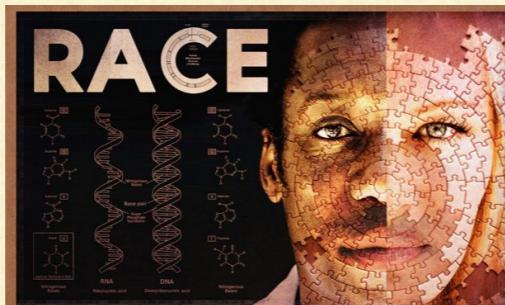
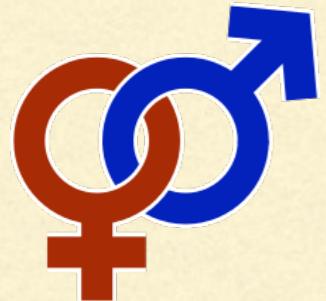


sensitive attributes A

features \mathcal{X}

label Y

Fairness Through Unawareness



male

white



female

black



male

black



sensitive attributes A

features \mathcal{X}

label Y

Fairness Through Unawareness



sensitive attributes A

features \mathcal{X}

label Y

Fairness Through Unawareness



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

sensitive attributes A

features \mathcal{X}

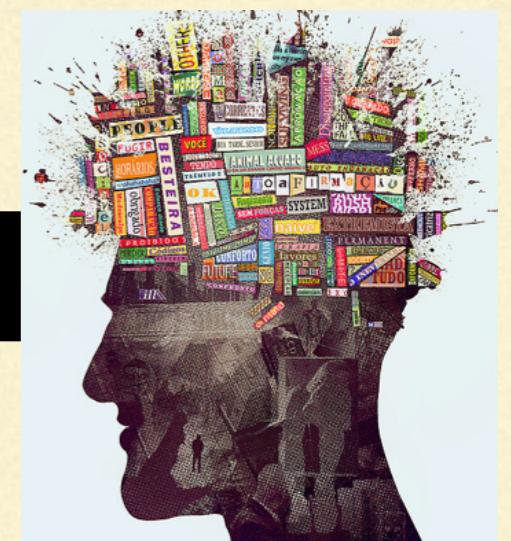
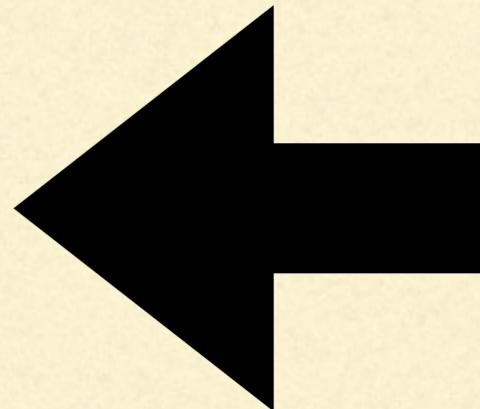
label Y

UNFAIR INFLUENCES

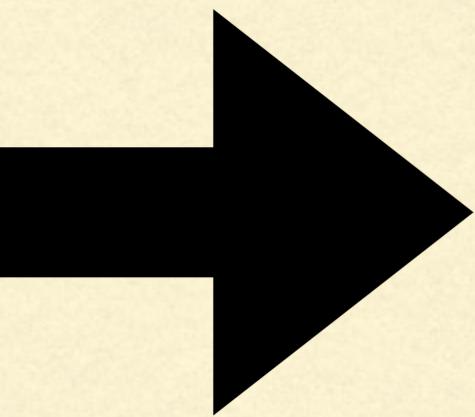
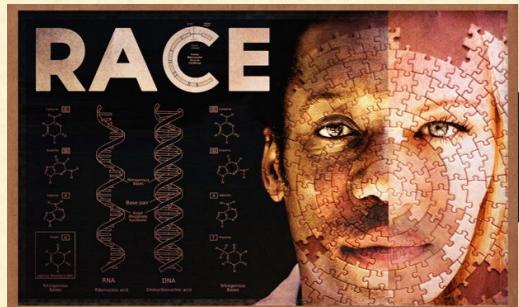


UNFAIR INFLUENCES

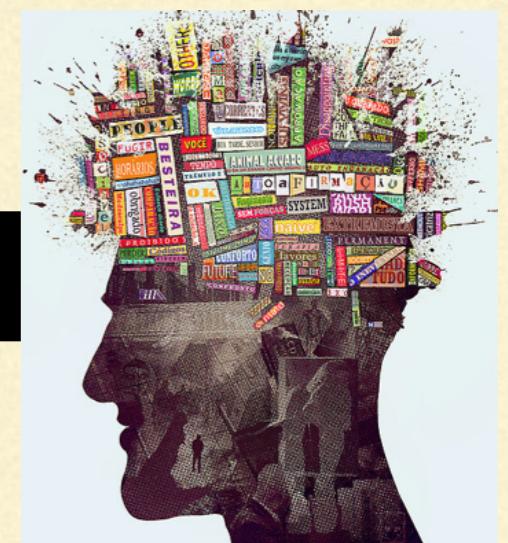
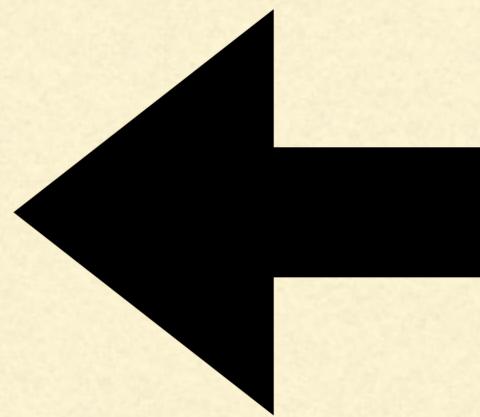
GPA



UNFAIR INFLUENCES

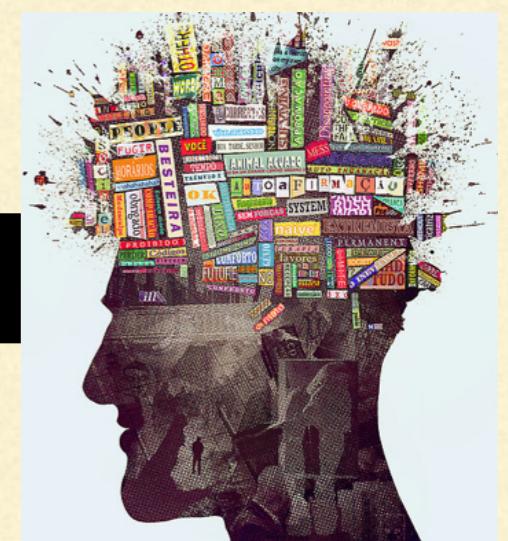
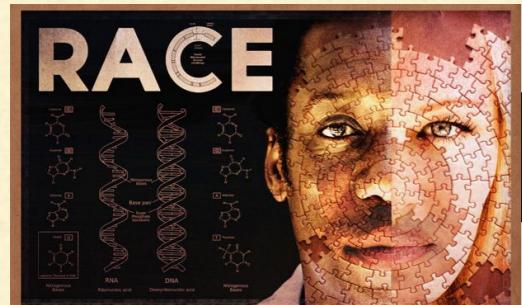


GPA



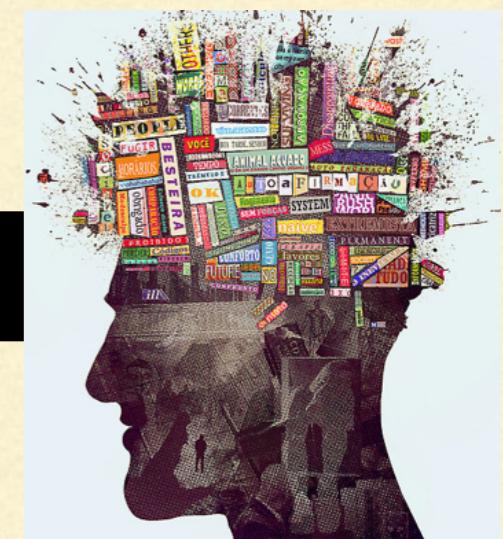
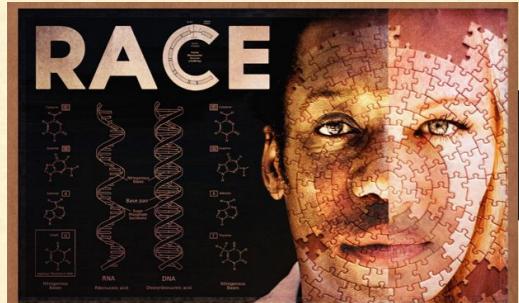
UNFAIR INFLUENCES

minority students
may feel teachers are
unsupportive
[Rowley, S. J., et al., 2014]



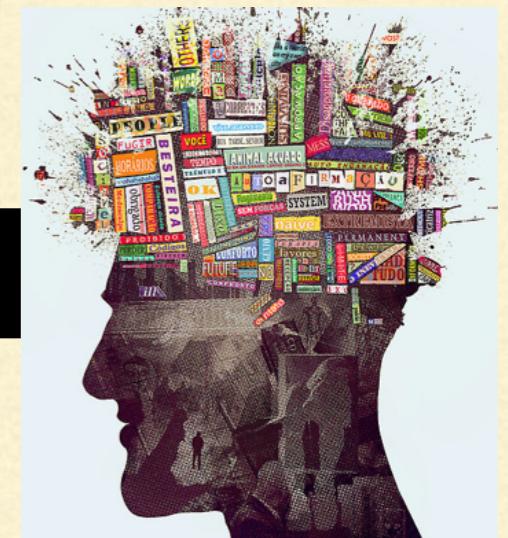
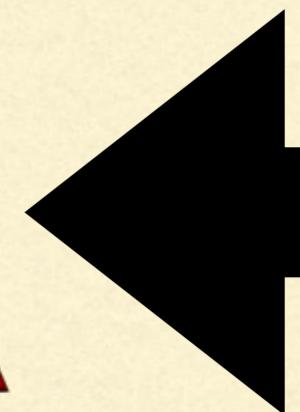
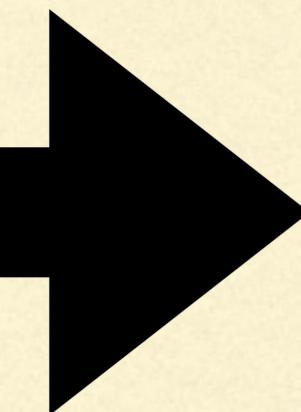
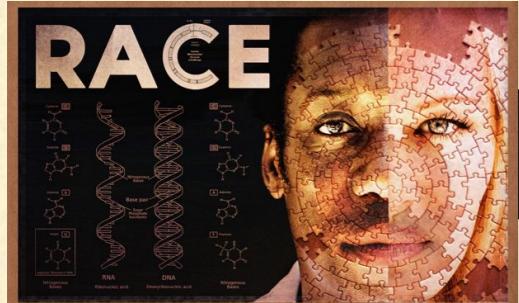
UNFAIR INFLUENCES

minority students
may feel teachers are
unsupportive
[Rowley, S. J., et al., 2014]



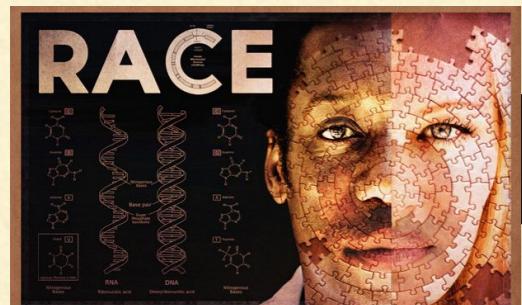
teachers may believe
minority students have
behavior issues
[Ferguson, 2003]

UNFAIR INFLUENCES

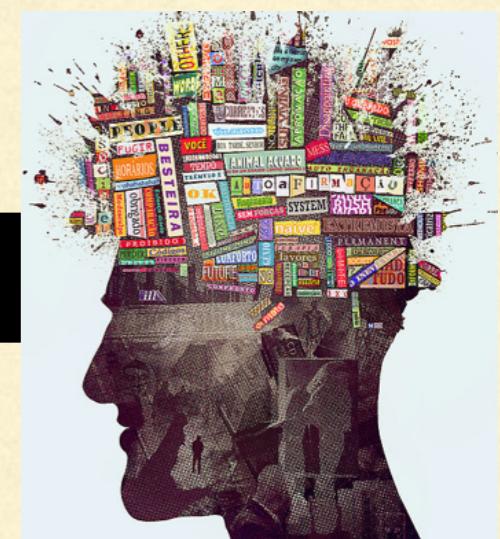
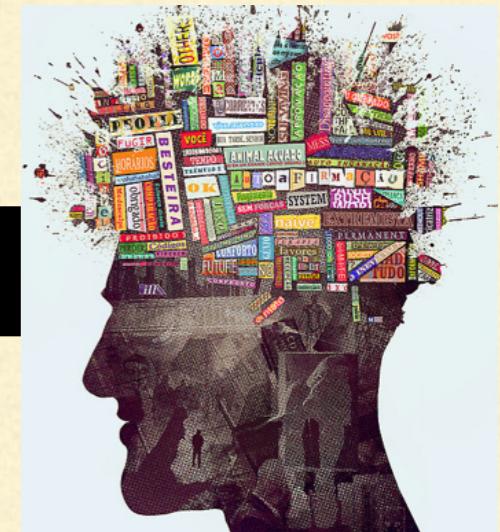


UNFAIR INFLUENCES

RACE
limited access to
academic institutions
**due to economic
history**
[Carter, 1973]



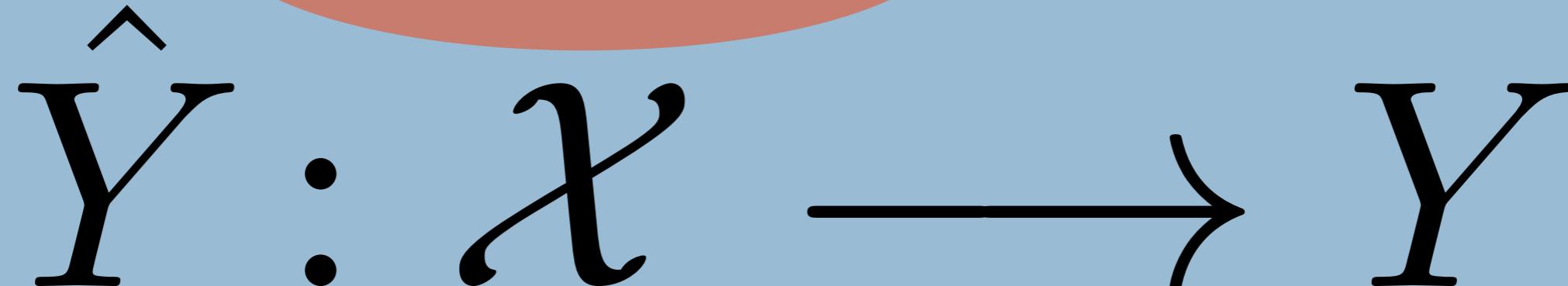
class placement
may be **different**
[Howard, 2003]



Fairness Through Unawareness



features are biased



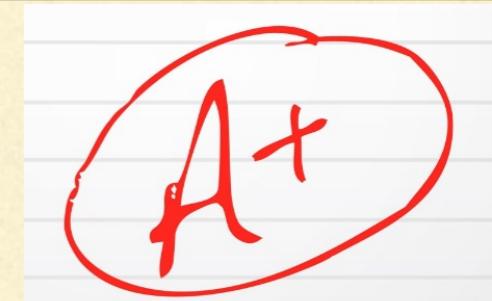
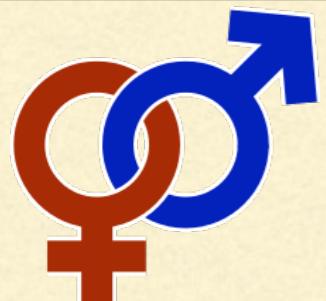
sensitive attributes A

features \mathcal{X}

label Y



WHO SHOULD BE ADMITTED?



male

white



female

black



male

black



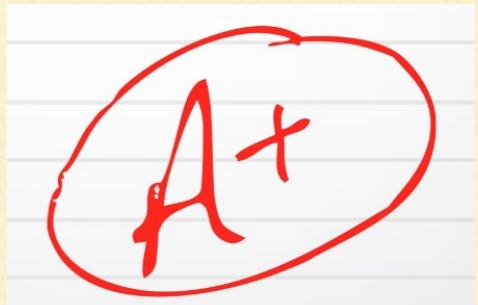
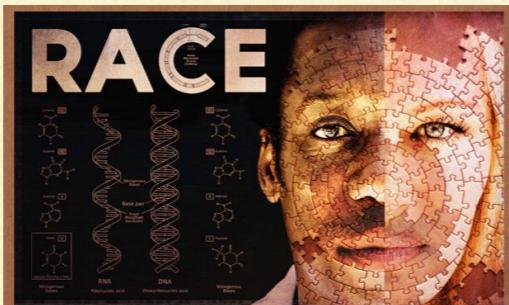
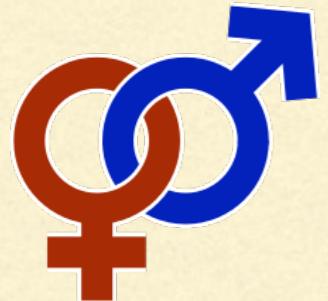
sensitive attributes A

features \mathcal{X}

label Y

Equality of Opportunity

[Hardt et al., 2016]



male

white



female

black



male

black



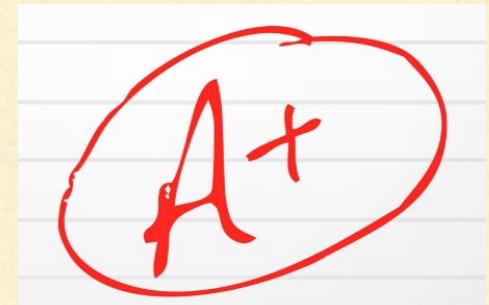
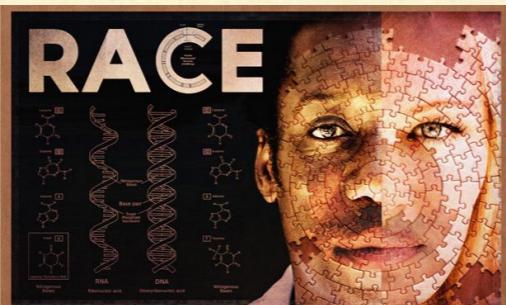
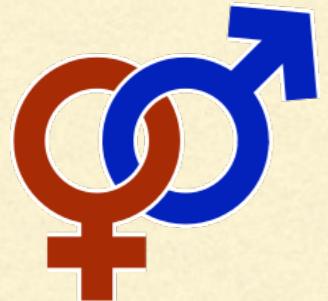
sensitive attributes A

features \mathcal{X}

label Y

Equality of Opportunity

[Hardt et al., 2016]



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

male

black

sensitive attributes A



features \mathcal{X}

label Y

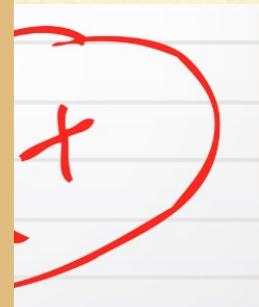
Equality of Opportunity

[Hardt et al., 2016]



black

$$\begin{aligned} P(\hat{Y} = 1 \mid A = a, Y = 1) &= \\ P(\hat{Y} = 1 \mid A = a', Y = 1) &= \end{aligned}$$



white

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

male

black

sensitive attributes A



features \mathcal{X}

label Y

Equality of Opportunity

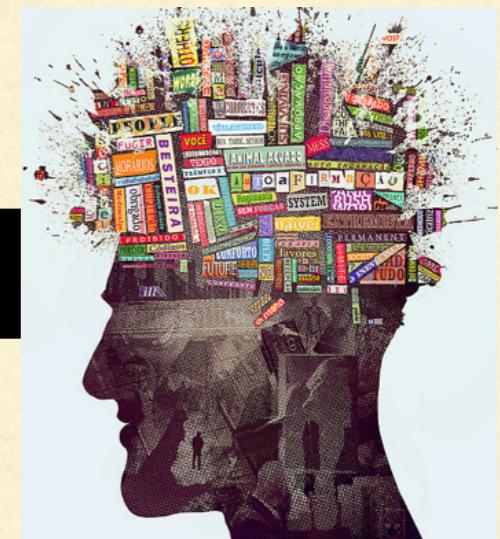
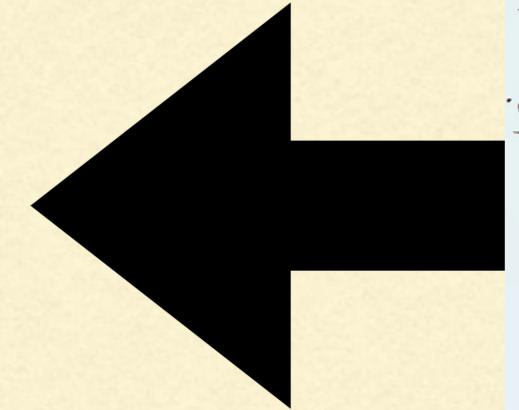
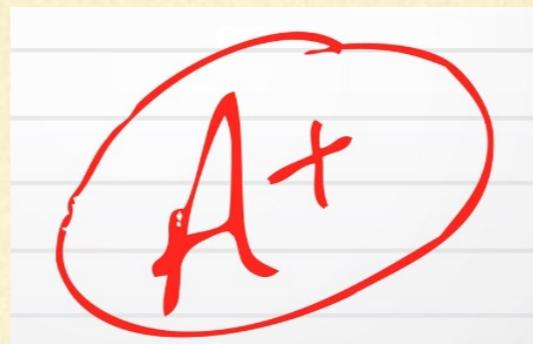
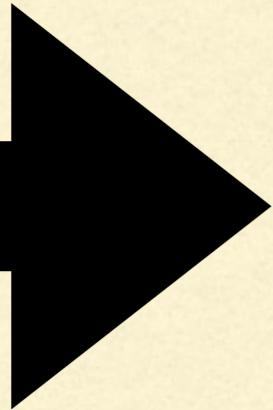
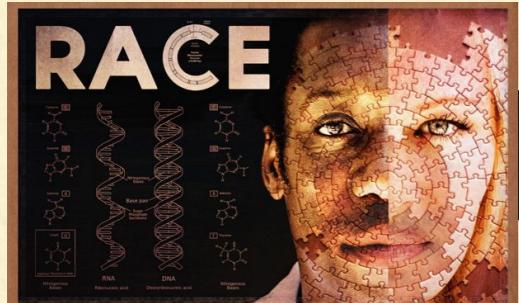
[Hardt et al., 2016]



$$\begin{aligned} \mathbb{P}(\hat{Y} = 1 \mid A = a, Y = 1) &= \\ \mathbb{P}(\hat{Y} = 1 \mid A = a', Y = 1) & \end{aligned}$$

black

white



Equality of Opportunity

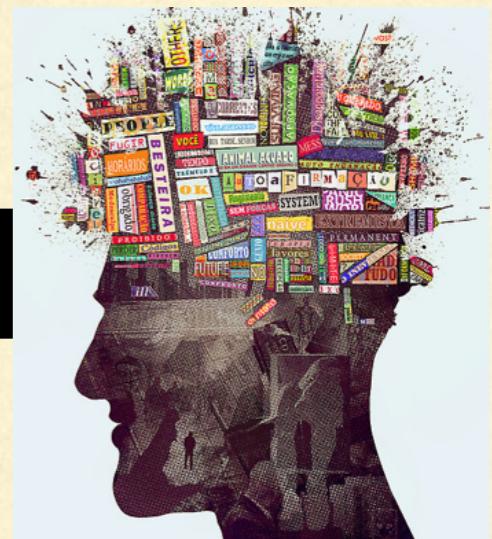
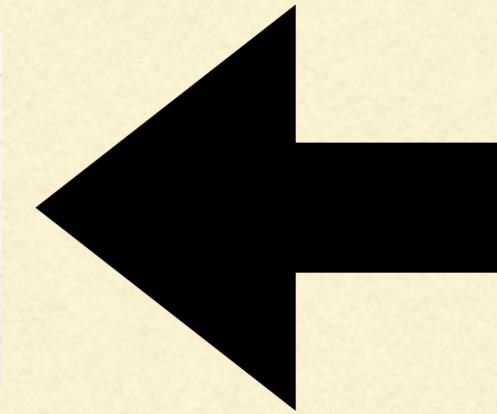
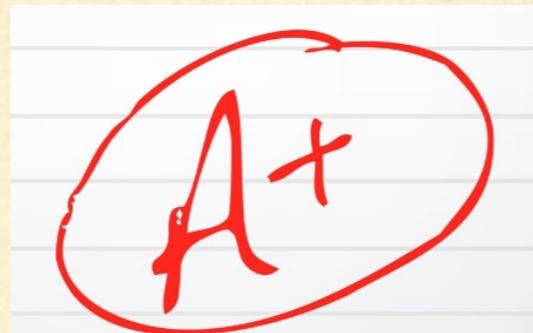
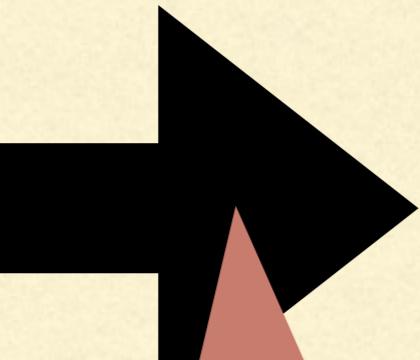
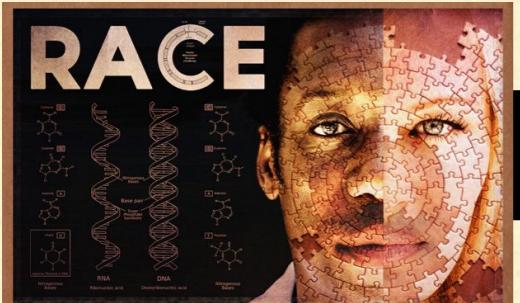
[Hardt et al., 2016]



$$\begin{aligned} \text{P}(\hat{Y} = 1 | A = a, Y = 1) &= \\ \text{P}(\hat{Y} = 1 | A = a', Y = 1) & \end{aligned}$$

black

white



lack of minority race teachers **affects**
minority student outcomes

[Birdsall et al., 2016]

Equality of Opportunity

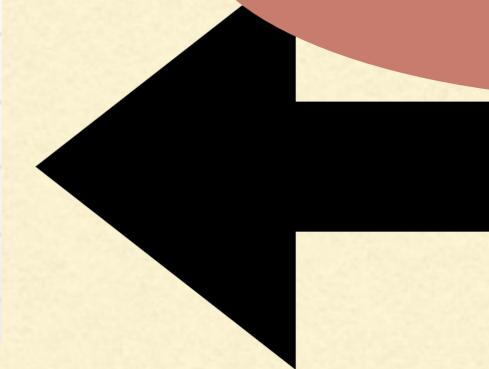
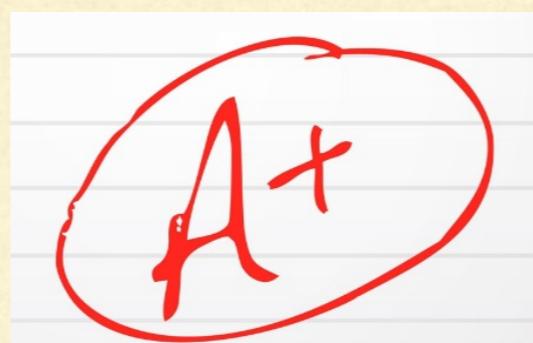
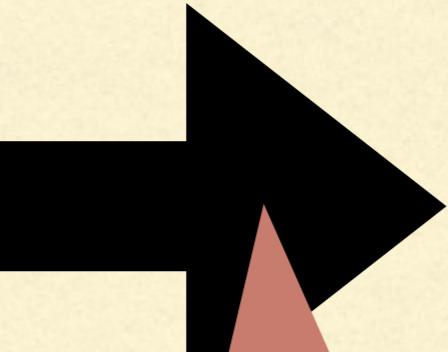
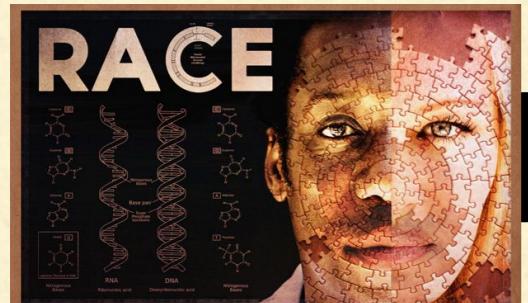
[Hardt et al., 2016]



$$\begin{aligned} \mathbb{P}(\hat{Y} = 1 \mid A = a, Y = 1) &= \\ \mathbb{P}(\hat{Y} = 1 \mid A = a', Y = 1) & \end{aligned}$$

white

50% of white students



lack of minority race teachers affects minority student outcomes

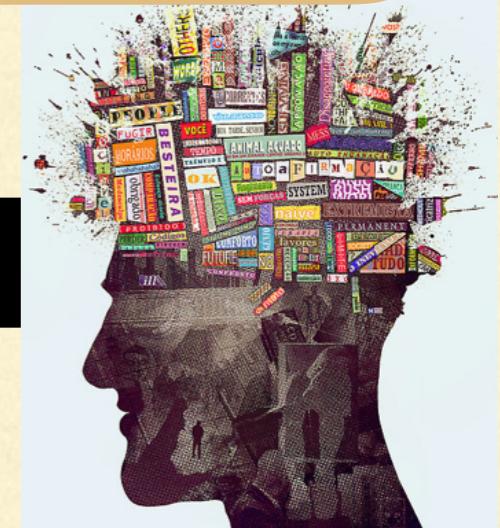
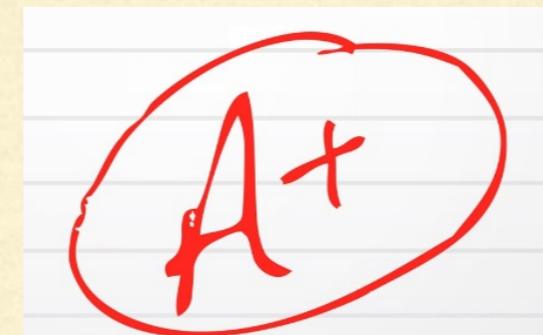
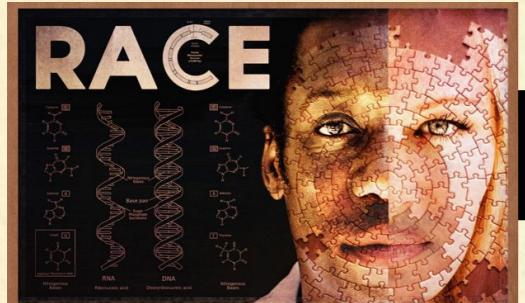
[Birdsall et al., 2016]

Equality of Opportunity

[Hardt et al., 2016]

label is biased

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

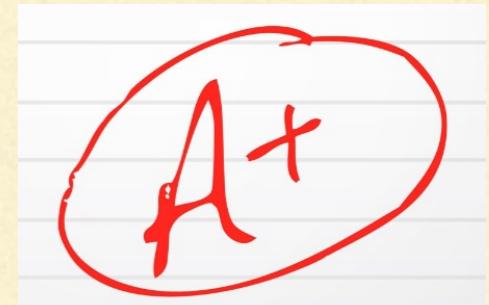
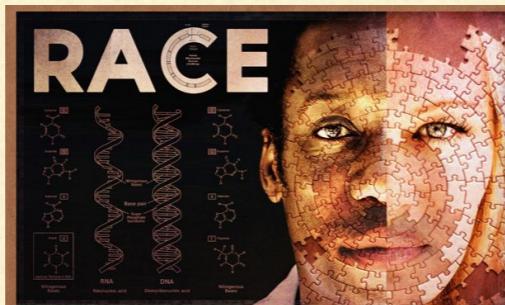
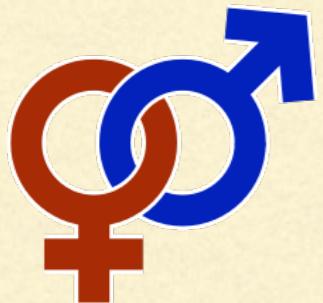


lack of minority race teachers **affects**
minority student outcomes

[Birdsall et al., 2016]

Individual Fairness

[Dwork et al., 2012]



male

white



female

black



male

black



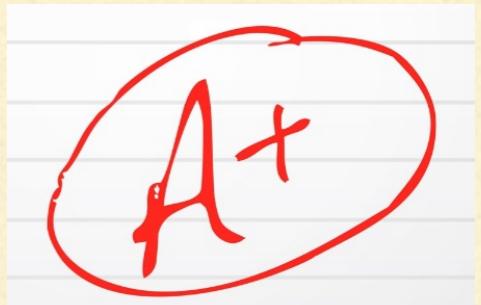
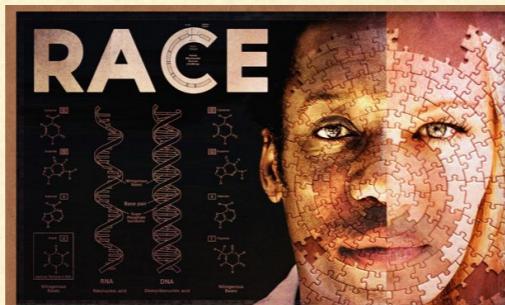
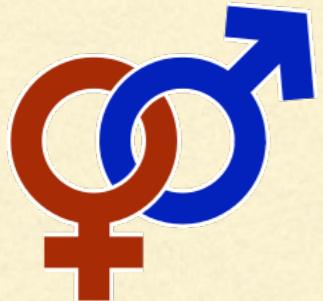
sensitive attributes A

features \mathcal{X}

label Y

Individual Fairness

[Dwork et al., 2012]

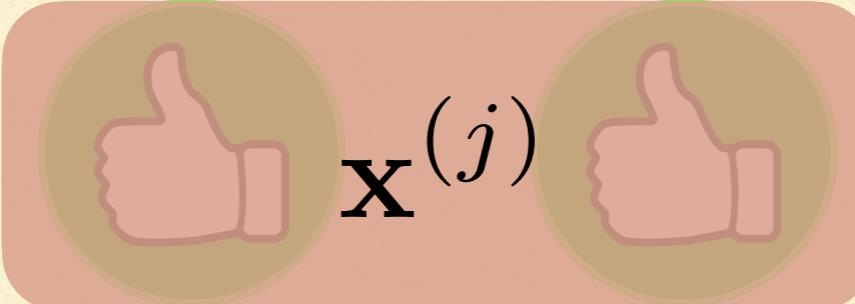
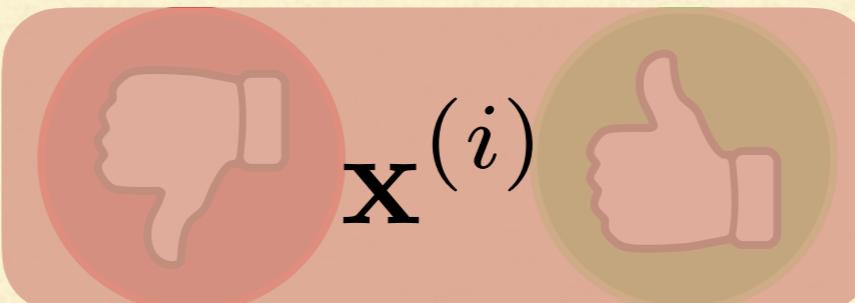


male $a^{(i)}$ white

female $a^{(j)}$ black

male black

sensitive attributes A



features \mathcal{X}

label Y

Individual Fairness

[Dwork et al., 2012]

assume ‘task-specific’ metric:

$$d(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$$

‘imposed by a regulatory body, or by a civil rights organization’

A randomized predictor \hat{Y} is fair if:

$$D(\hat{Y}(\mathbf{x}^{(i)}, \mathbf{a}^{(i)}), \hat{Y}(\mathbf{x}^{(j)}, \mathbf{a}^{(j)})) \leq d(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$$

‘distance between distributions’

Individual Fairness

[Dwork et al., 2012]

assume ‘task-specific’ metric:

$$d(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$$

‘imposed by a regulatory body, or by a civil rights organization’

How do we get this?

A randomized predictor \hat{Y} is fair if:

$$D(\hat{Y}(\mathbf{x}^{(i)}, \mathbf{a}^{(i)}), \hat{Y}(\mathbf{x}^{(j)}, \mathbf{a}^{(j)})) \leq d(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$$

‘distance between distributions’

Demographic Parity

[Calders et al., 2009; Zemel et al., 2013; Zliobaite, 2015; Zafar et al., 2015]

A predictor \hat{Y} is fair if:

Demographic Parity

[Calders et al., 2009; Zemel et al., 2013; Zliobaite, 2015; Zafar et al., 2015]

A predictor \hat{Y} is fair if:

$$\mathbb{P}(\hat{Y}|A=1) = \mathbb{P}(\hat{Y}|A=0)$$

e.g.

black

white

STATISTICAL FAIRNESS PROPOSALS

Fairness Through
Unawareness

Equality of
Opportunity
[Hardt et al., 2016]

Individual Fairness
[Dwork et al., 2012]

Demographic Parity
[Zemel et al., 2013; Zliobaite, 2015]

Fair Calibration
[Pleiss et al., 2017]

Preference Fairness
[Zafar et al., 2017]

WHICH DEFINITION SHOULD WE CHOOSE?

Fairness Through
Unawareness

Equality of
Opportunity
[Hardt et al., 2016]

Individual Fairness
[Dwork et al., 2012]

Demographic Parity
[Zemel et al., 2013; Zliobaite, 2015]

Fair Calibration
[Pleiss et al., 2017]

Preference Fairness
[Zafar et al., 2017]

WHICH DEFINITION SHOULD WE CHOOSE?

Fairness Through
Unawareness

Equality of
Opportunity
[Hardt et al., 2016]

Individual Fairness
[Dwork et al., 2012]

Demographic Parity
[Zemel et al., 2013; Zliobaite, 2015]

Fair Calibration
[Pleiss et al., 2017]

Preference Fairness
[Zafar et al., 2017]

None of them!

Claim: Any fairness notion based
on **observation alone** is
unable to ‘fix’ discrimination

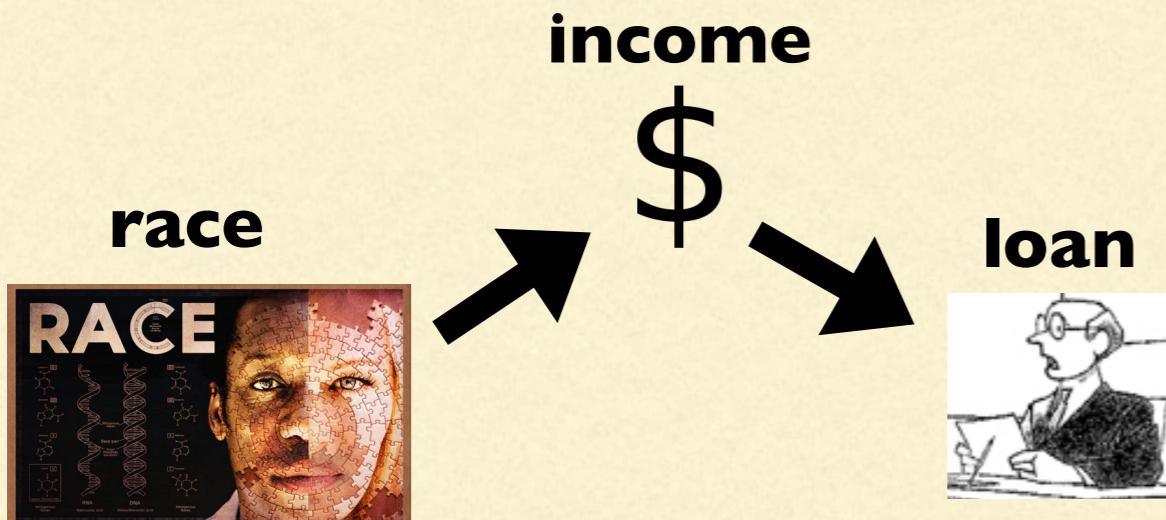
Why?: They cannot model
how discrimination happens

STATISTICS CAN'T MODEL HOW

“Race is *correlated* with loan decision!”

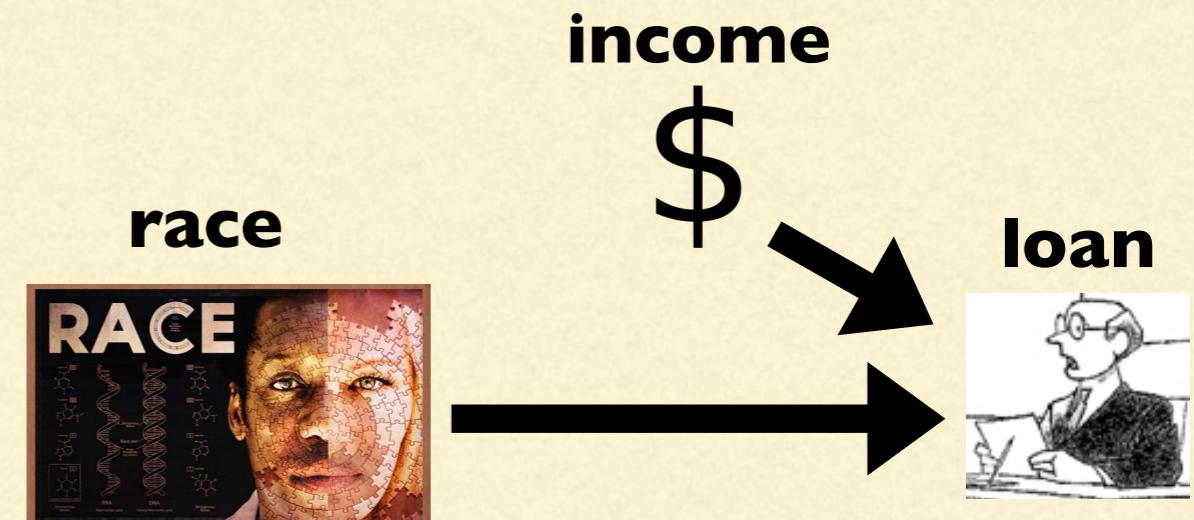
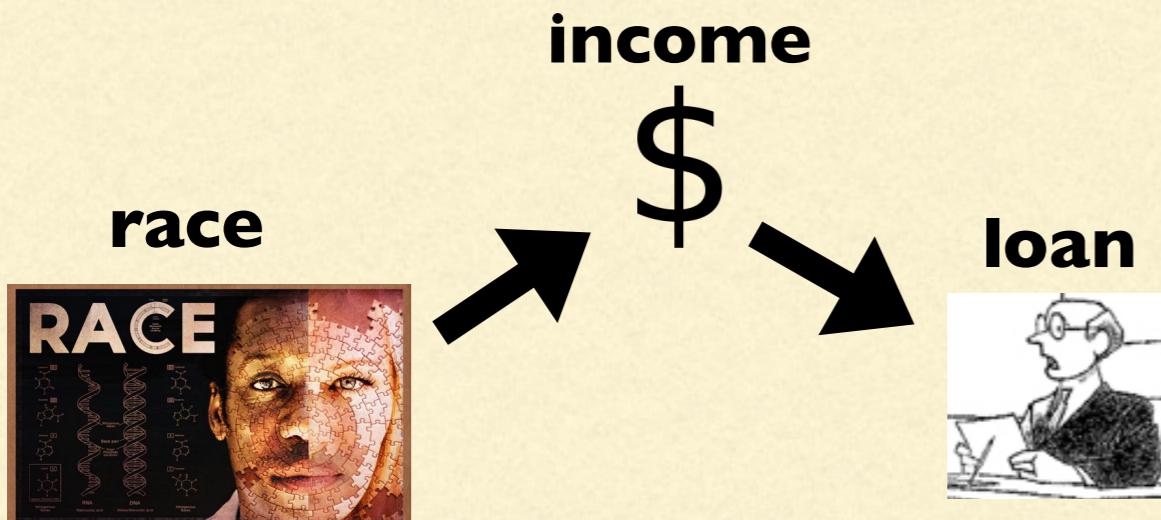
STATISTICS CAN'T MODEL HOW

“Race is *correlated* with loan decision!”



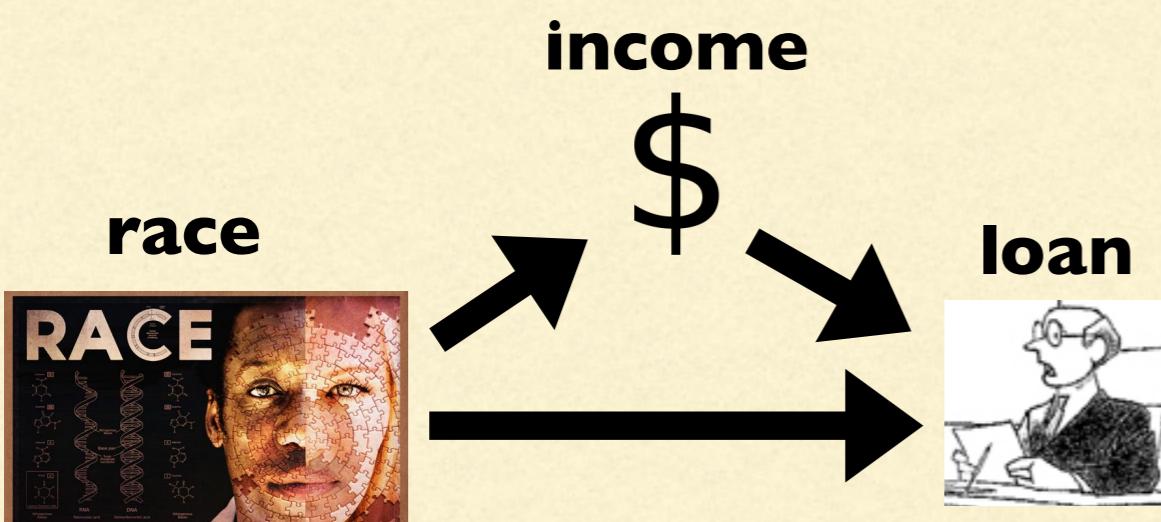
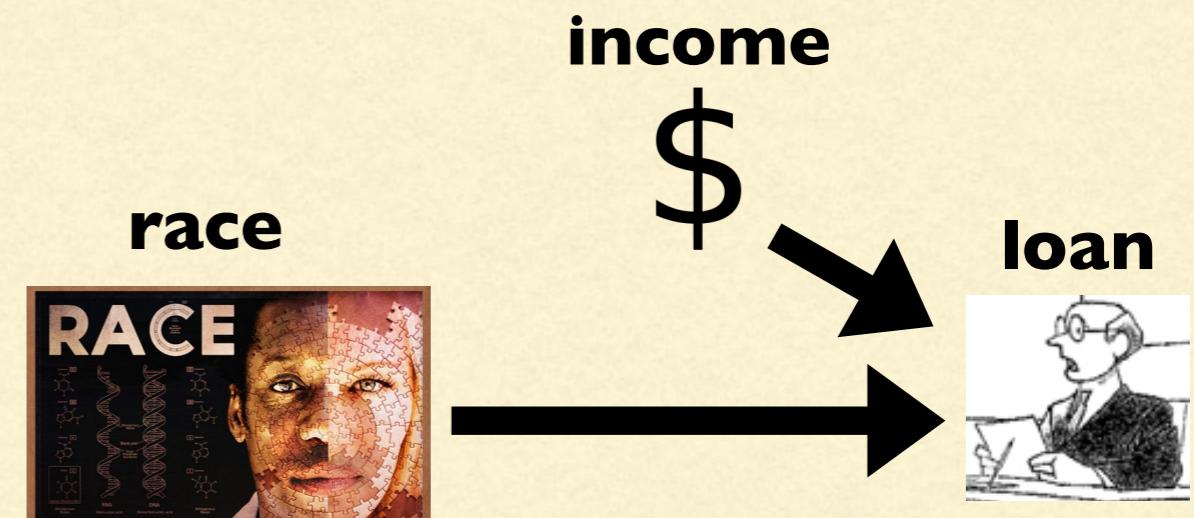
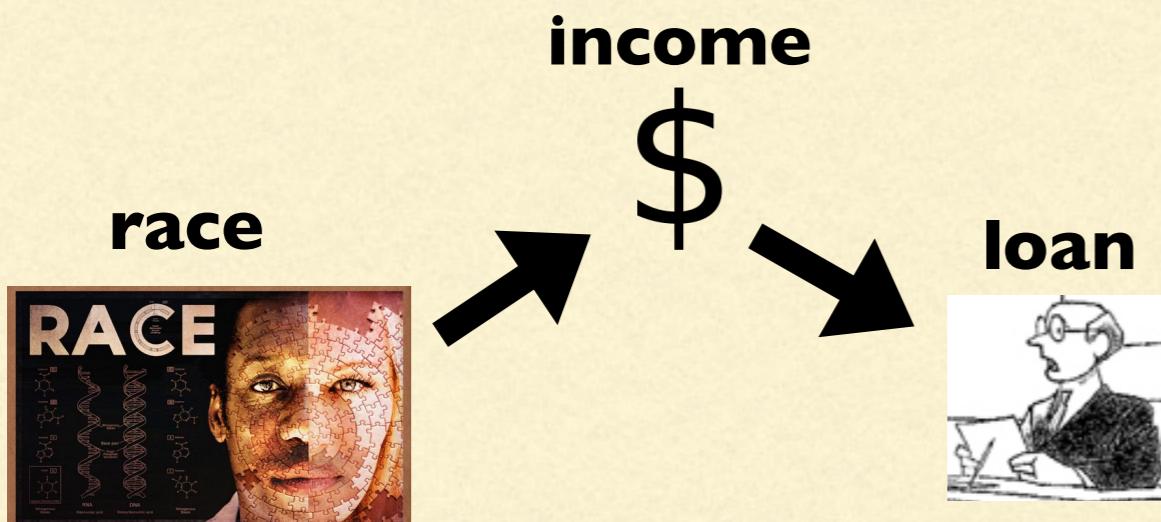
STATISTICS CAN'T MODEL HOW

“Race is *correlated* with loan decision!”



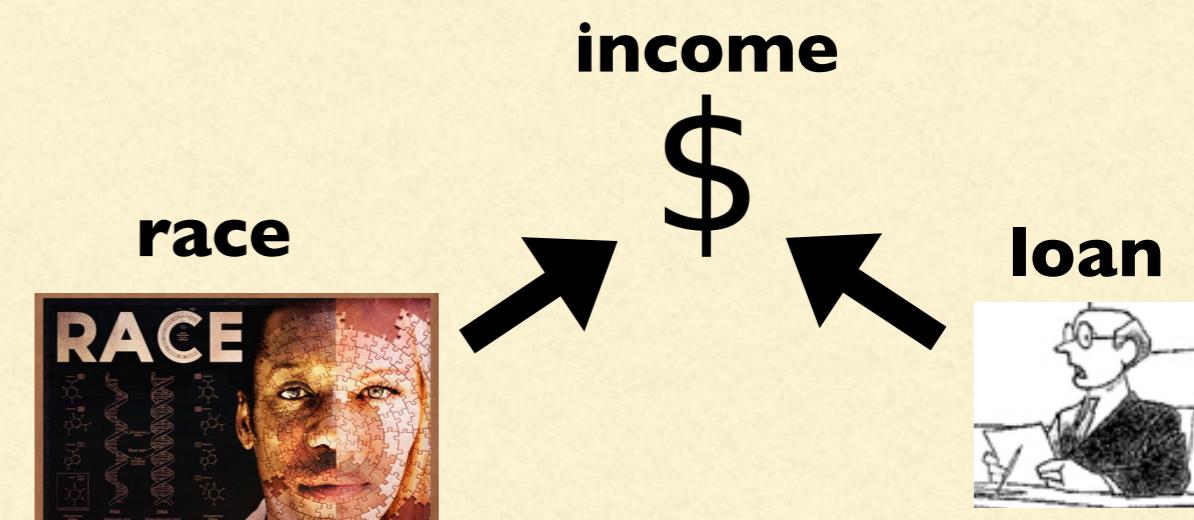
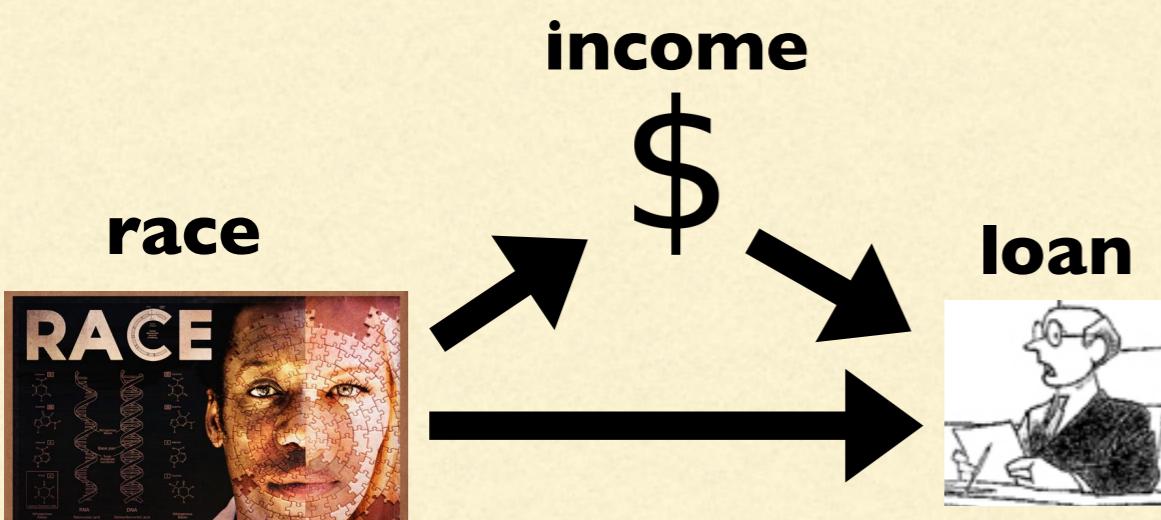
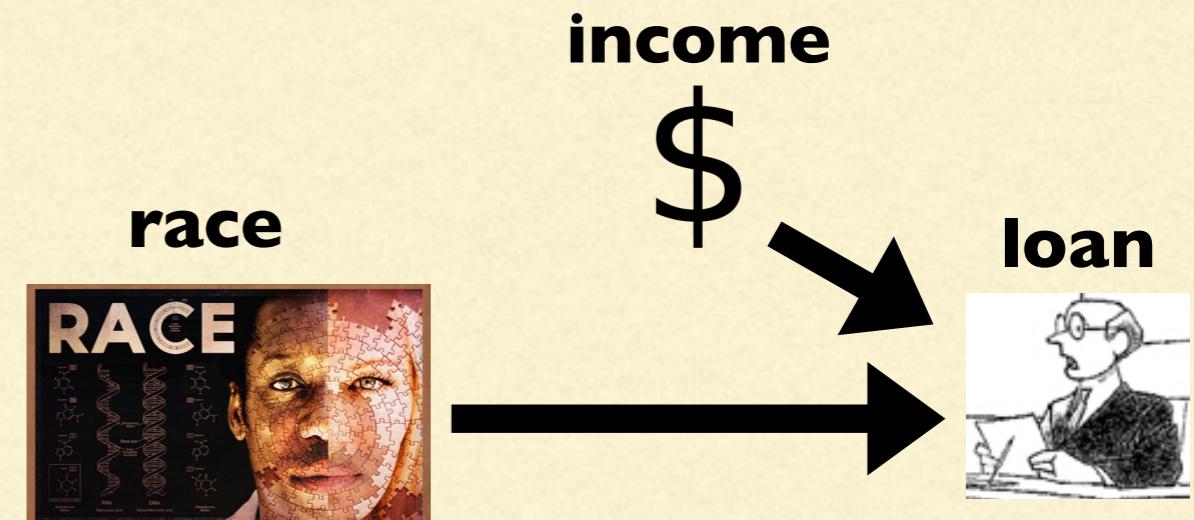
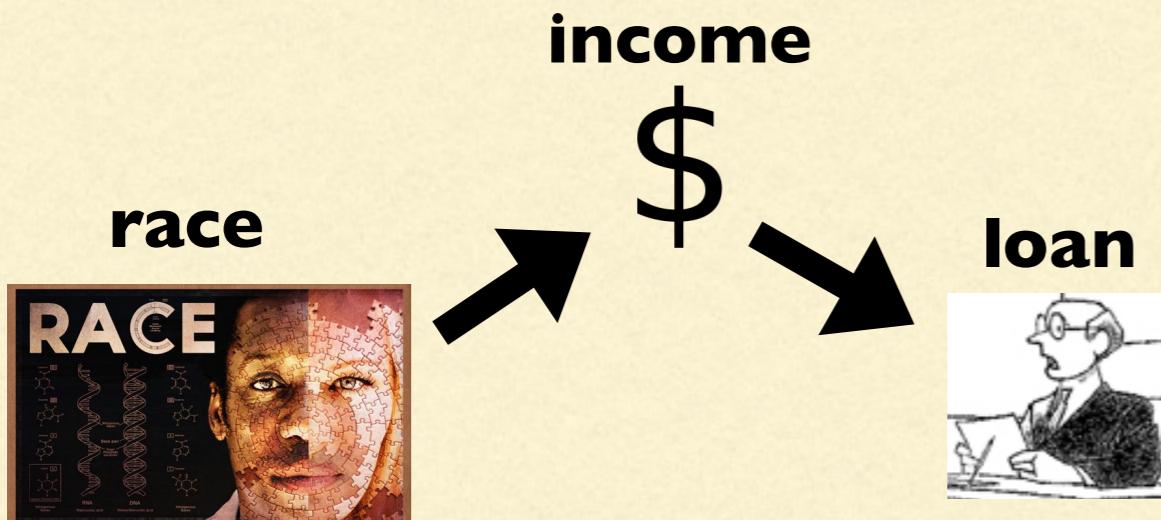
STATISTICS CAN'T MODEL HOW

“Race is *correlated* with loan decision!”



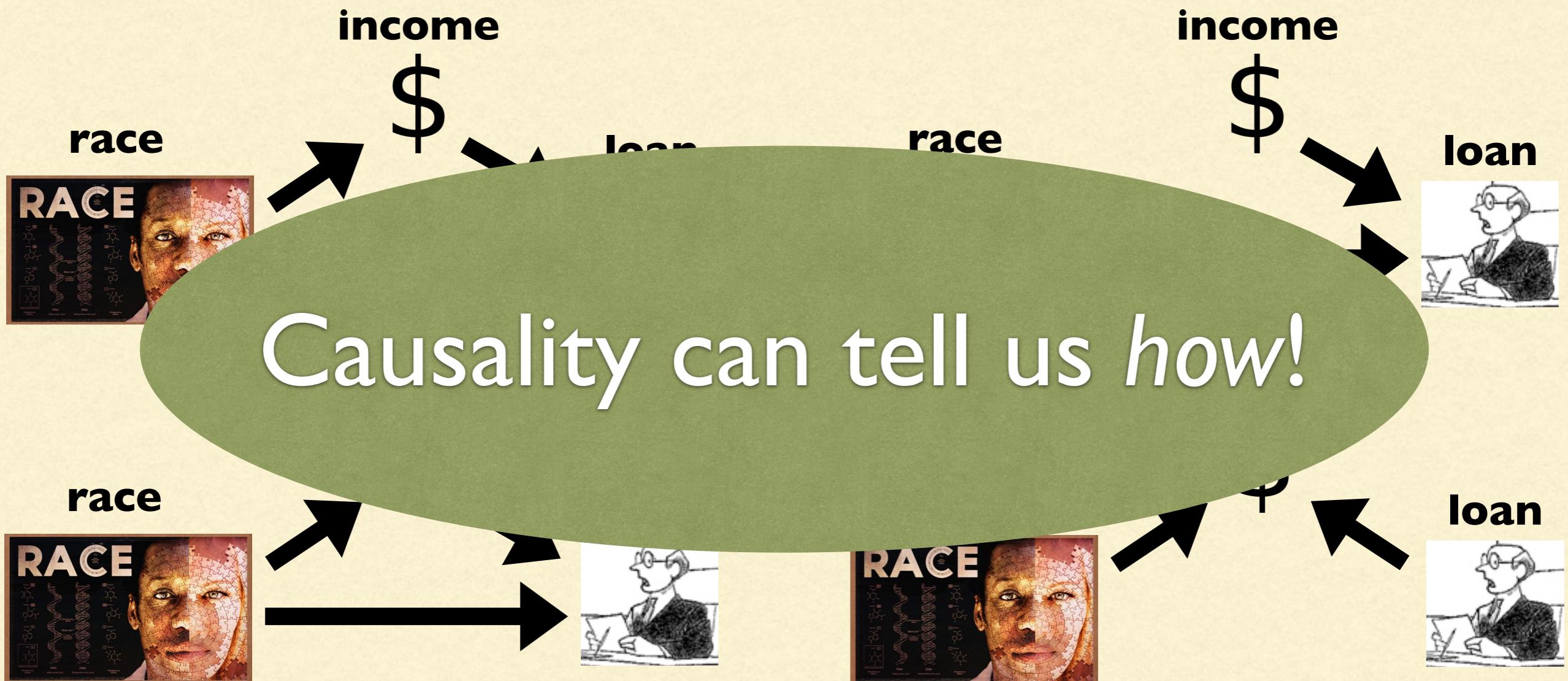
STATISTICS CAN'T MODEL HOW

“Race is correlated with loan decision!”

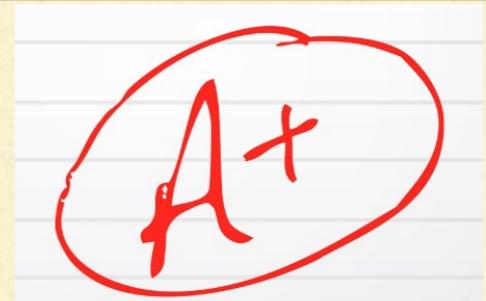
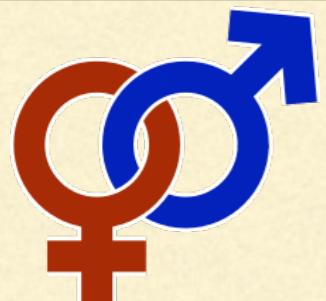


STATISTICS CAN'T MODEL HOW

“Race is *correlated* with loan decision!”



WHO SHOULD BE ADMITTED?



male

white



female

black



male

black

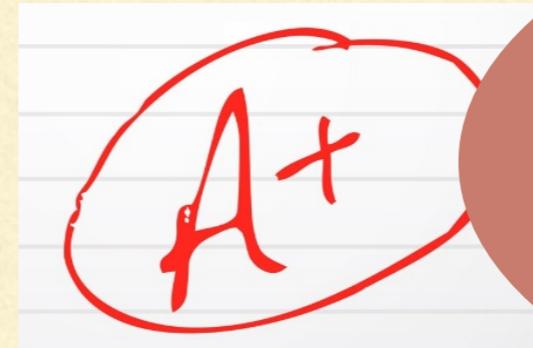
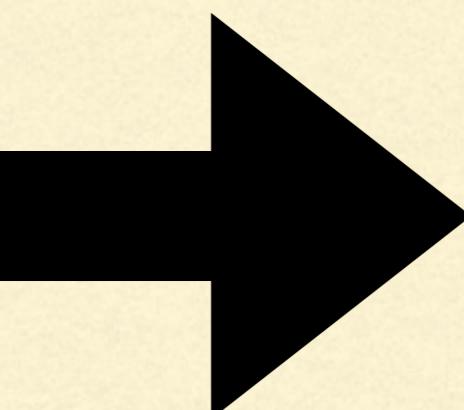
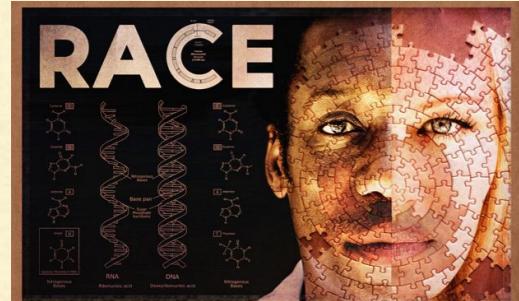
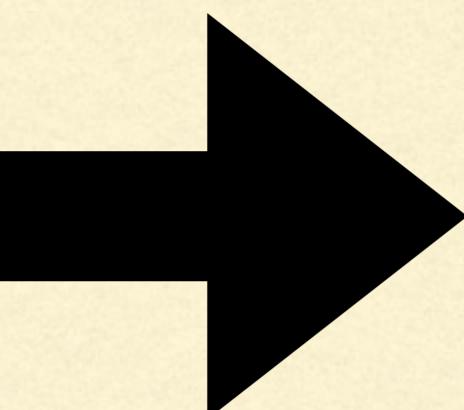
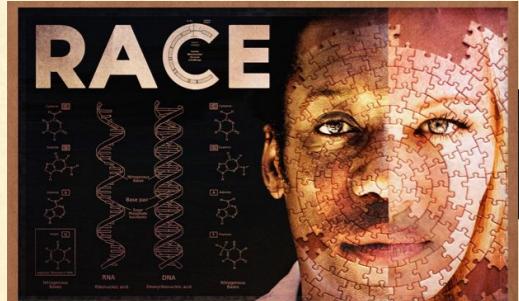
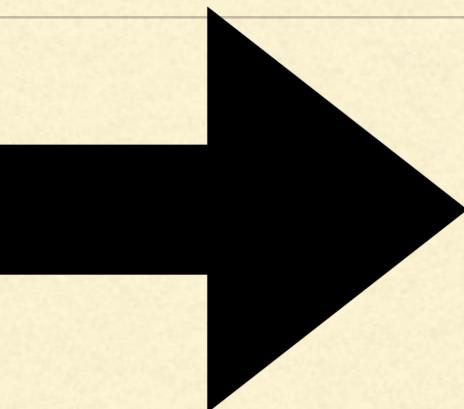
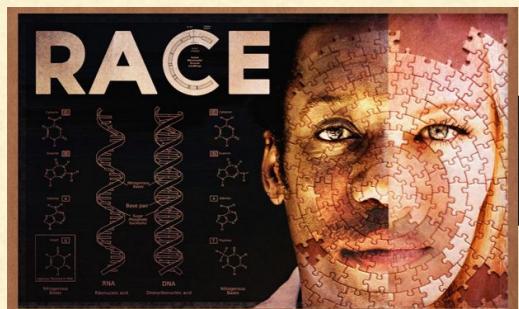


sensitive attributes A

features \mathcal{X}

label Y

UNFAIR INFLUENCES

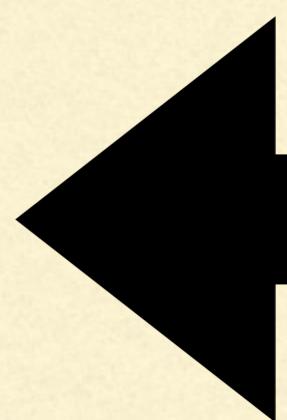
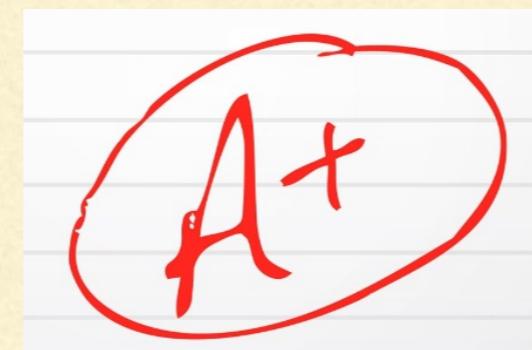
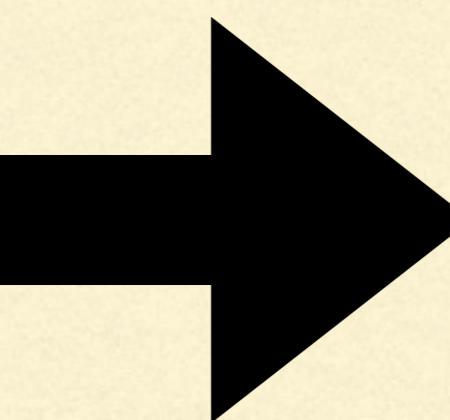
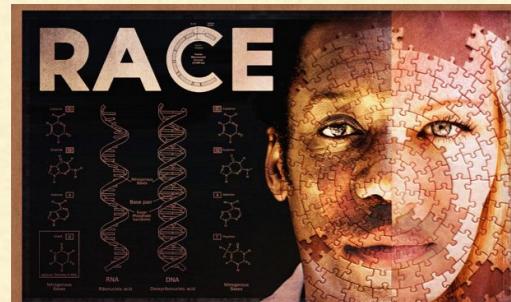
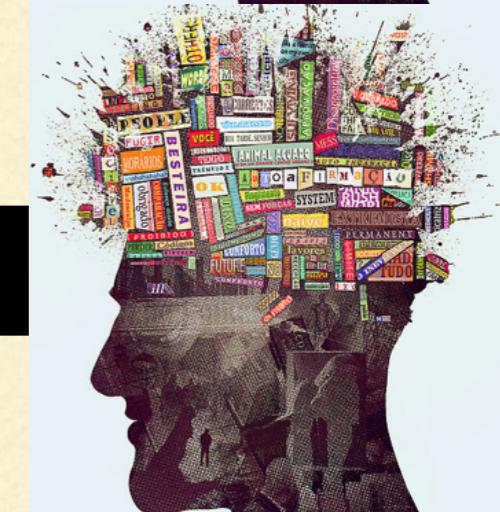
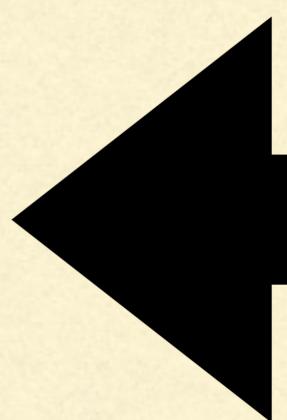
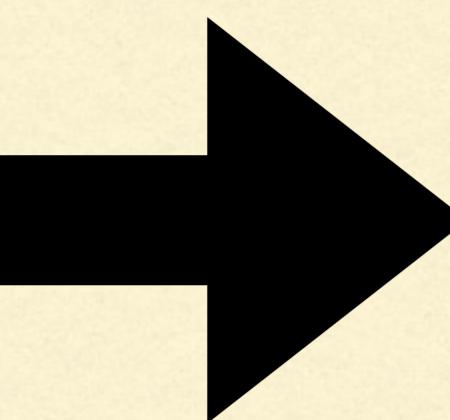
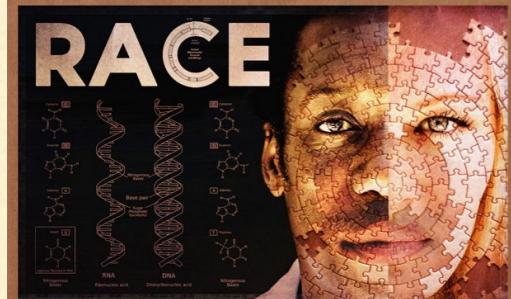
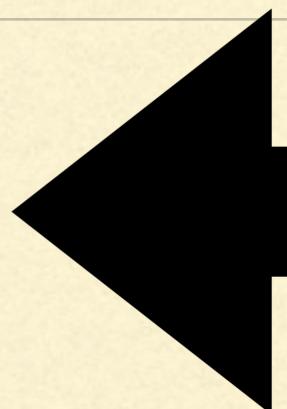
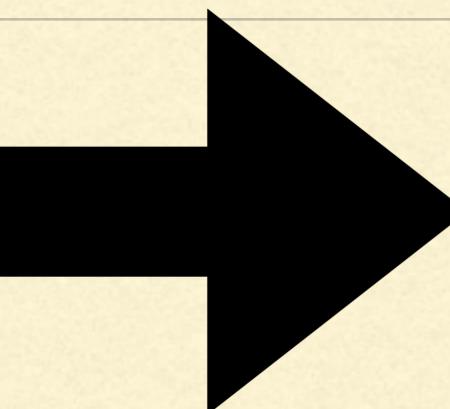
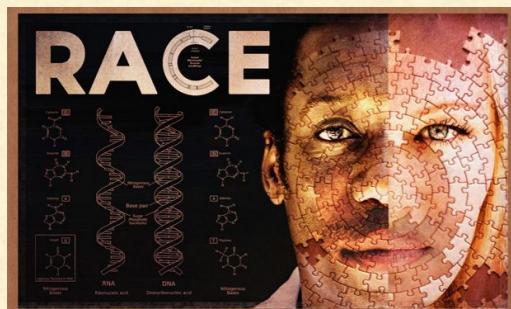


limited access to academic institutions
due to economic history
[Carter, 1973]

class placement may be **different**
[Howard, 2003]

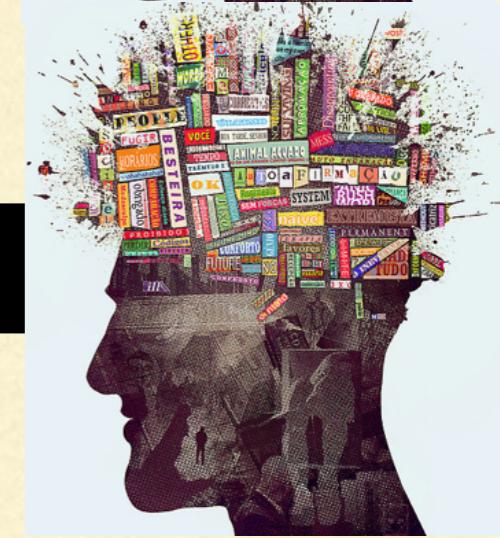
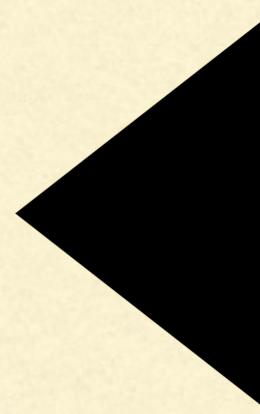
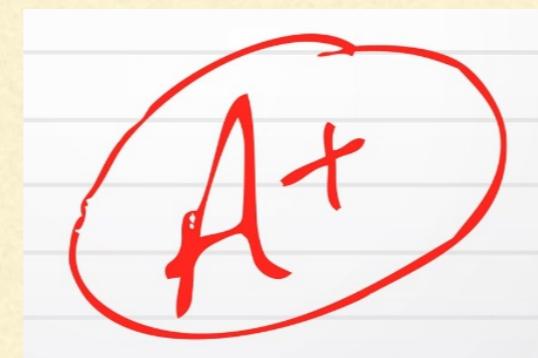
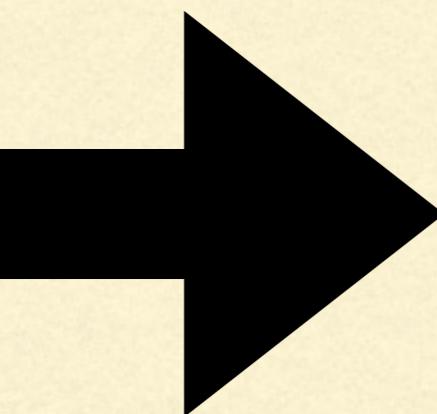
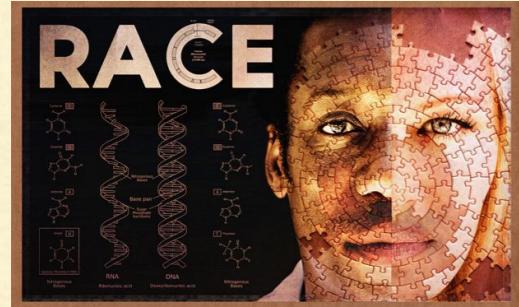
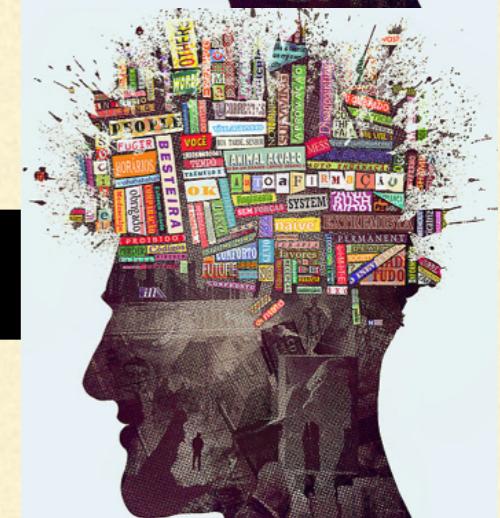
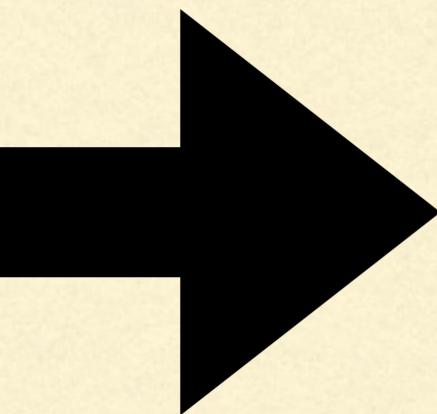
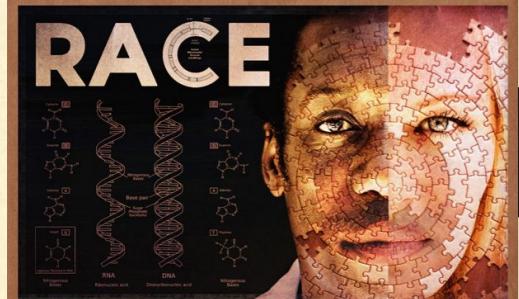
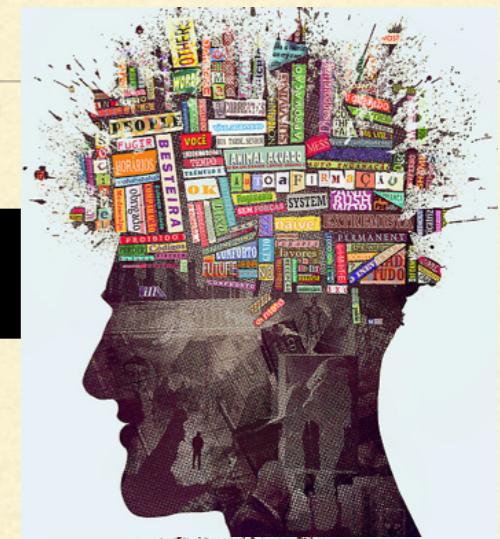
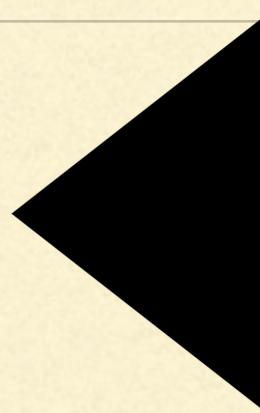
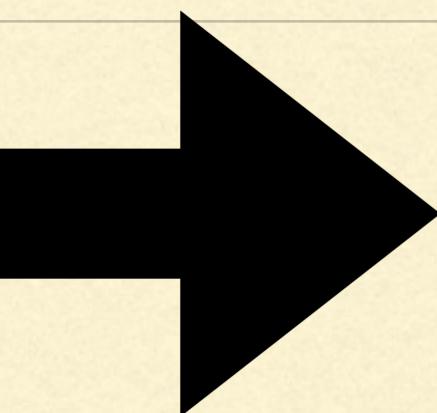
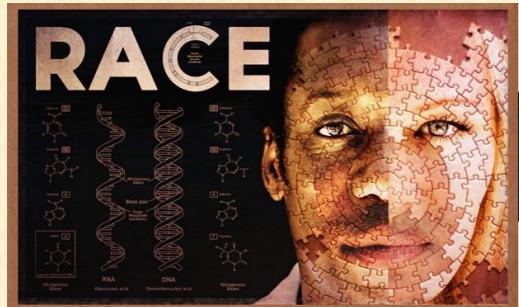
lack of minority race teachers **affects** minority student outcomes
[Birdsall et al., 2016]

UNFAIR INFLUENCES



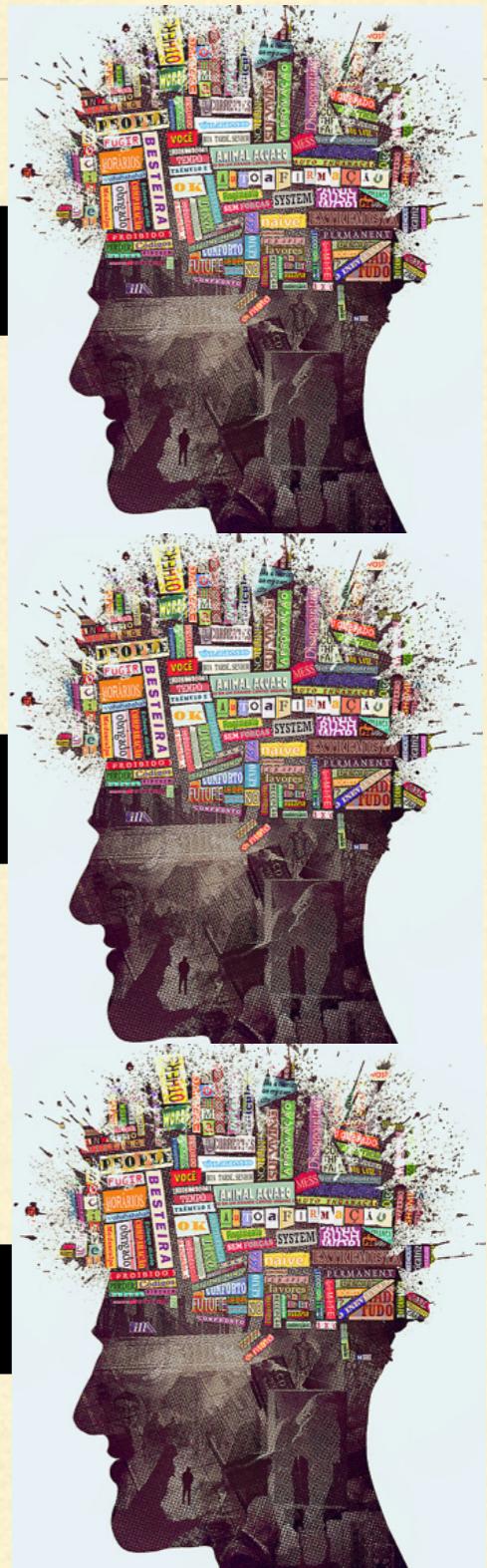
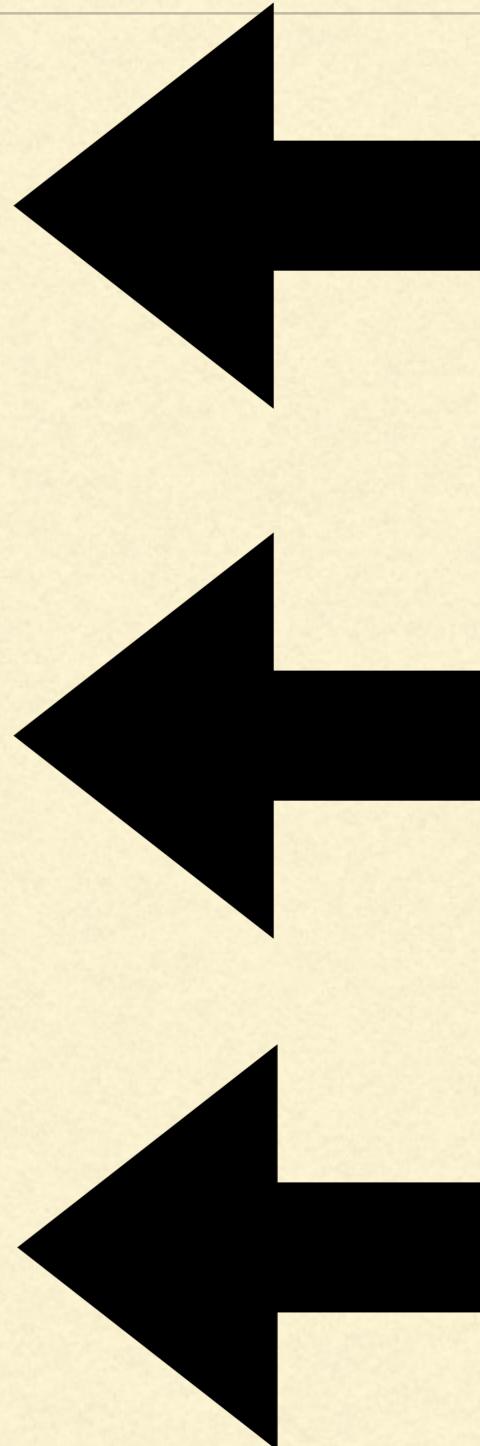
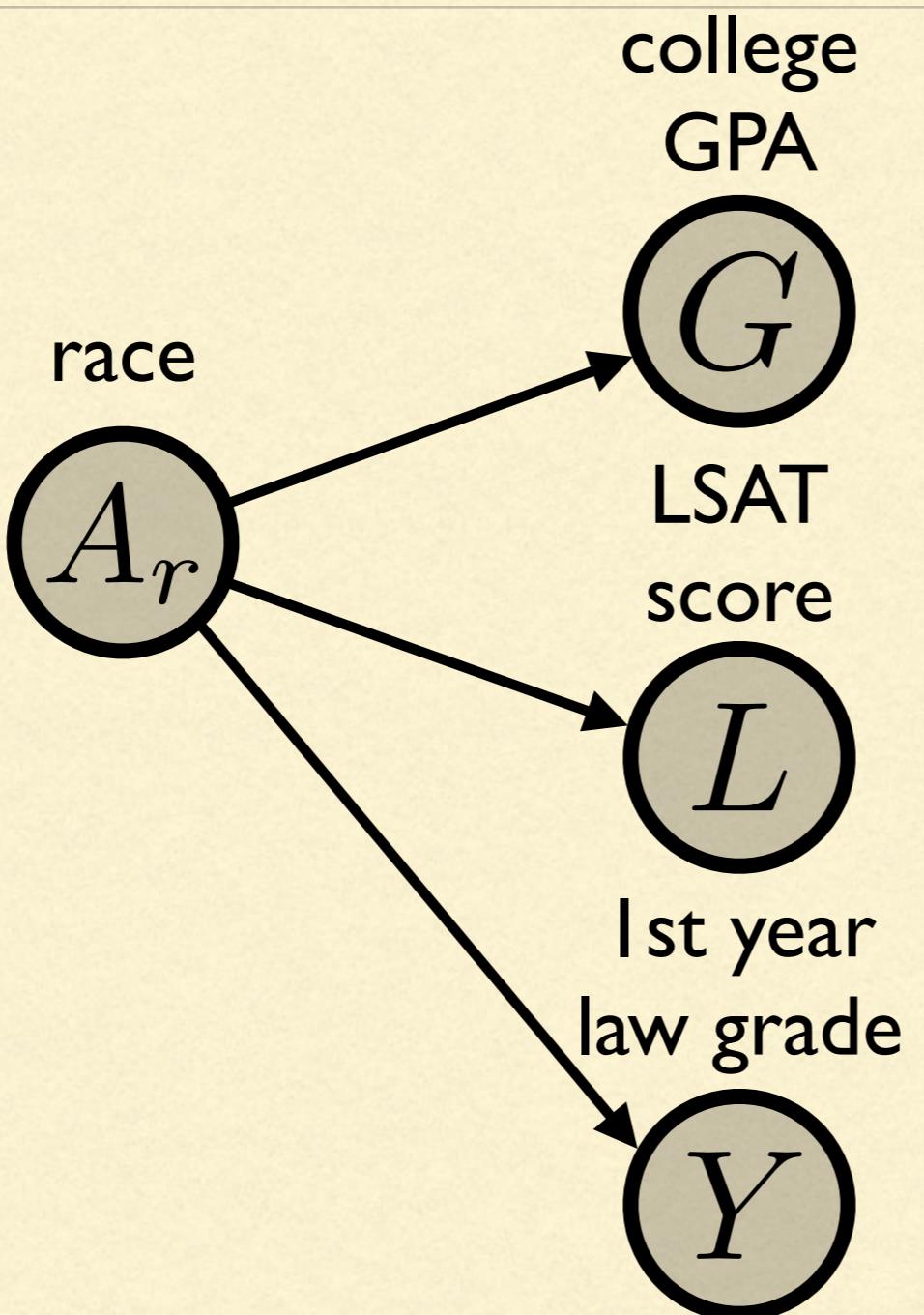
CAUSALITY

[Pearl, 2000; Pearl, 2009; Pearl et al., 2016]



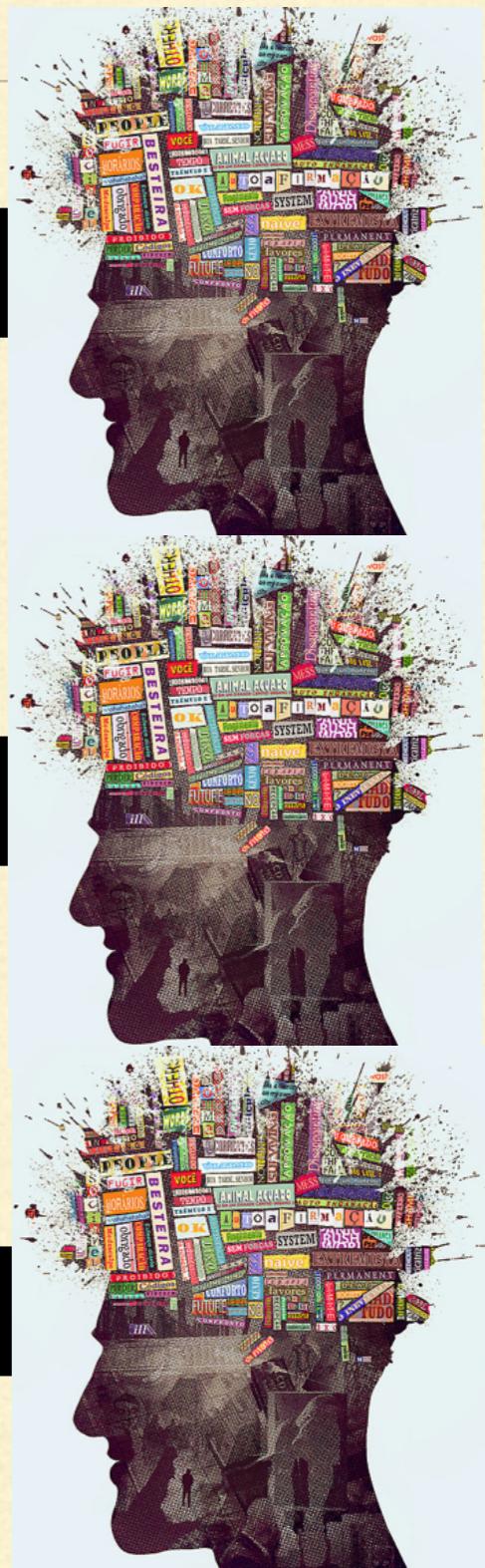
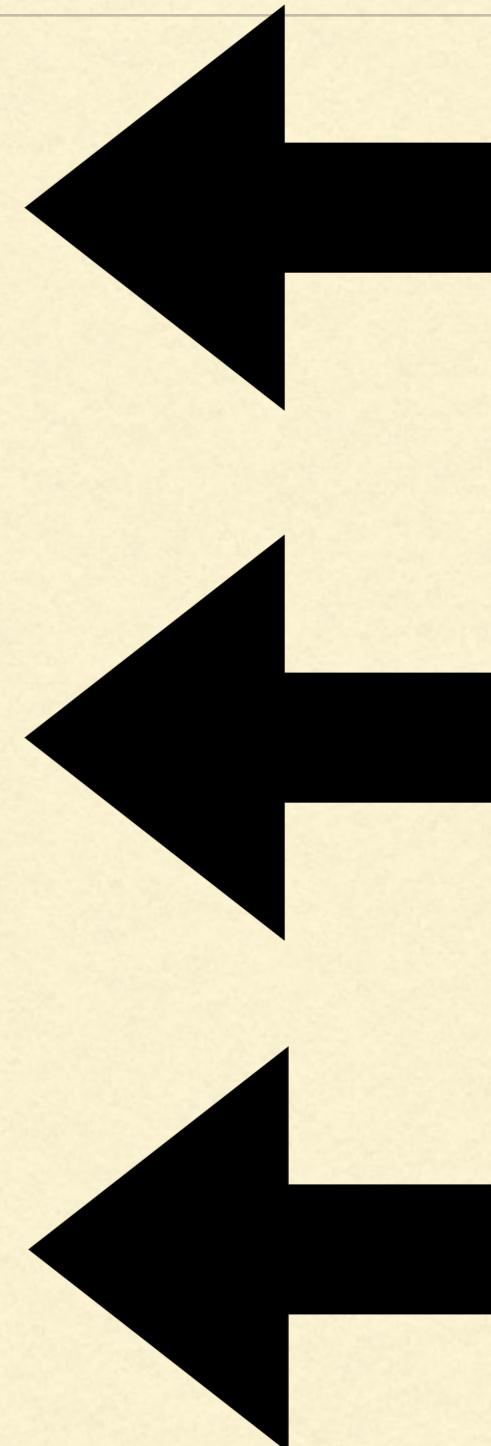
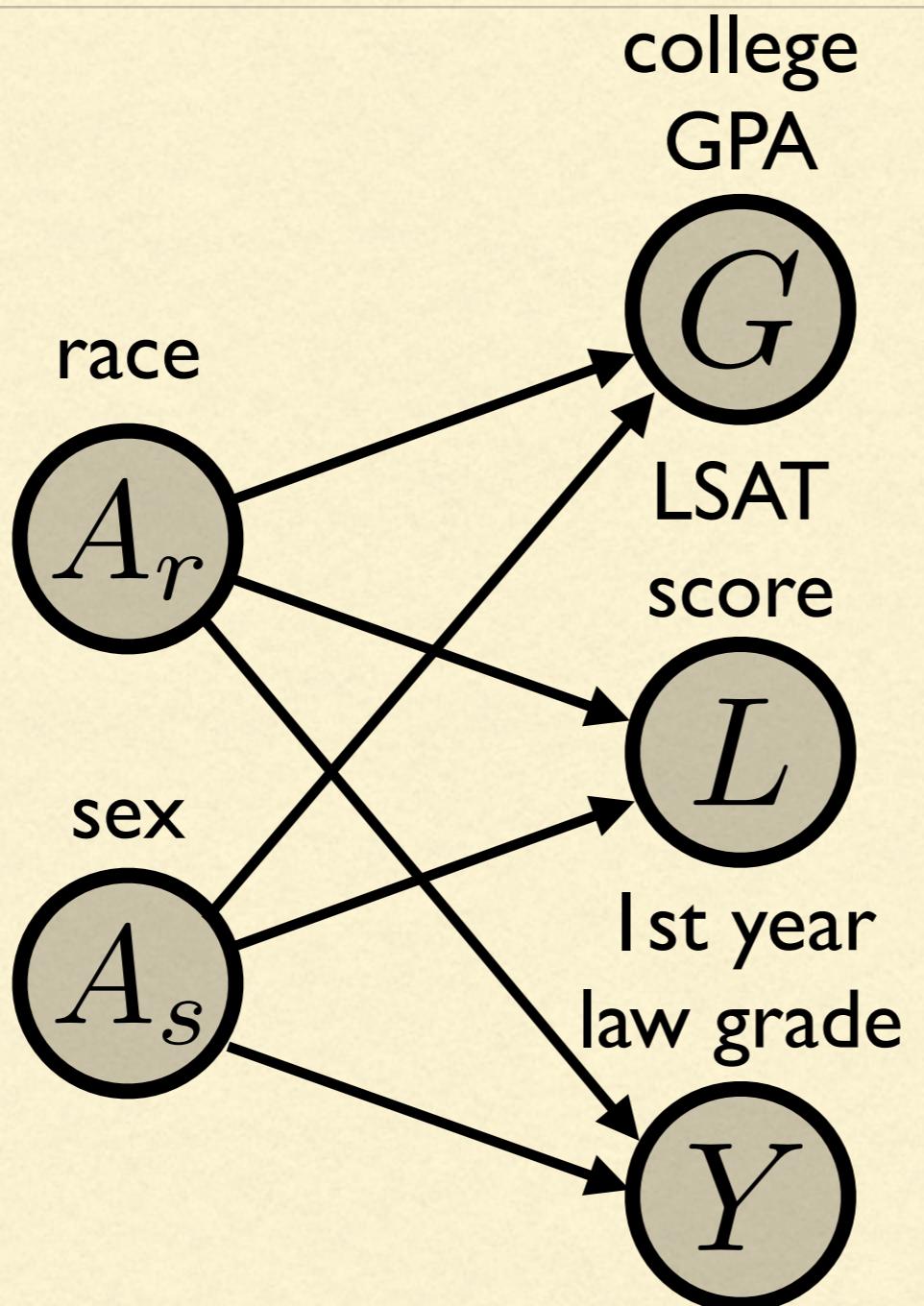
CAUSALITY

[Pearl, 2000; Pearl, 2009; Pearl et al., 2016]



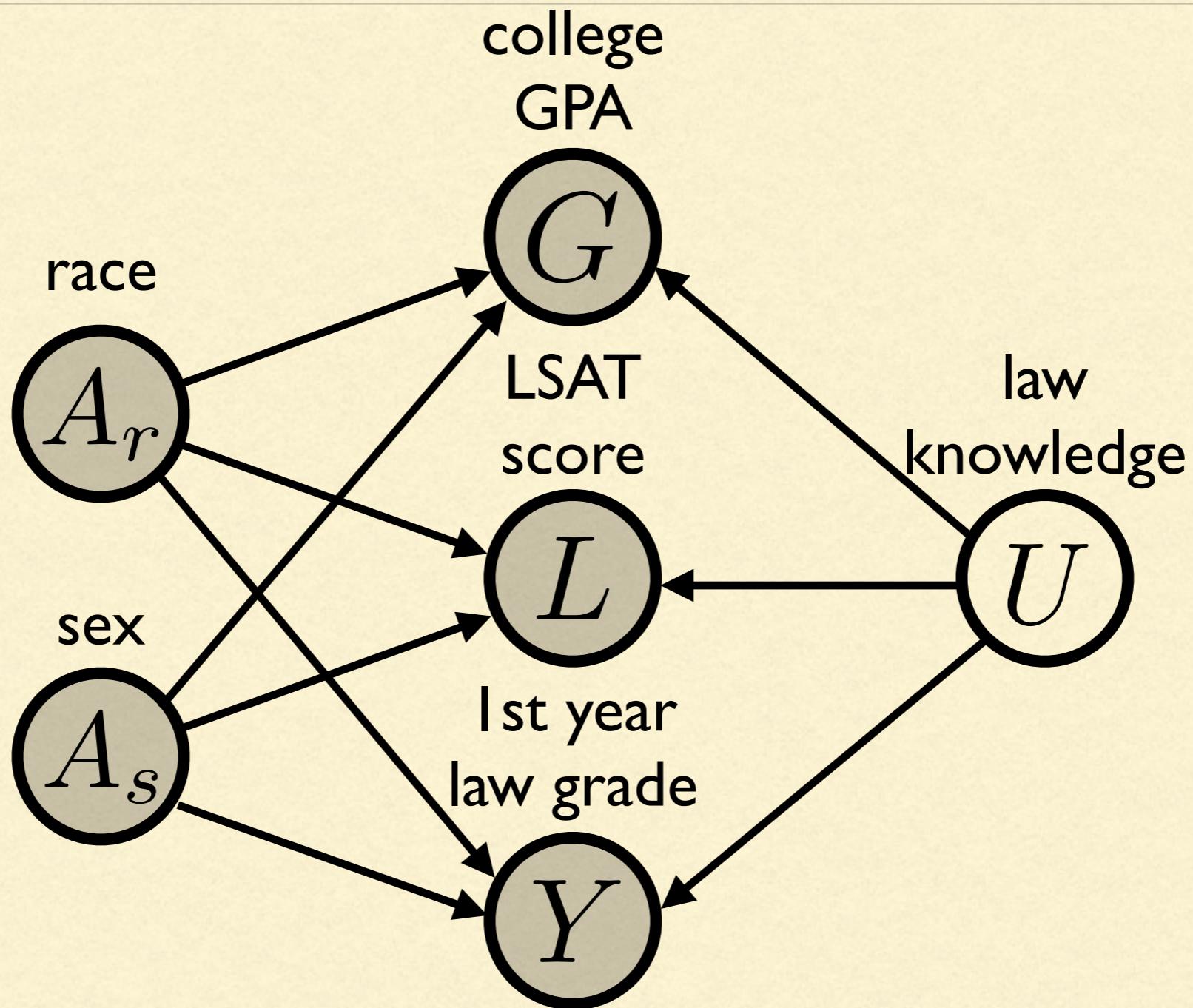
CAUSALITY

[Pearl, 2000; Pearl, 2009; Pearl et al., 2016]



CAUSALITY

[Pearl, 2000; Pearl, 2009; Pearl et al., 2016]



COUNTERFACTUALS

How would someone have been different if they had been a different race/sex?

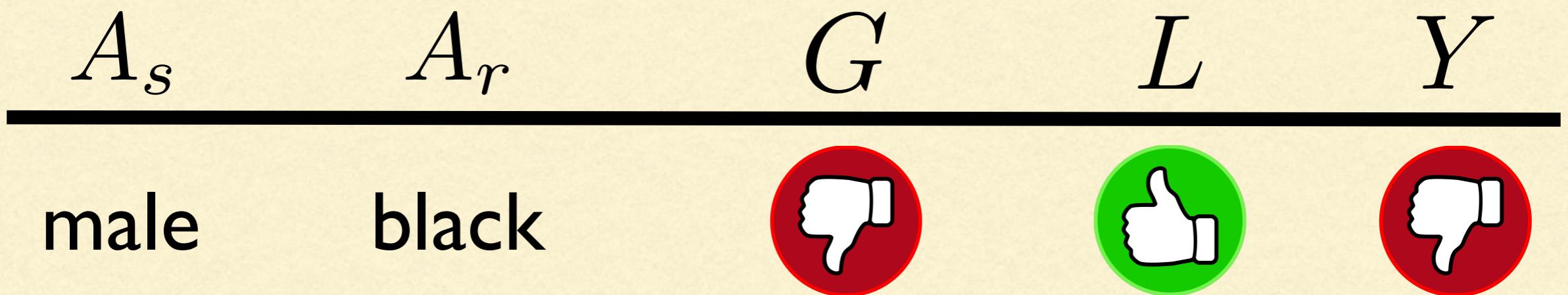
COUNTERFACTUALS

[Pearl, 2000; Pearl, 2009; Pearl et al., 2016]



3-STEP PROCEDURE

[Pearl, 2000; Pearl, 2009; Pearl et al., 2016]



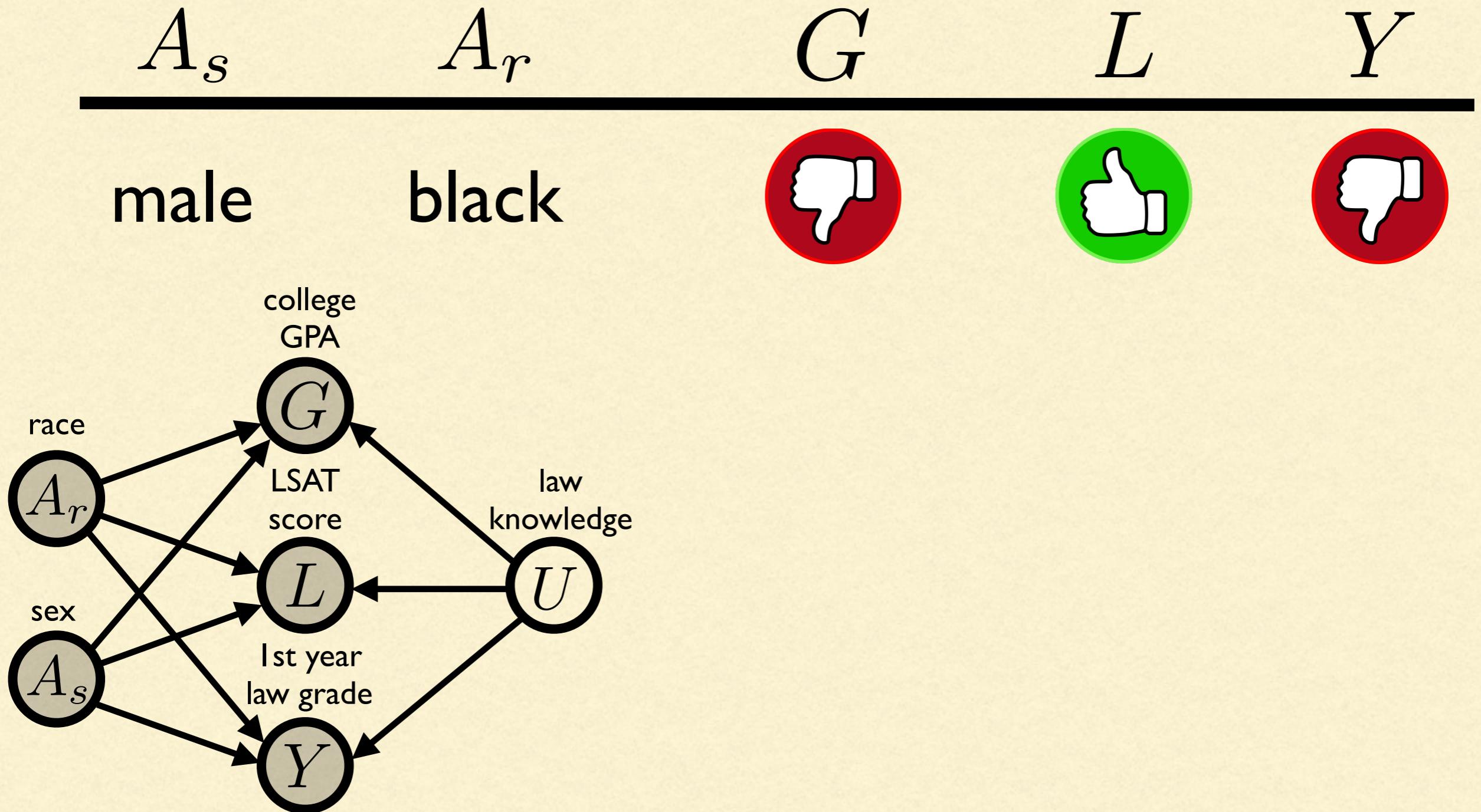
I. COMPUTE UNOBSERVED VARIABLES IN CAUSAL MODEL

[Pearl, 2000; Pearl, 2009; Pearl et al., 2016]

A_s	A_r	G	L	Y
male	black			

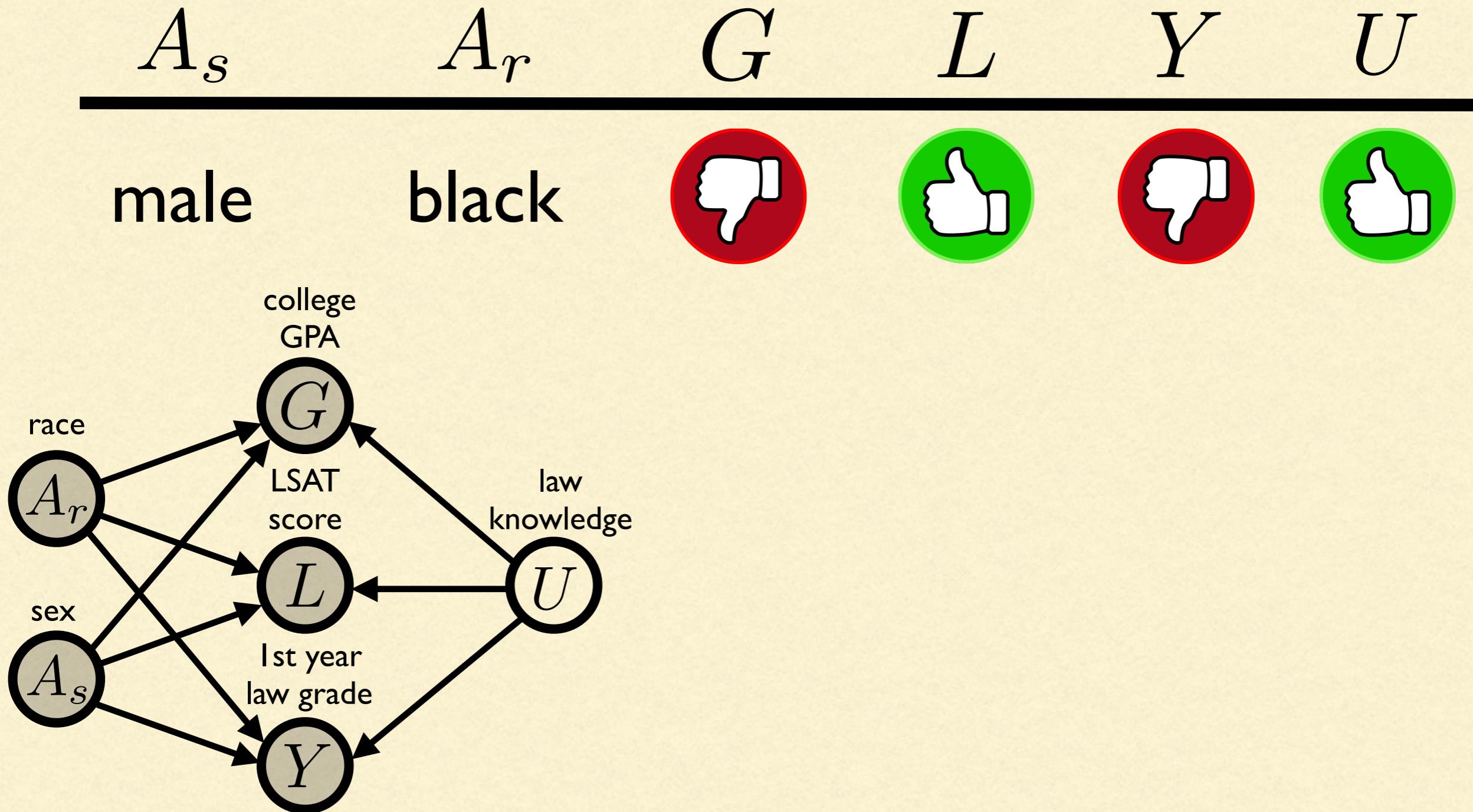
I. COMPUTE UNOBSERVED VARIABLES IN CAUSAL MODEL

[Pearl, 2000; Pearl, 2009; Pearl et al., 2016]



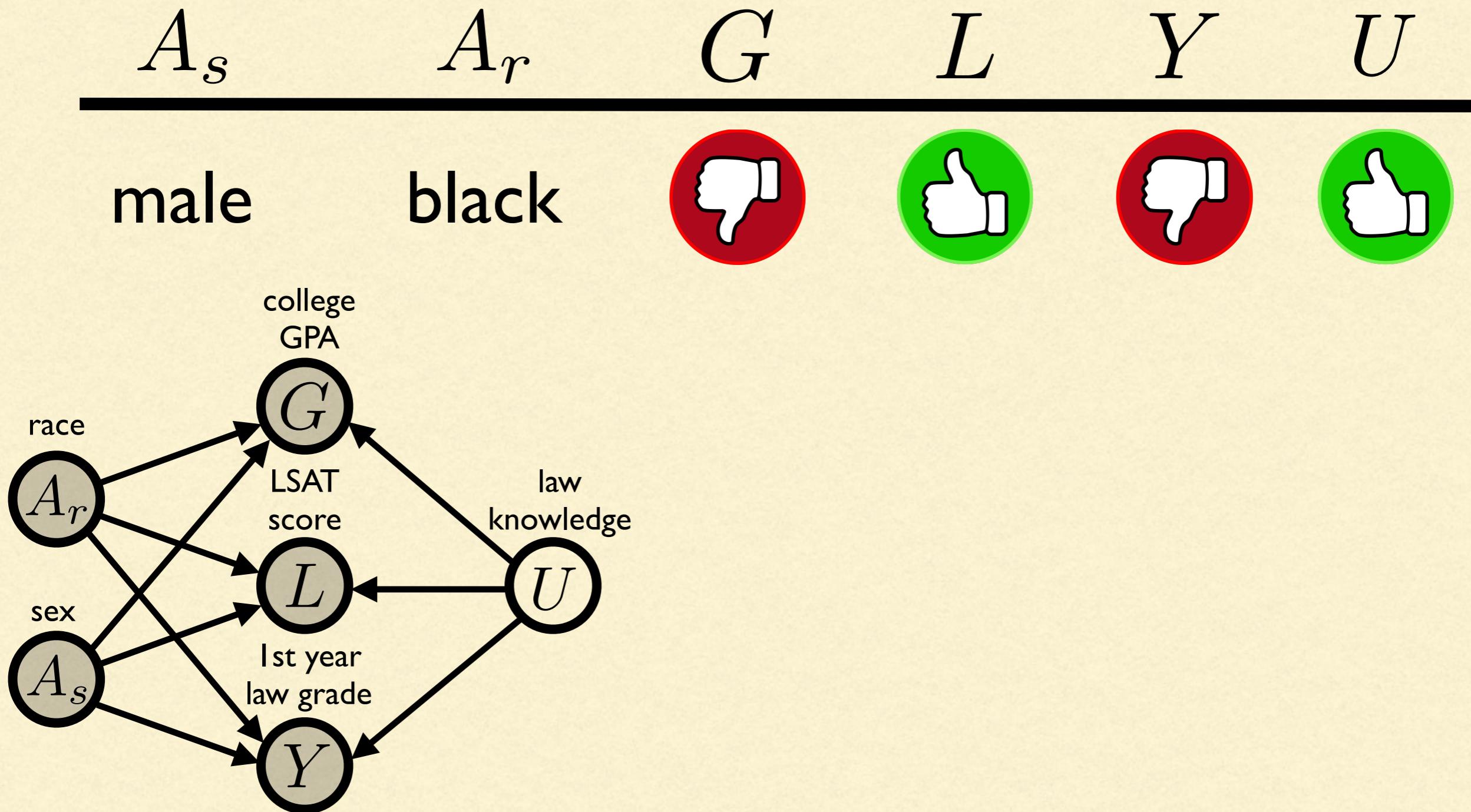
I. COMPUTE UNOBSERVED VARIABLES IN CAUSAL MODEL

[Pearl, 2000; Pearl, 2009; Pearl et al., 2016]



2. CHANGE FACTUAL A

[Pearl, 2000; Pearl, 2009; Pearl et al., 2016]



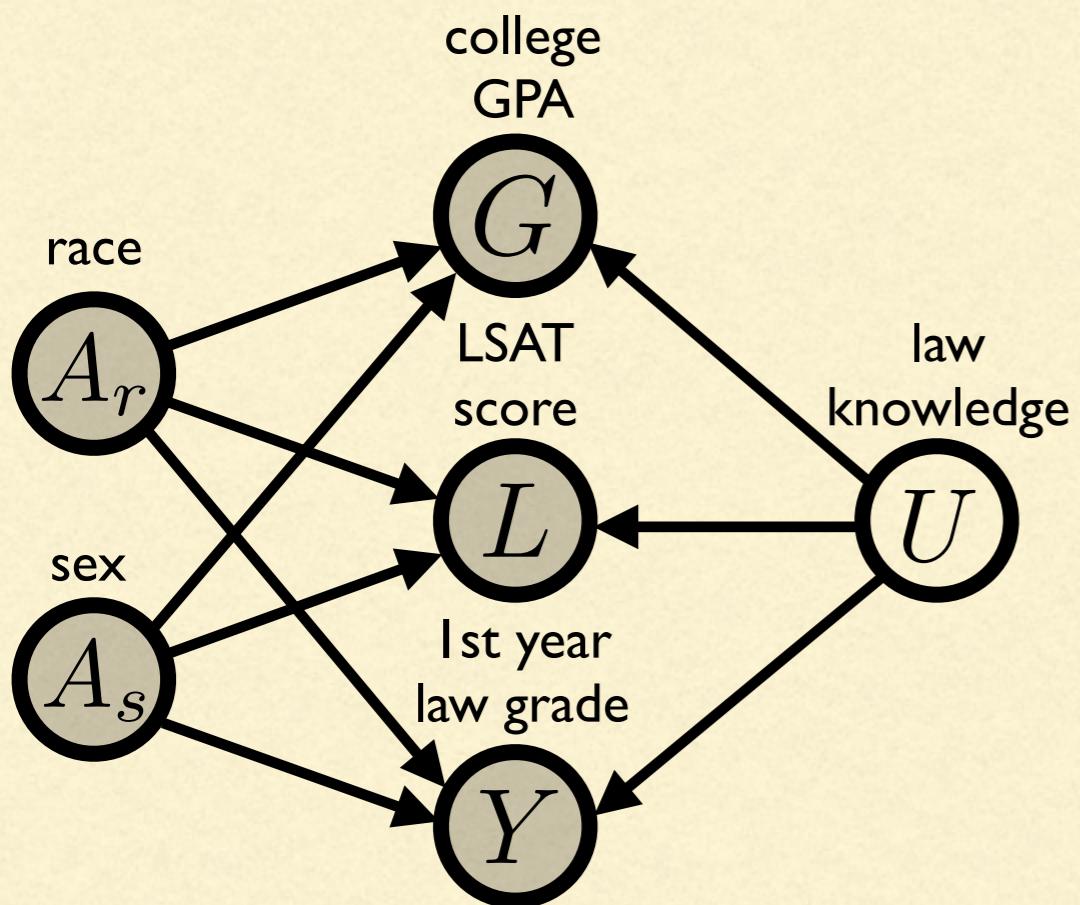
2. CHANGE FACTUAL A

[Pearl, 2000; Pearl, 2009; Pearl et al., 2016]

$$A_s \quad A_r \leftarrow a' \quad G \quad L \quad Y \quad U$$

male

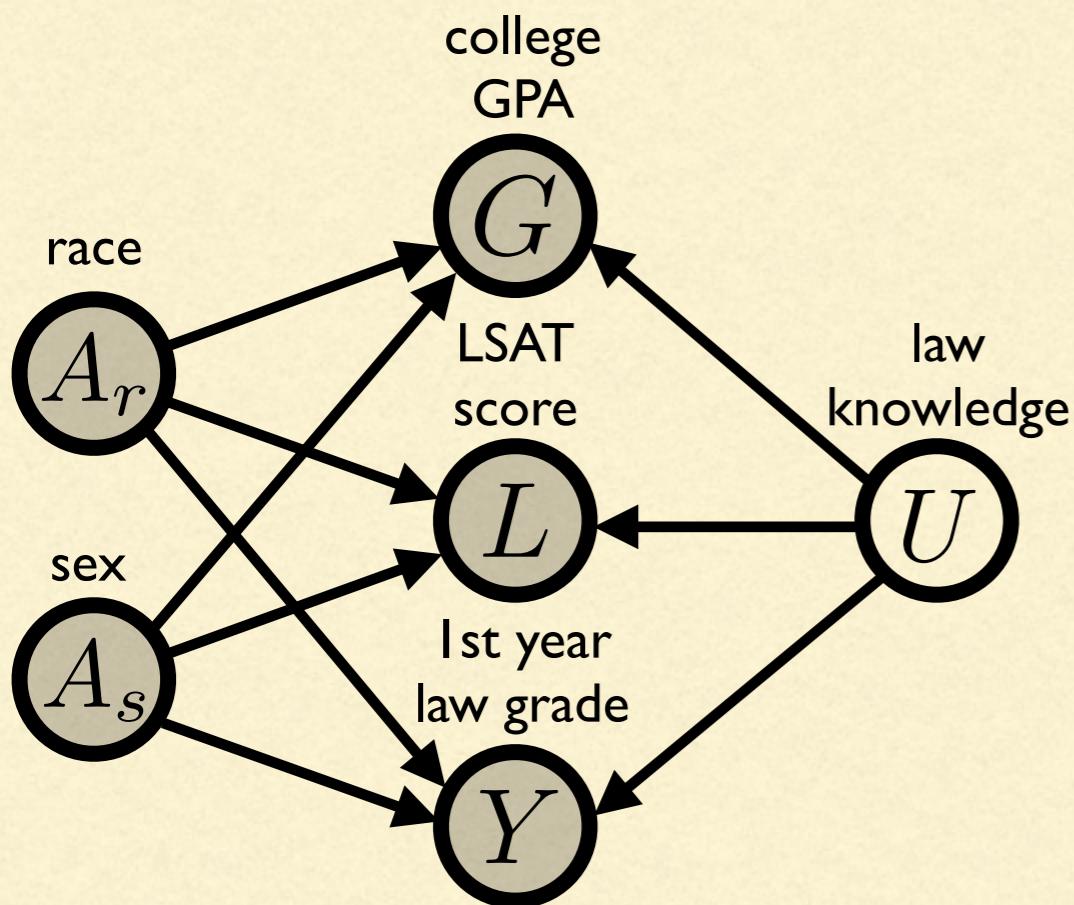
white



3. RECOMPUTE OBSERVED VARIABLES IN CAUSAL MODEL

[Pearl, 2000; Pearl, 2009; Pearl et al., 2016]

$$\frac{A_s \quad A_r \leftarrow a' \quad G \quad L \quad Y \quad U}{\text{male} \quad \text{white}}$$

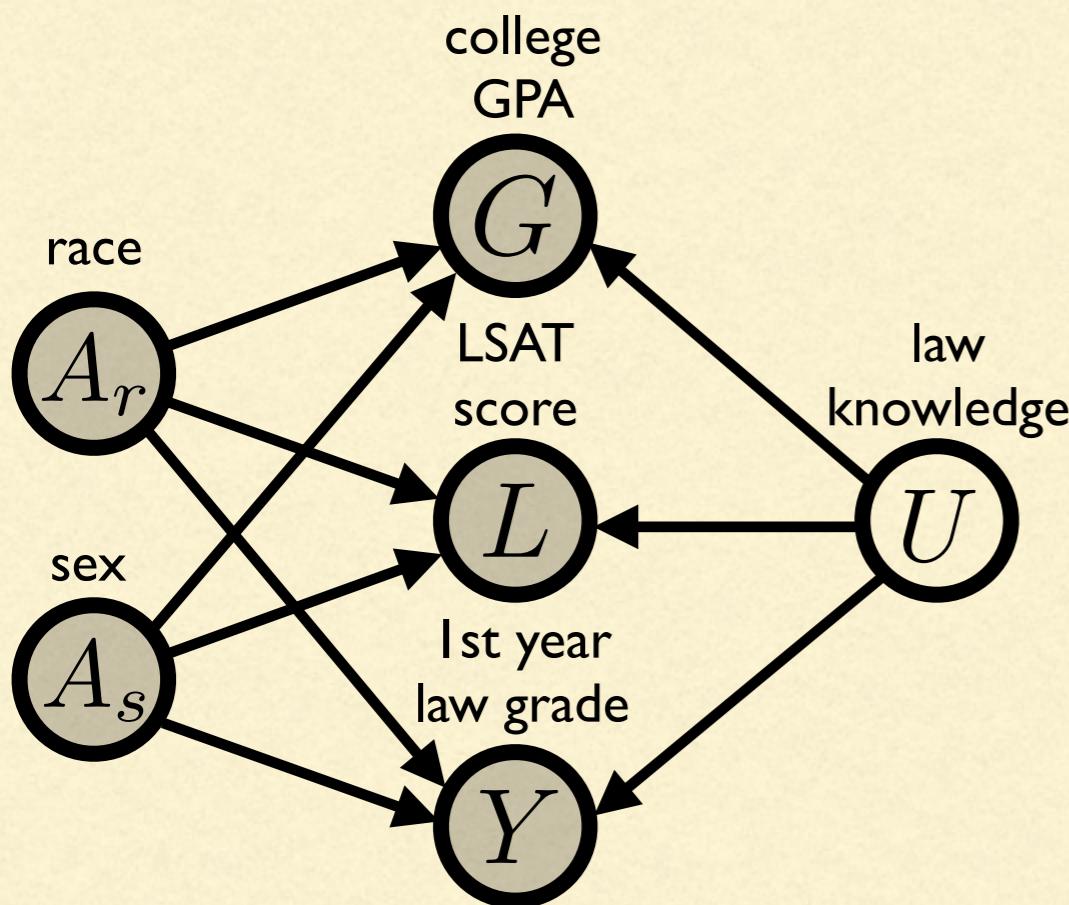


3. RECOMPUTE OBSERVED VARIABLES IN CAUSAL MODEL

[Pearl, 2000; Pearl, 2009; Pearl et al., 2016]

$$A_s \ A_r \leftarrow a' G_{A_r \leftarrow a'} \ L_{A_r \leftarrow a'} \ Y_{A_r \leftarrow a'} \ U$$

male white

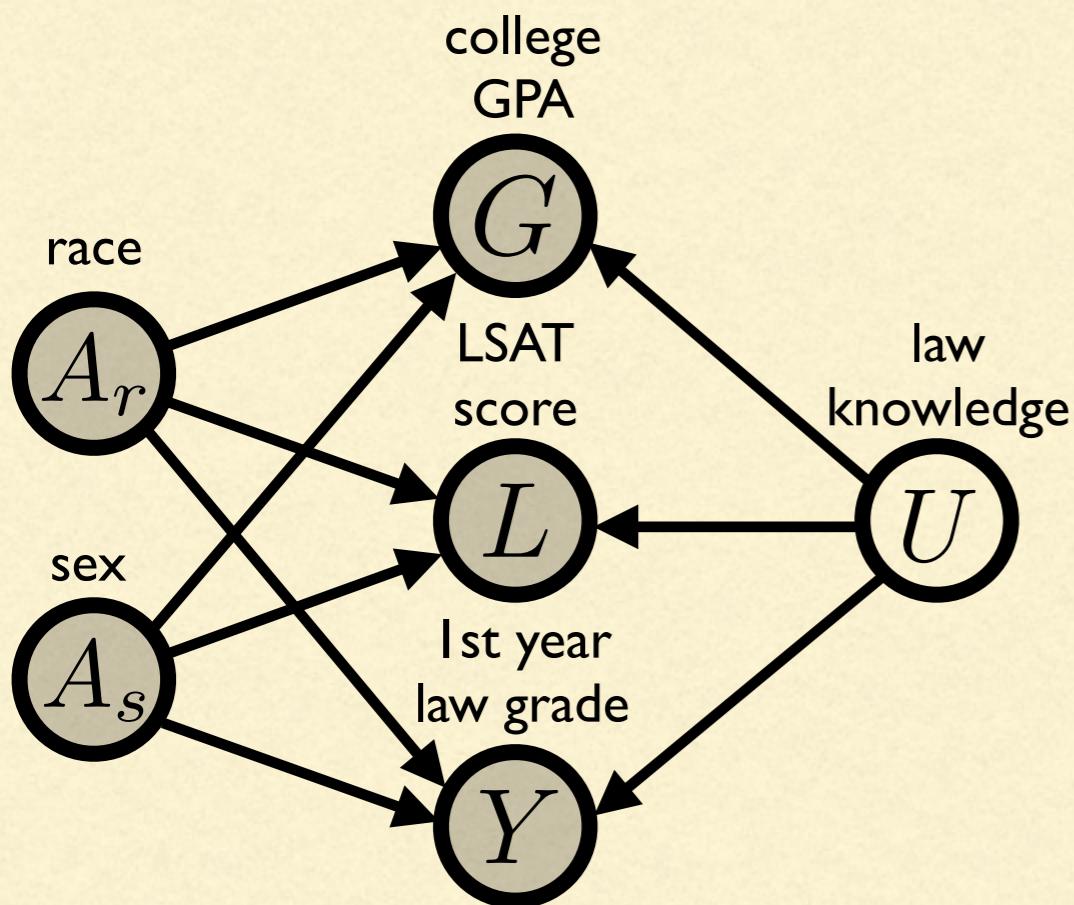


3. RECOMPUTE OBSERVED VARIABLES IN CAUSAL MODEL

[Pearl, 2000; Pearl, 2009; Pearl et al., 2016]

$$A_s \ A_r \leftarrow a' G_{A_r \leftarrow a'} \ L_{A_r \leftarrow a'} \ Y_{A_r \leftarrow a'} \ U$$

male white



COUNTERFACTUAL FAIRNESS

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

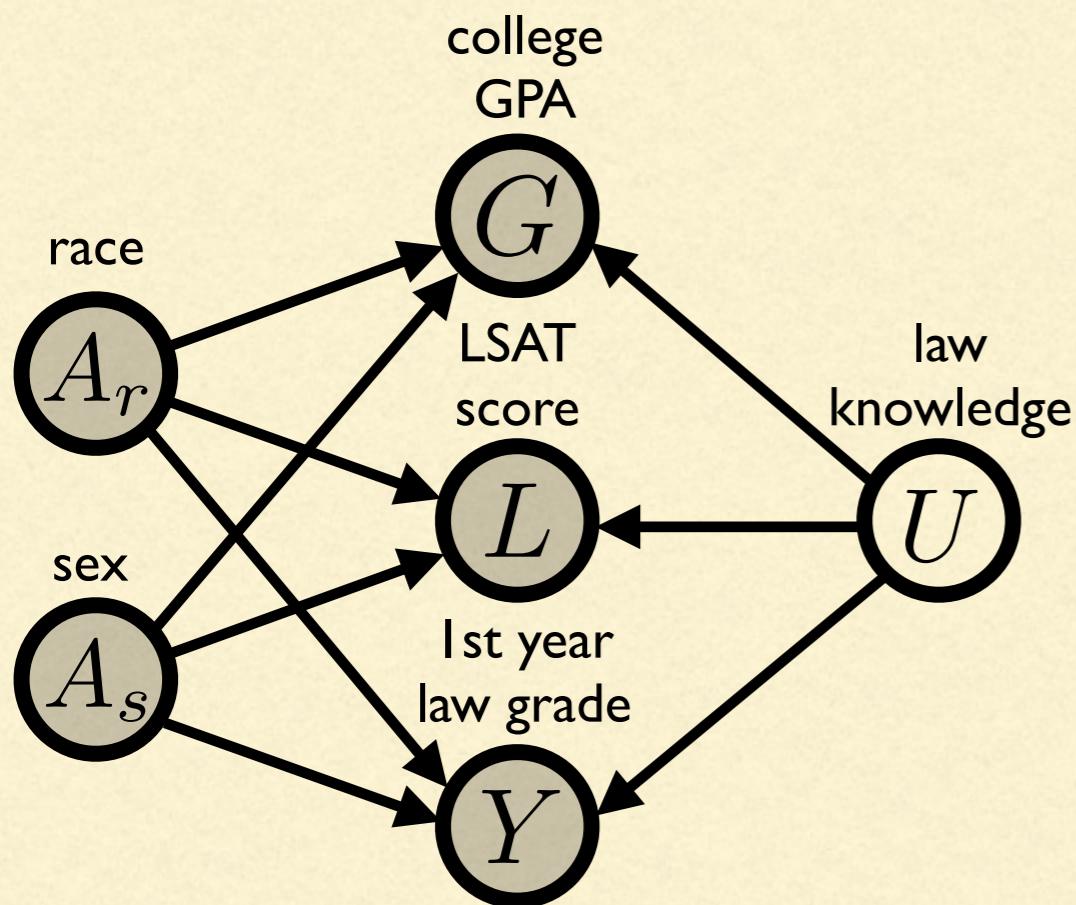
IS THIS CLASSIFIER FAIR?

$$A_s \ A_r \leftarrow a' G_{A_r \leftarrow a'} \ L_{A_r \leftarrow a'} \ Y_{A_r \leftarrow a'} \ U$$

male white

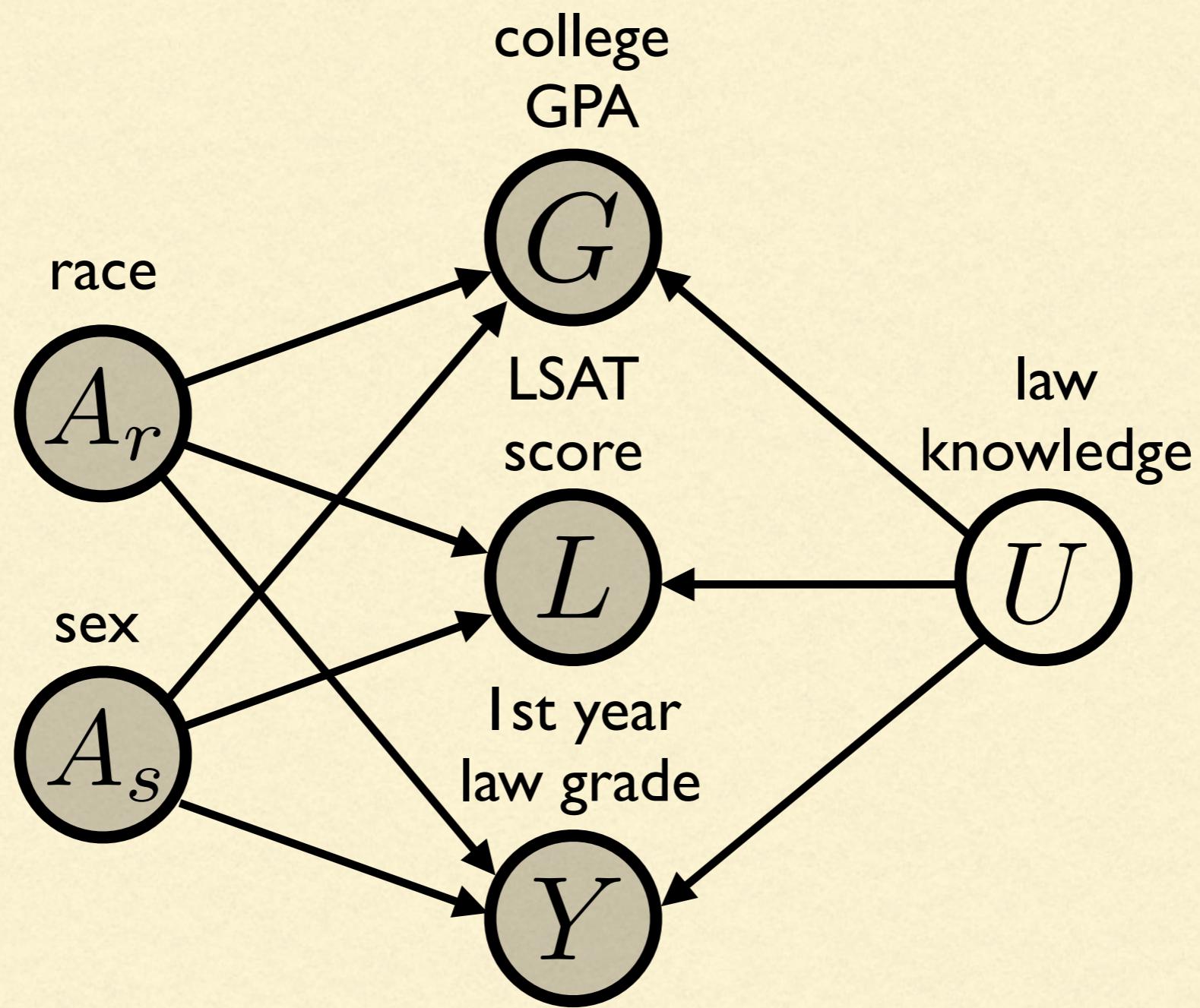


fair classifier

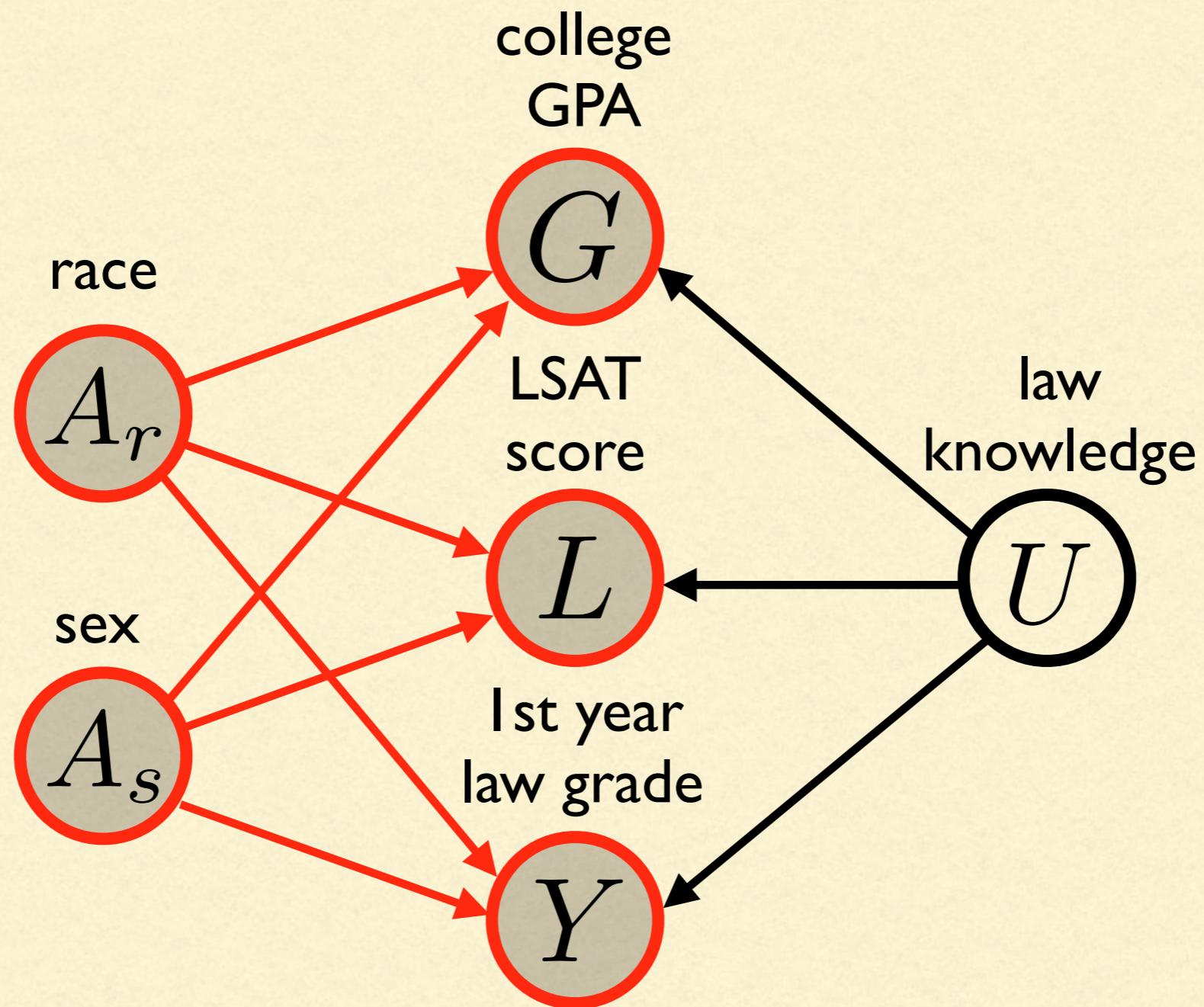


$$\begin{aligned}\hat{Y}(G, L) = \\ \hat{Y}(G_{A_r \leftarrow a'}, L_{A_r \leftarrow a'})\end{aligned}$$

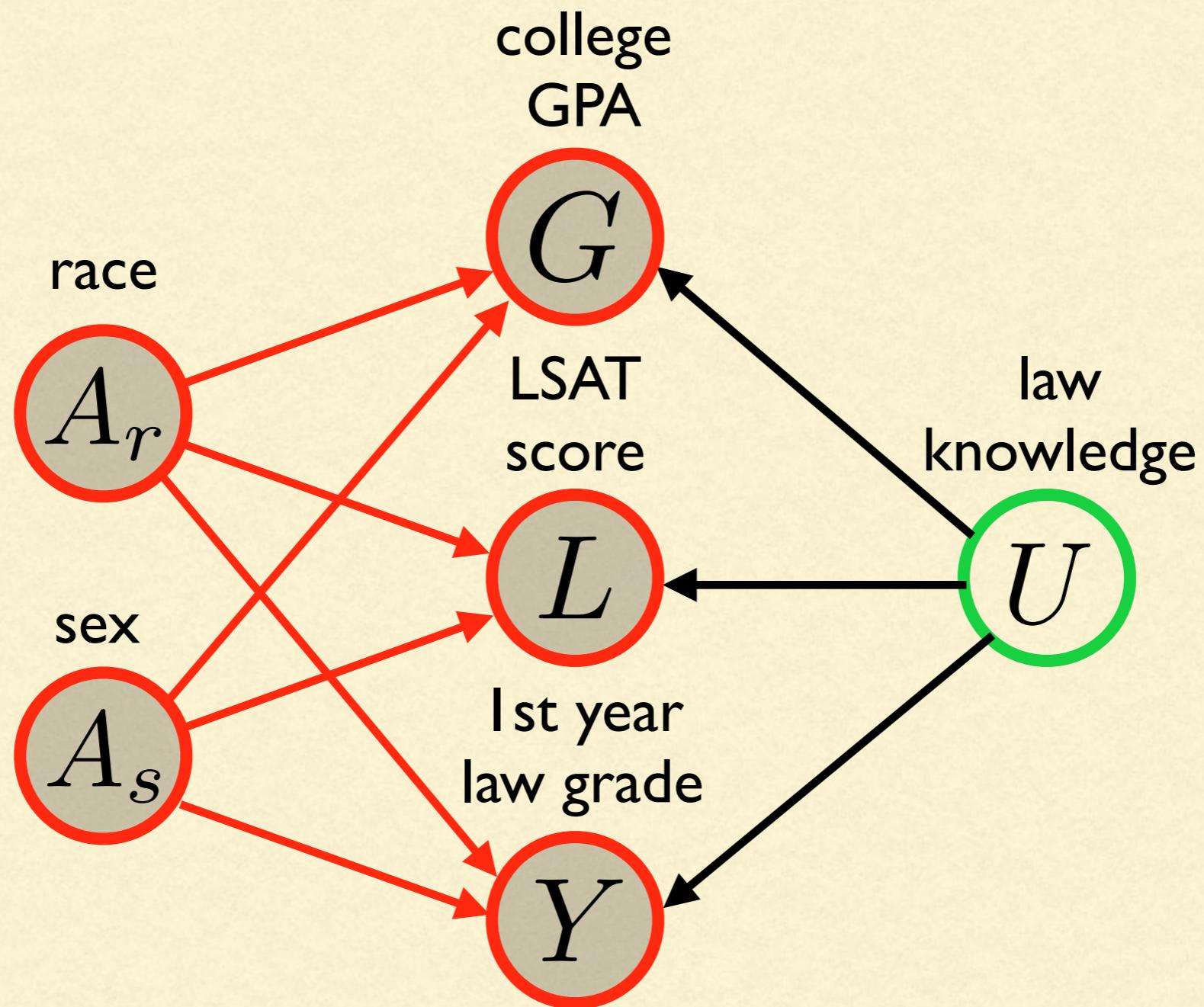
MAKING COUNTERFACTUALLY FAIR PREDICTIONS



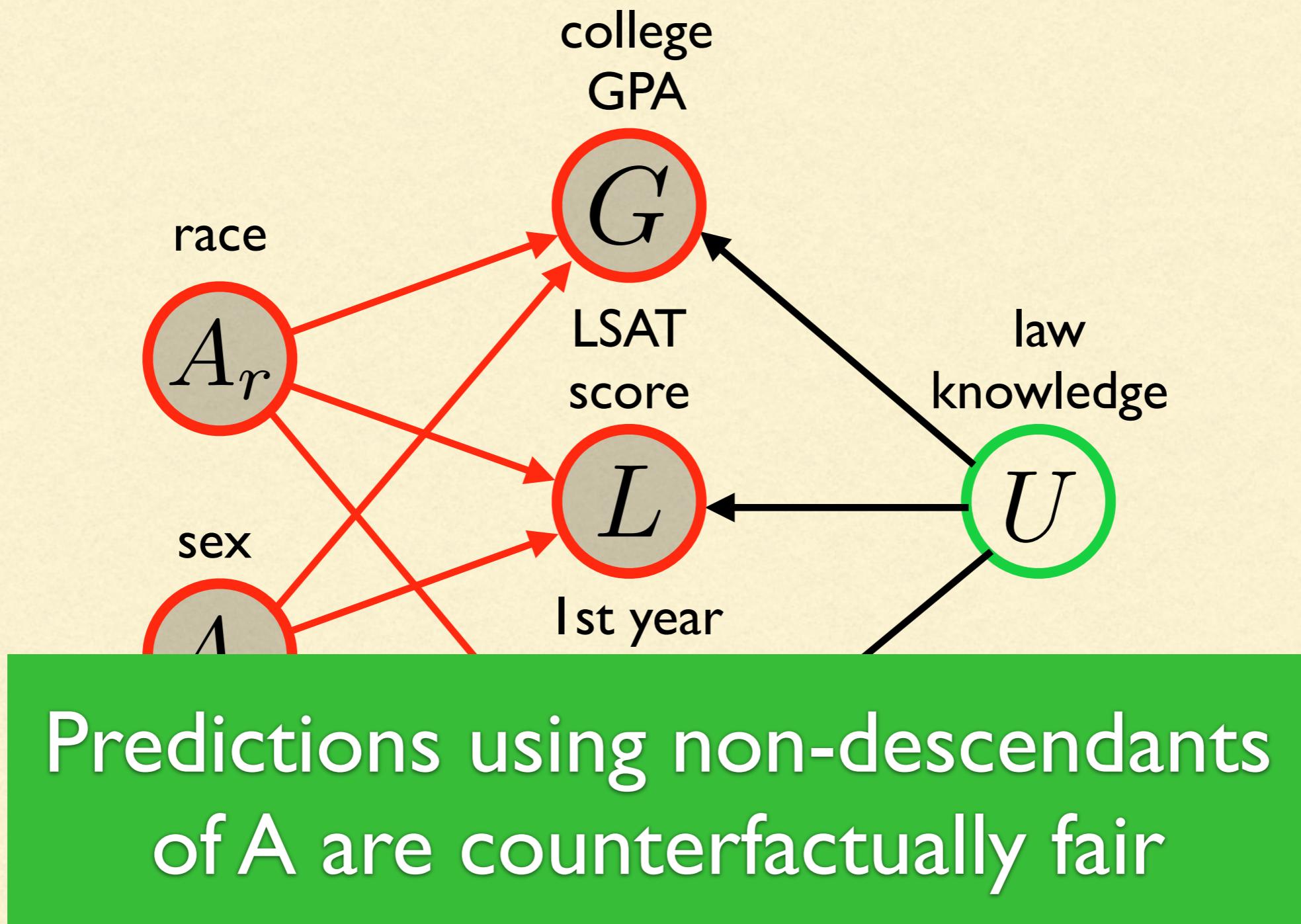
MAKING COUNTERFACTUALLY FAIR PREDICTIONS



MAKING COUNTERFACTUALLY FAIR PREDICTIONS



MAKING COUNTERFACTUALLY FAIR PREDICTIONS



How does this work in practice?

US LAW SCHOOL GRADES

[Wightman, 1998]

21,790 students
163 law schools



US LAW SCHOOL GRADES



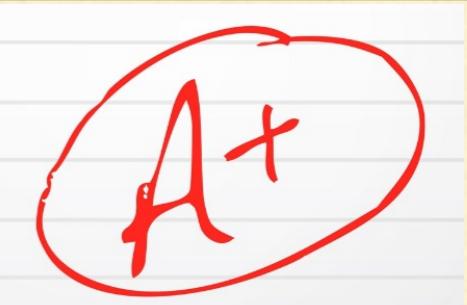
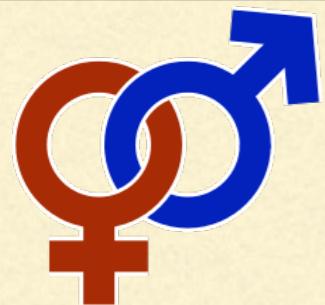
GPA

A large, stylized red 3D text "GPA" with a blue graduation cap resting on top of the letter "G".

Full

RMSE 0.873

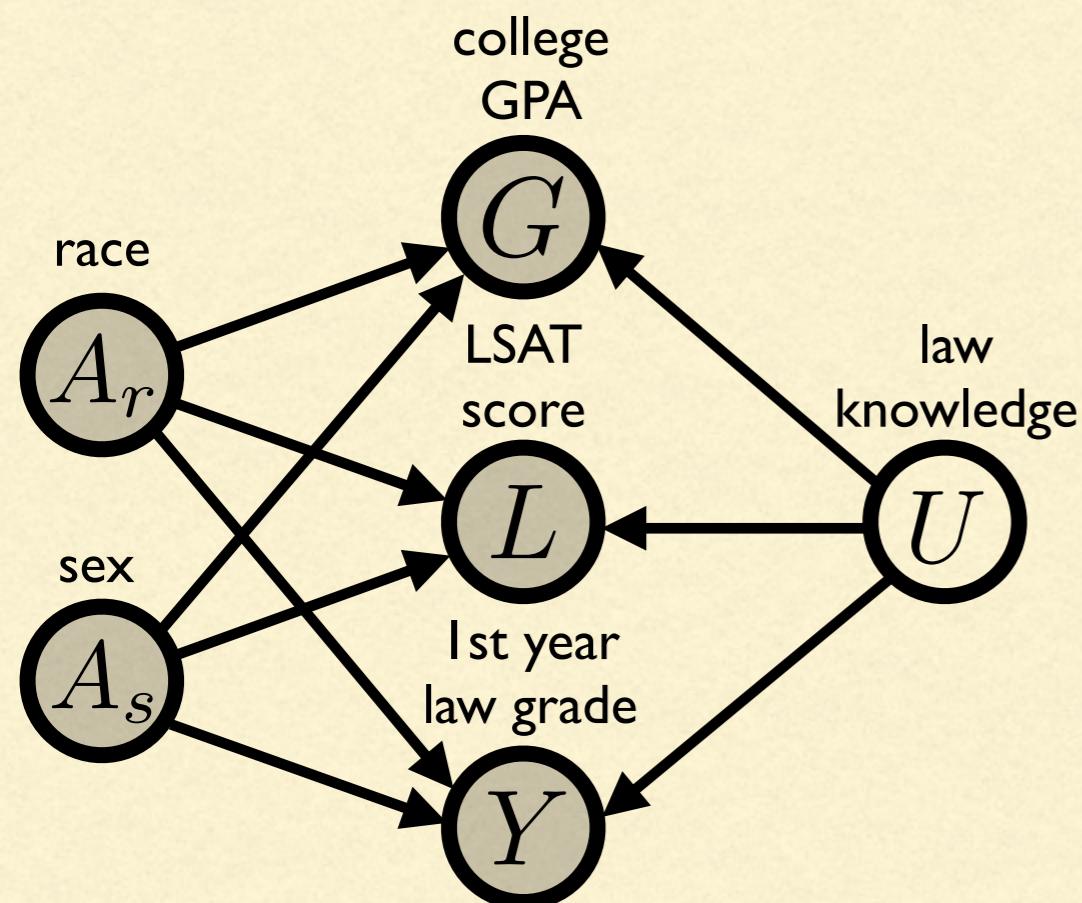
US LAW SCHOOL GRADES



Full Unaware

RMSE	0.873	0.894
------	-------	-------

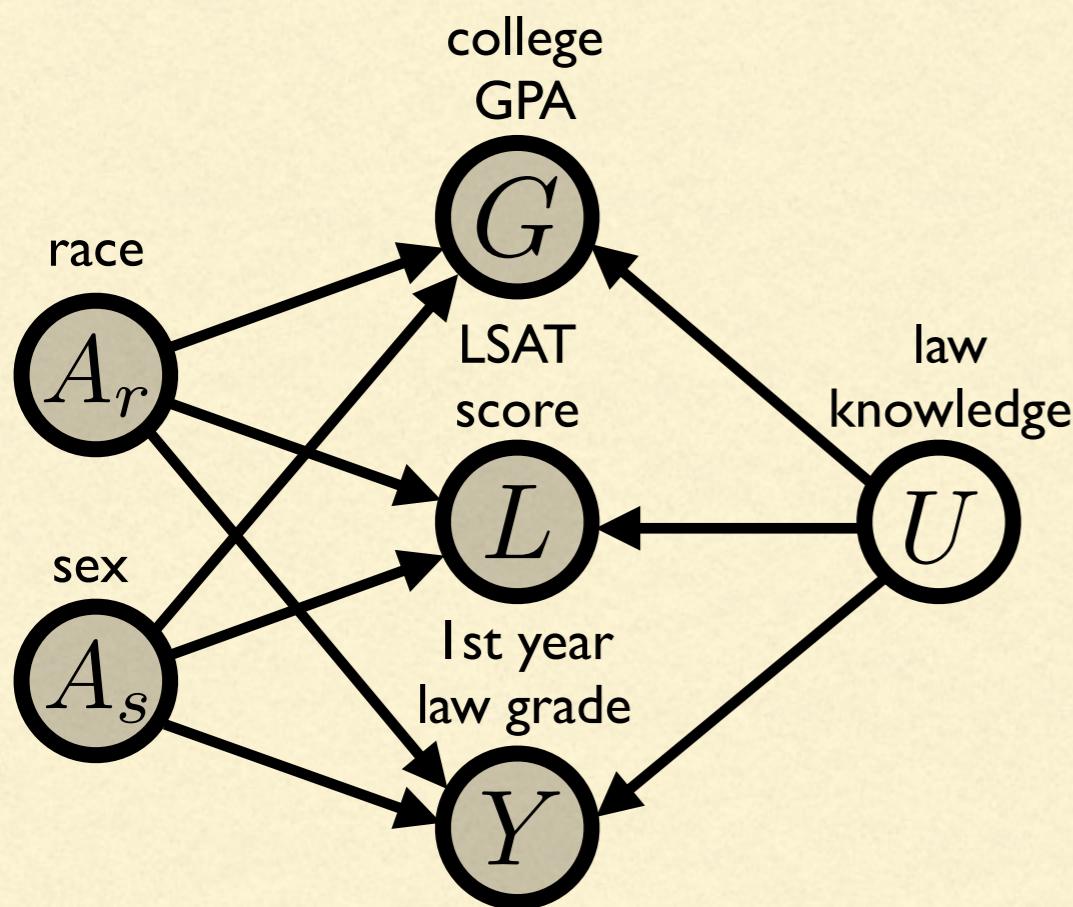
US LAW SCHOOL GRADES



structural causal model (non-deterministic)

$$G \sim \text{Normal}(b_G + [U, A_r, A_s] \mathbf{w}_G, \sigma_G)$$
$$L \sim \text{Poisson}(\exp(b_L + [U, A_r, A_s] \mathbf{w}_L))$$
$$Y \sim \text{Normal}([U, A_r, A_s] \mathbf{w}_Y, 1)$$
$$U \sim \text{Normal}(0, 1)$$

US LAW SCHOOL GRADES



structural causal model (non-deterministic)

$$G \sim \text{Normal}(b_G + [U, A_r, A_s] \mathbf{w}_G, \sigma_G)$$
$$L \sim \text{Poisson}(\exp(b_L + [U, A_r, A_s] \mathbf{w}_L))$$
$$Y \sim \text{Normal}([U, A_r, A_s] \mathbf{w}_Y, 1)$$
$$U \sim \text{Normal}(0, 1)$$

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

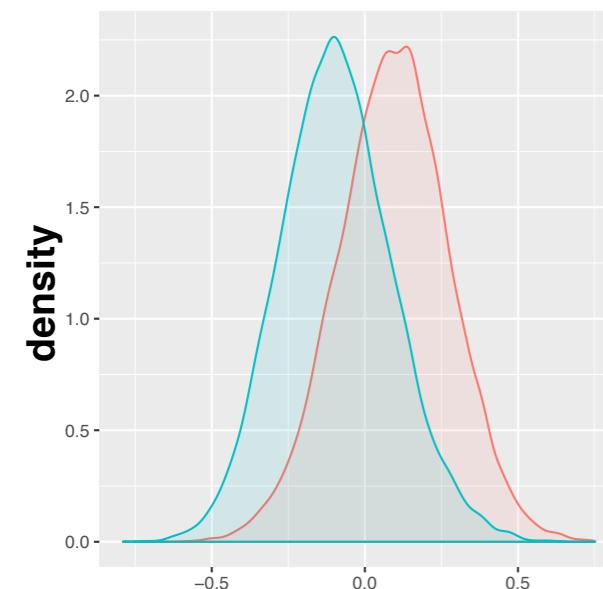
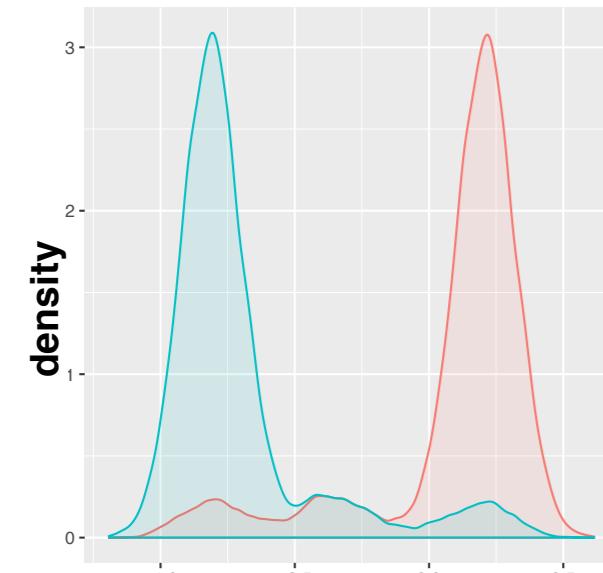
LAW SCHOOL COUNTERFACTUALS

Unaware Full

 \hat{Y} \hat{Y} \hat{Y} \hat{Y} 

LAW SCHOOL COUNTERFACTUALS

black \leftrightarrow white



- original data
- counter-factual

\hat{Y}

\hat{Y}

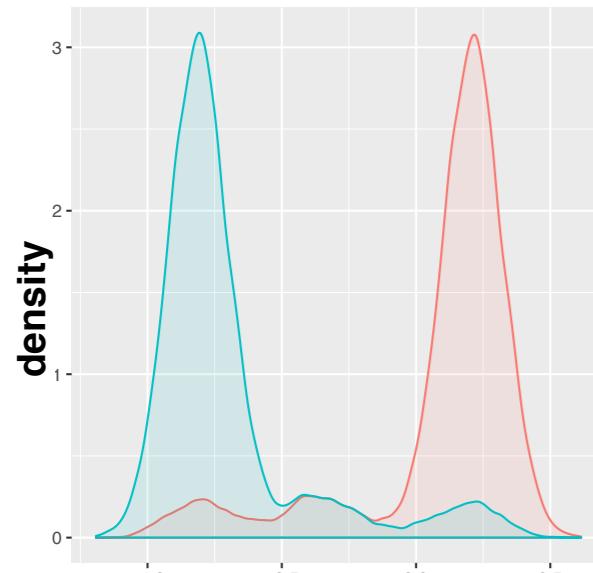
\hat{Y}

\hat{Y}

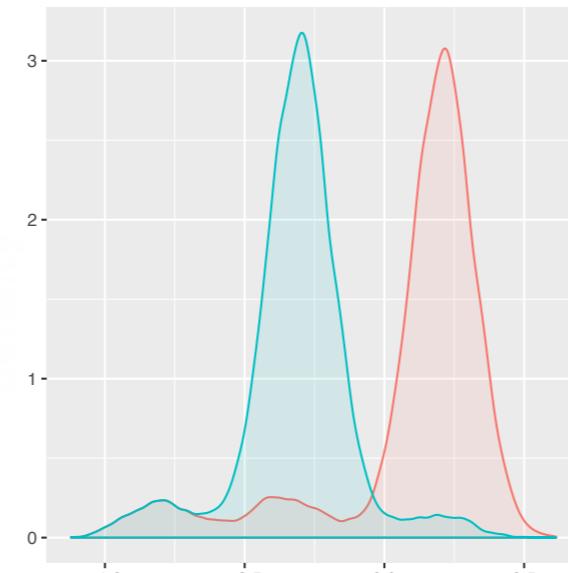
LAW SCHOOL COUNTERFACTUALS

Full

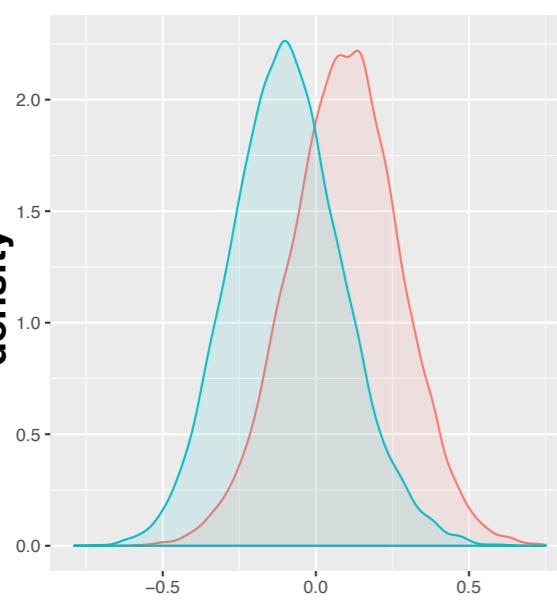
black \leftrightarrow white



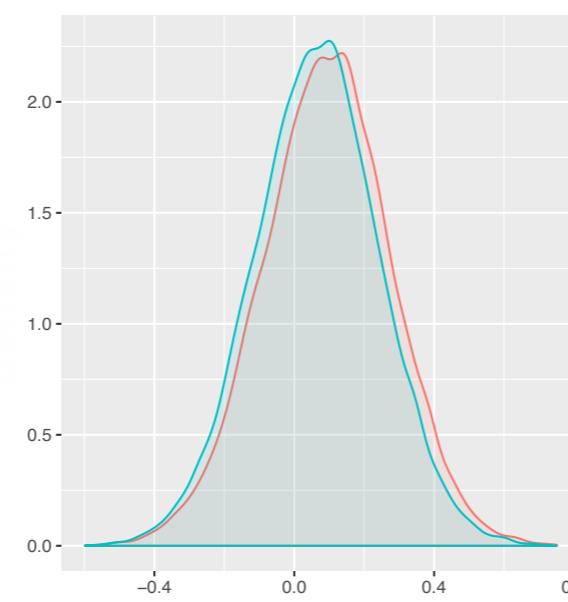
asian \leftrightarrow white



Unaware



\hat{Y}



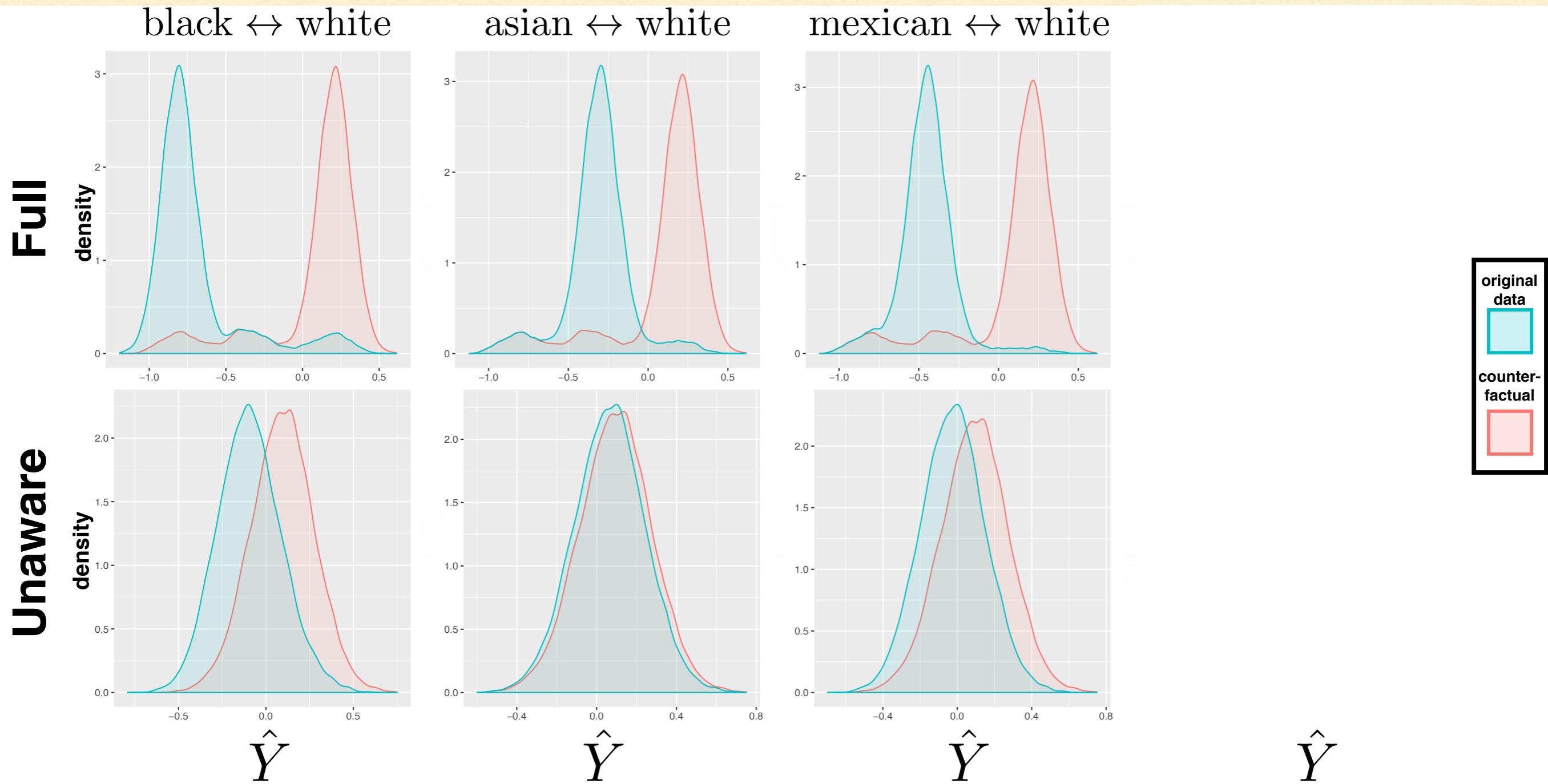
\hat{Y}

\hat{Y}

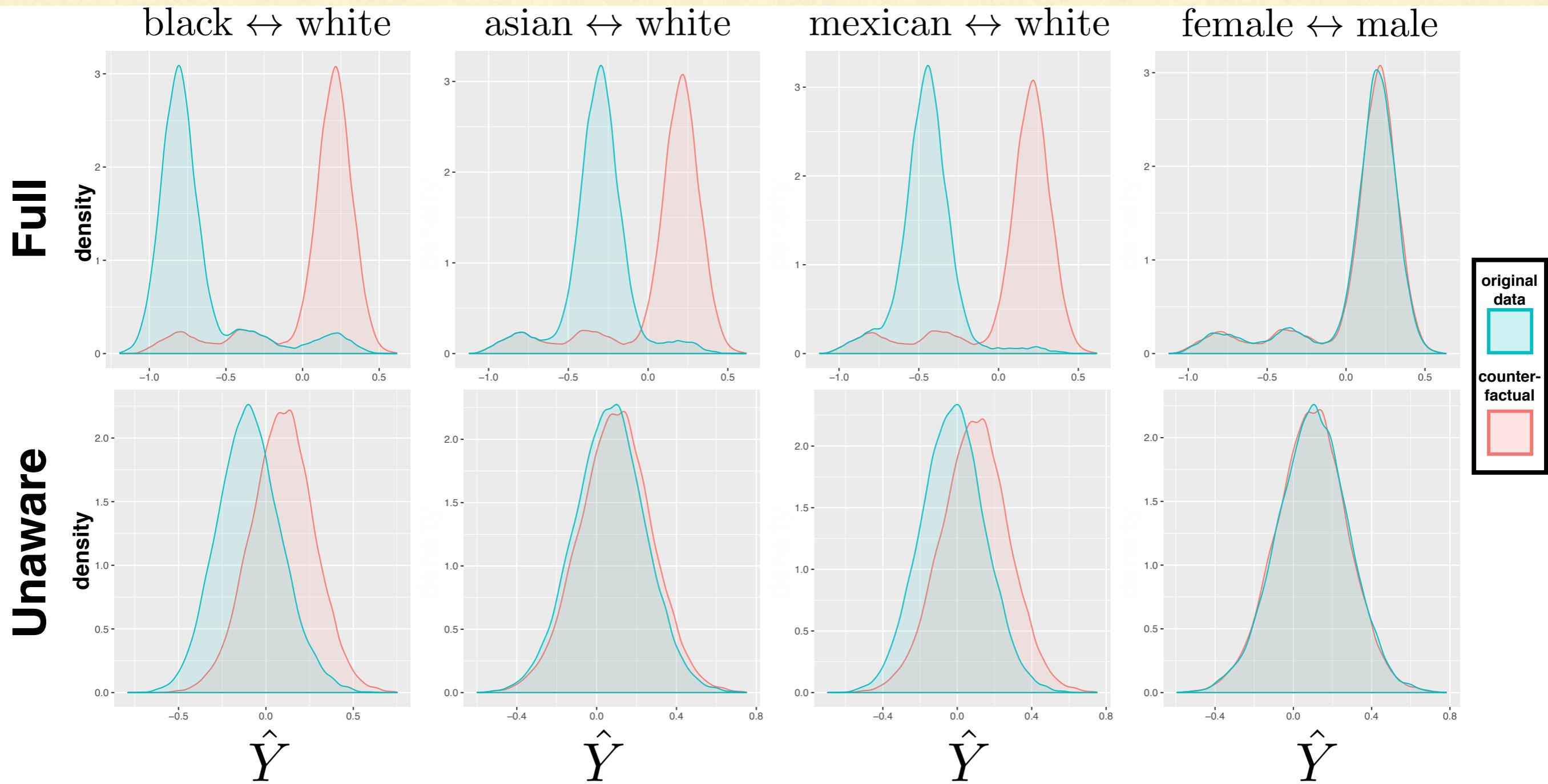
\hat{Y}

original data
counter-factual

LAW SCHOOL COUNTERFACTUALS



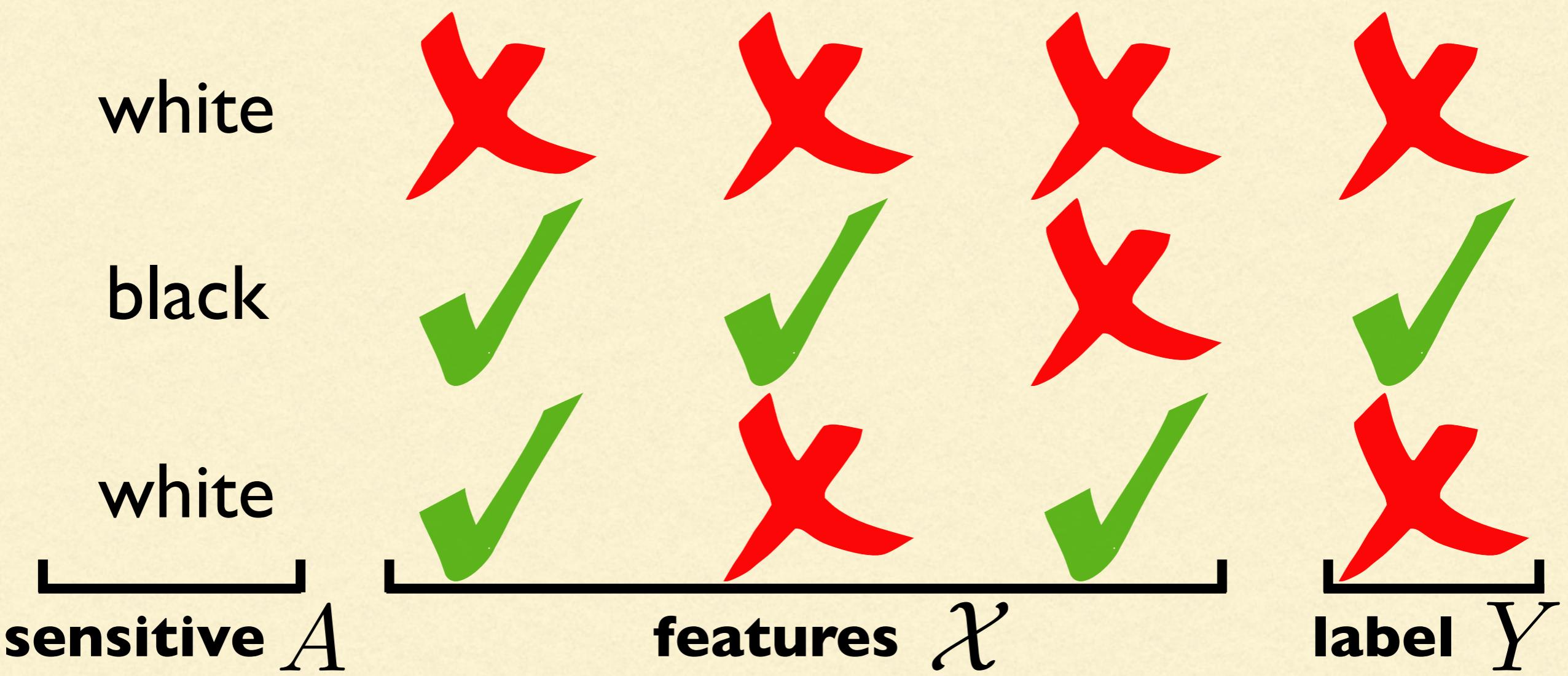
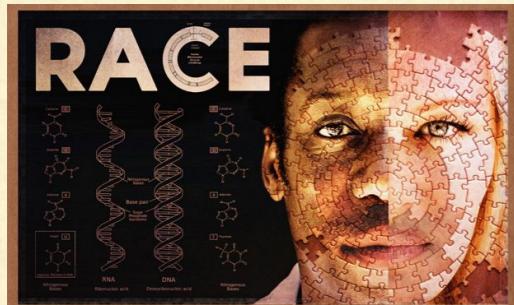
LAW SCHOOL COUNTERFACTUALS



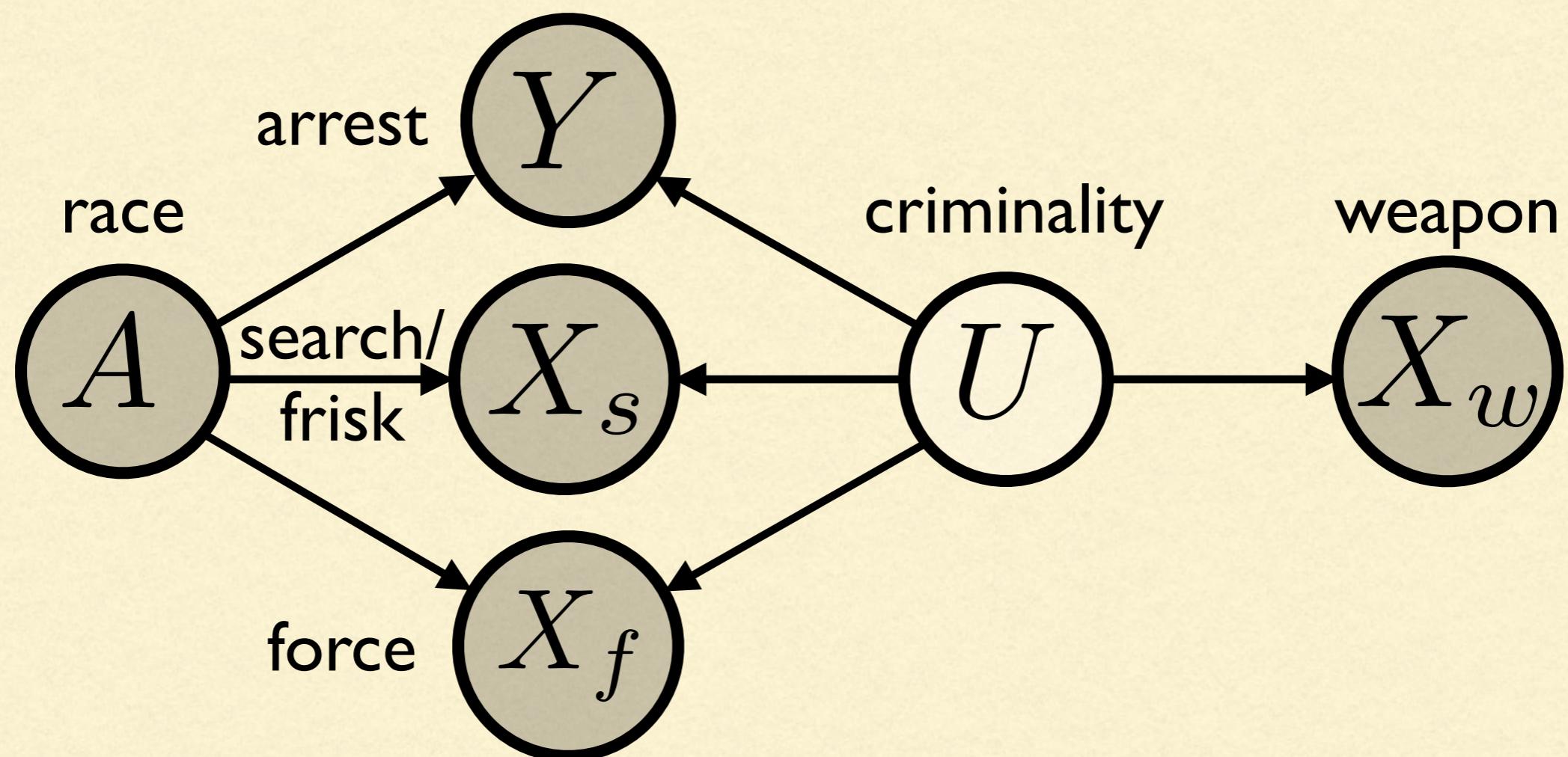
NYC STOP-AND-FRISK DATA



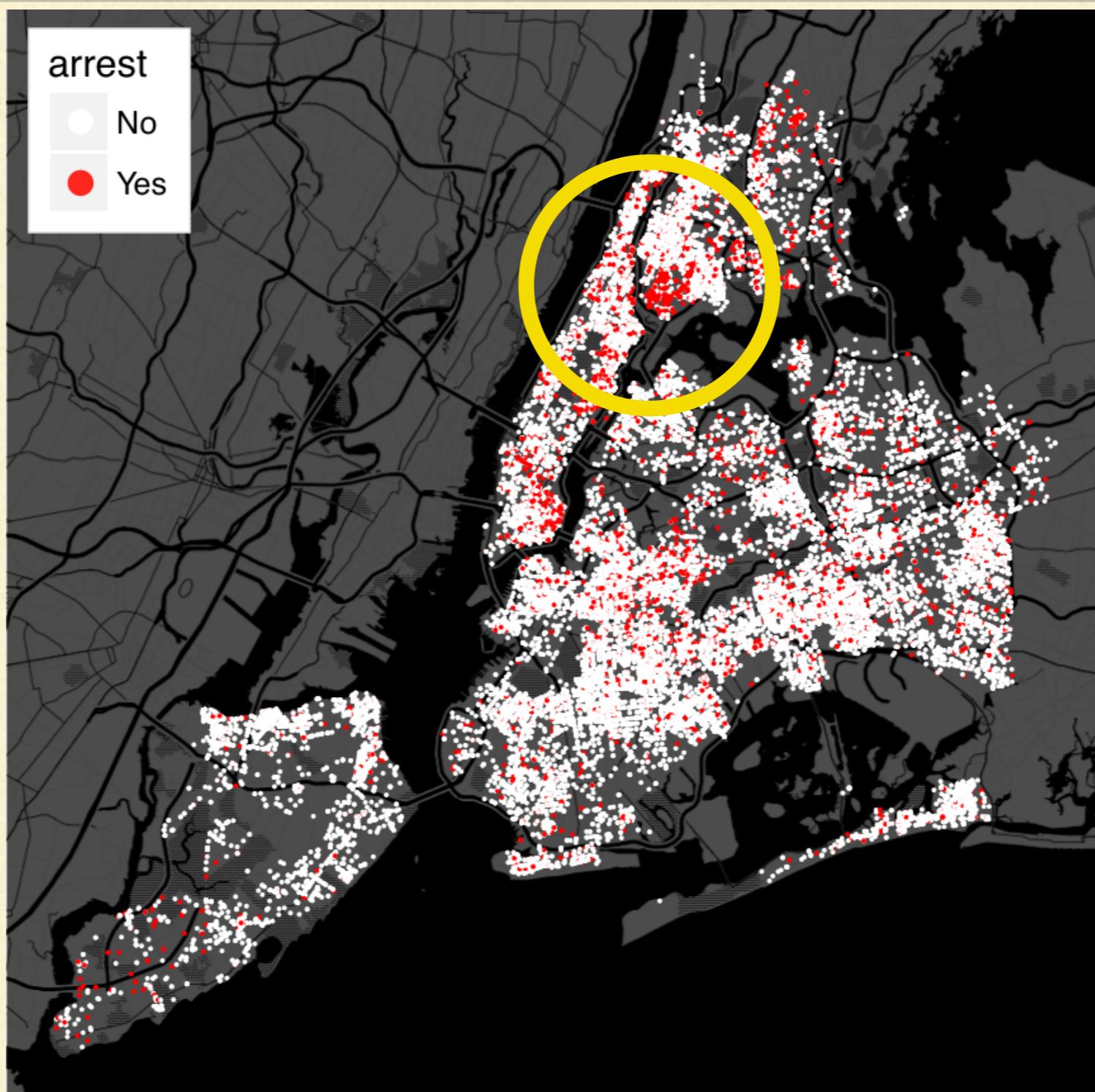
NYC STOP-AND-FRISK DATA



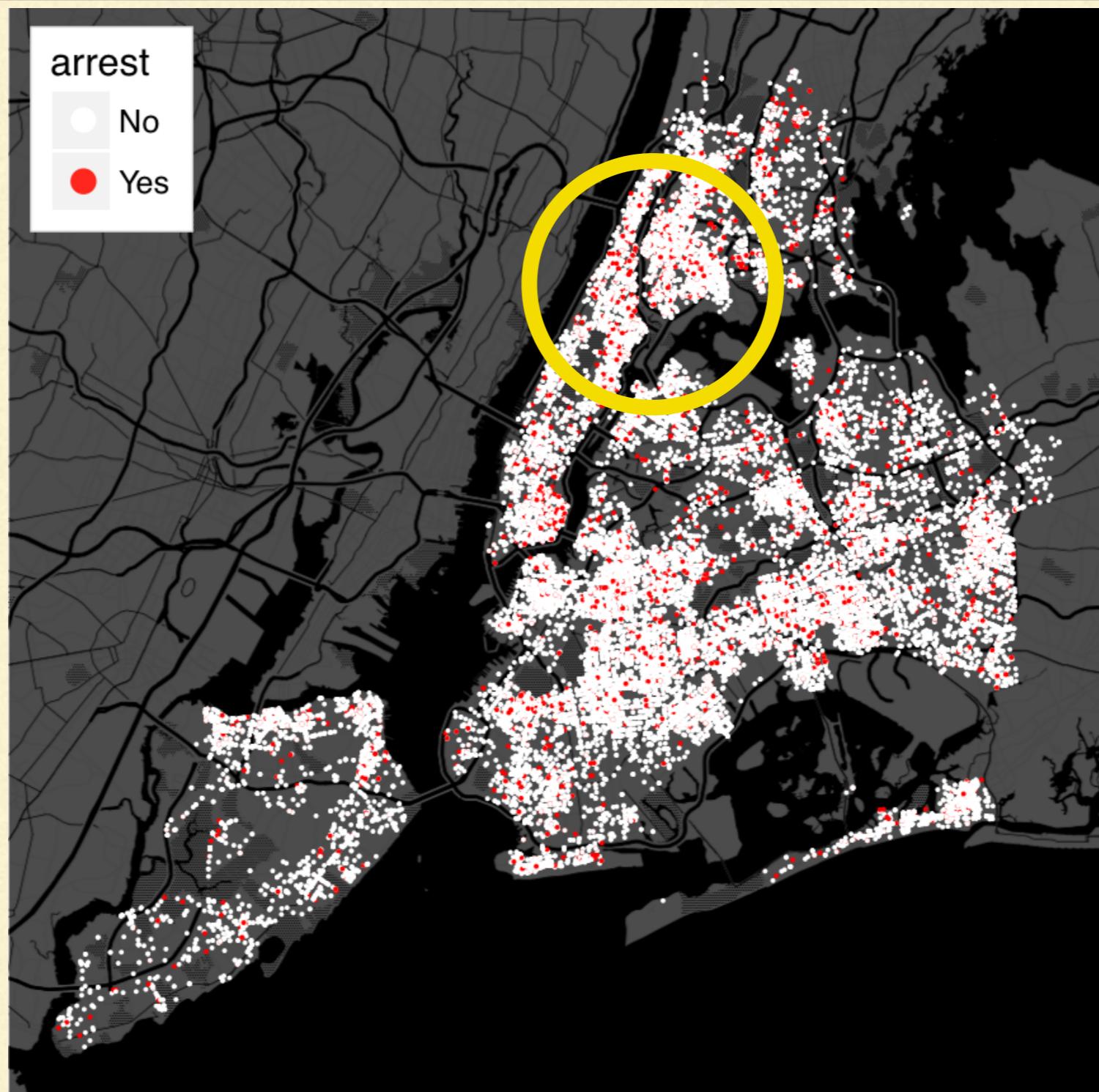
NYC STOP-AND-FRISK DATA



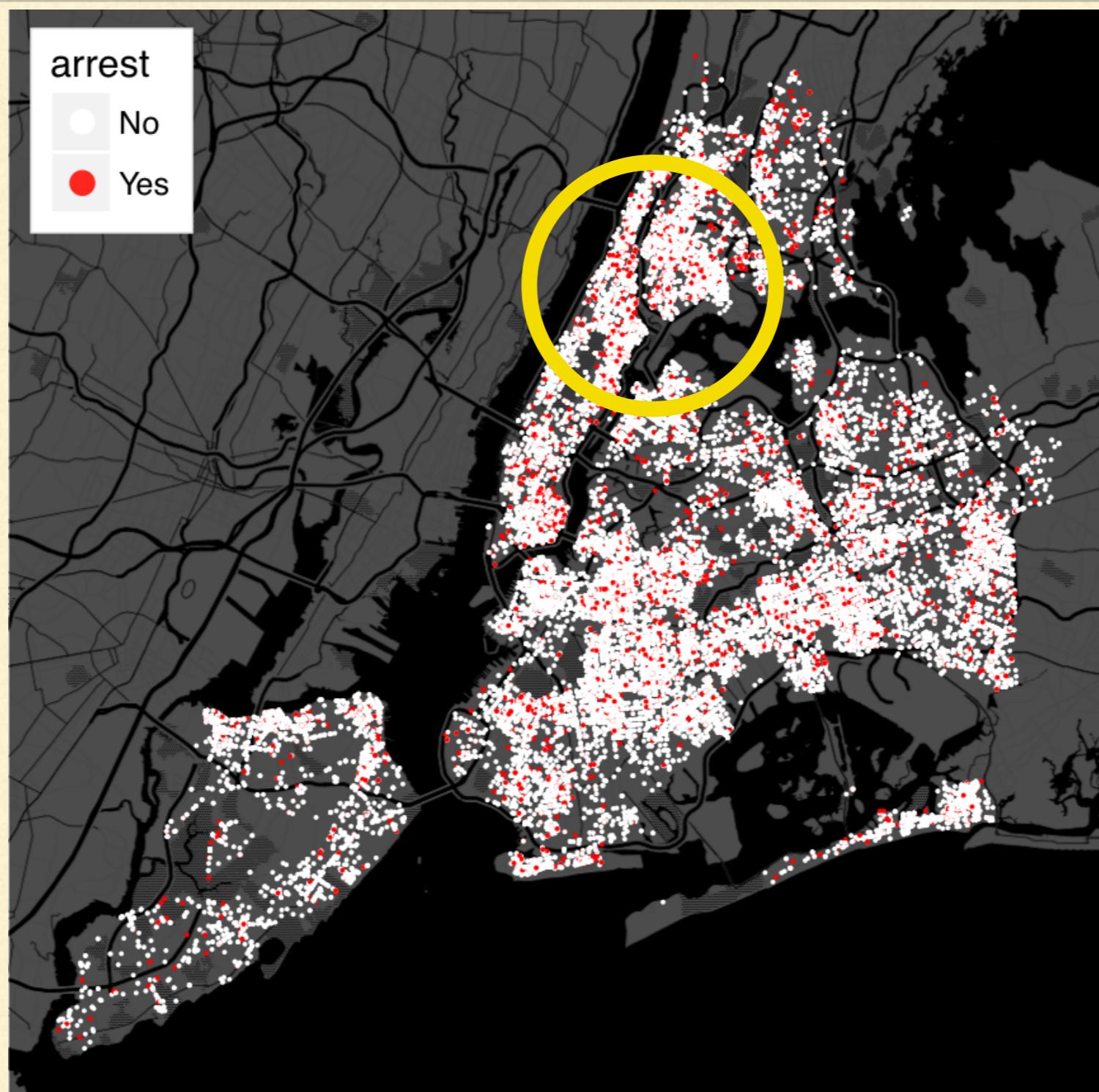
NYC STOP-AND-FRISK DATA



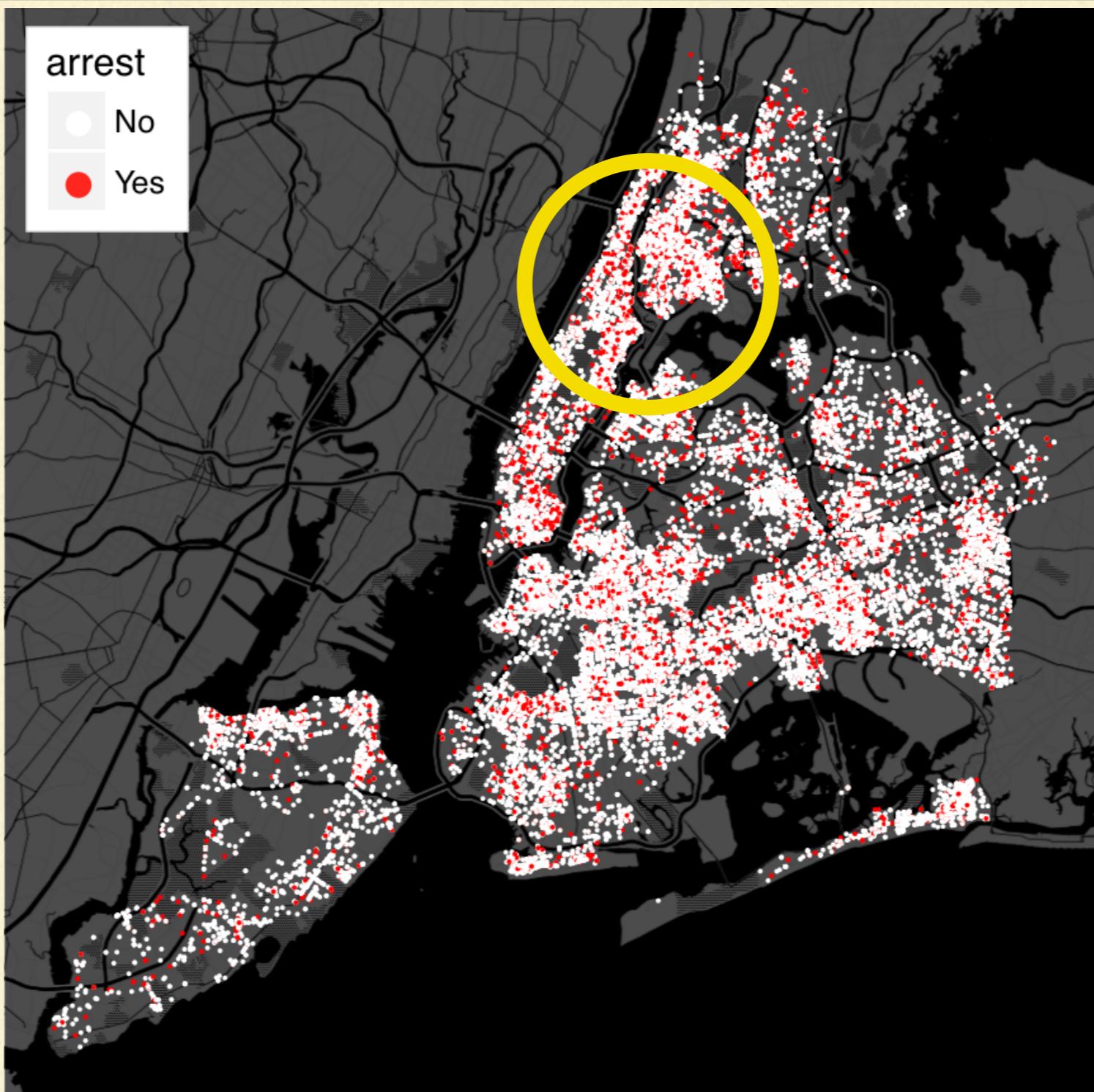
WHITE COUNTERFACTUAL



BLACK COUNTERFACTUAL



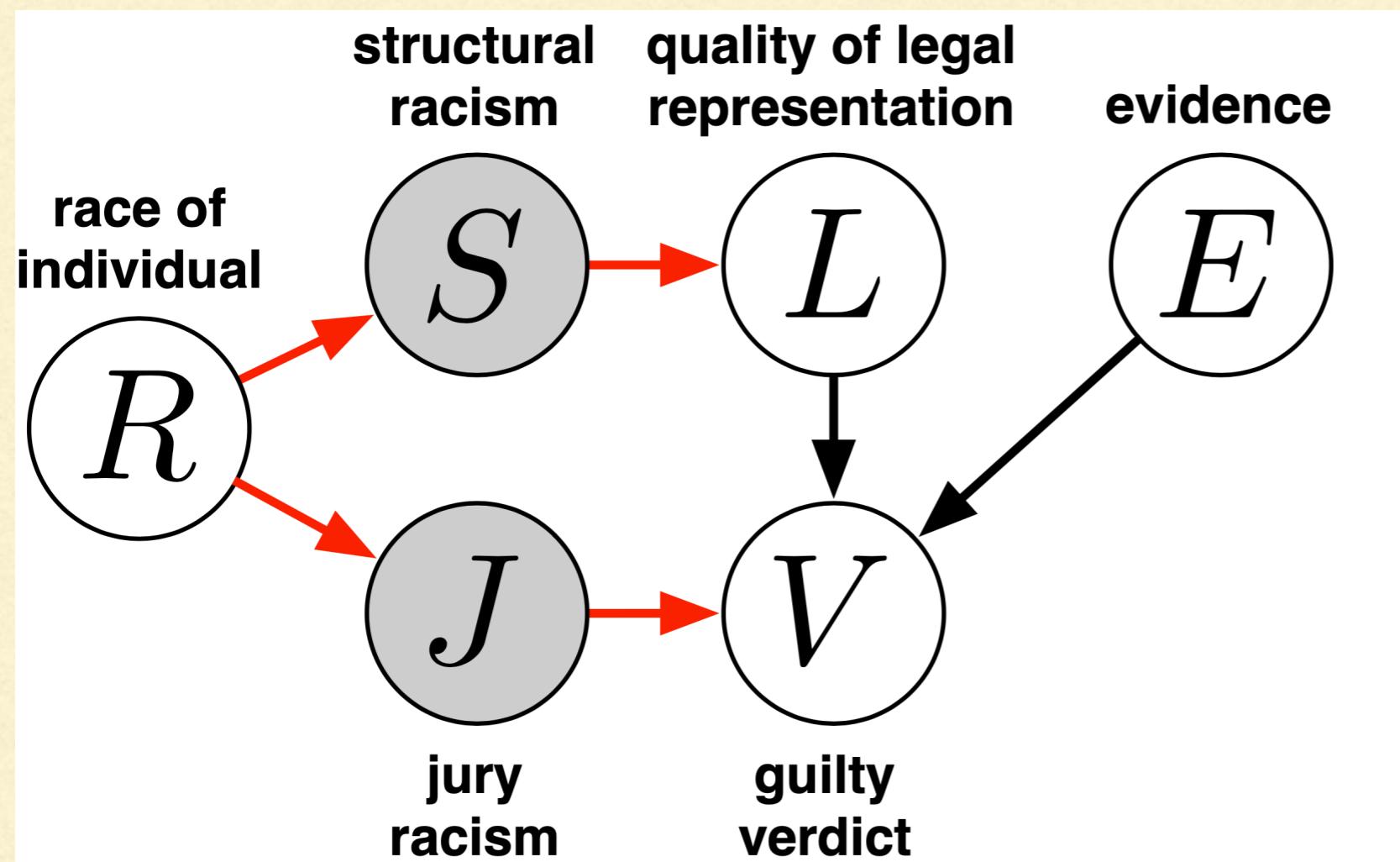
BLACK-HISPANIC COUNTERFACTUAL



OTHER CAUSAL TOOLS

Sensitivity Analysis

[Kilbertus et al., 2019]

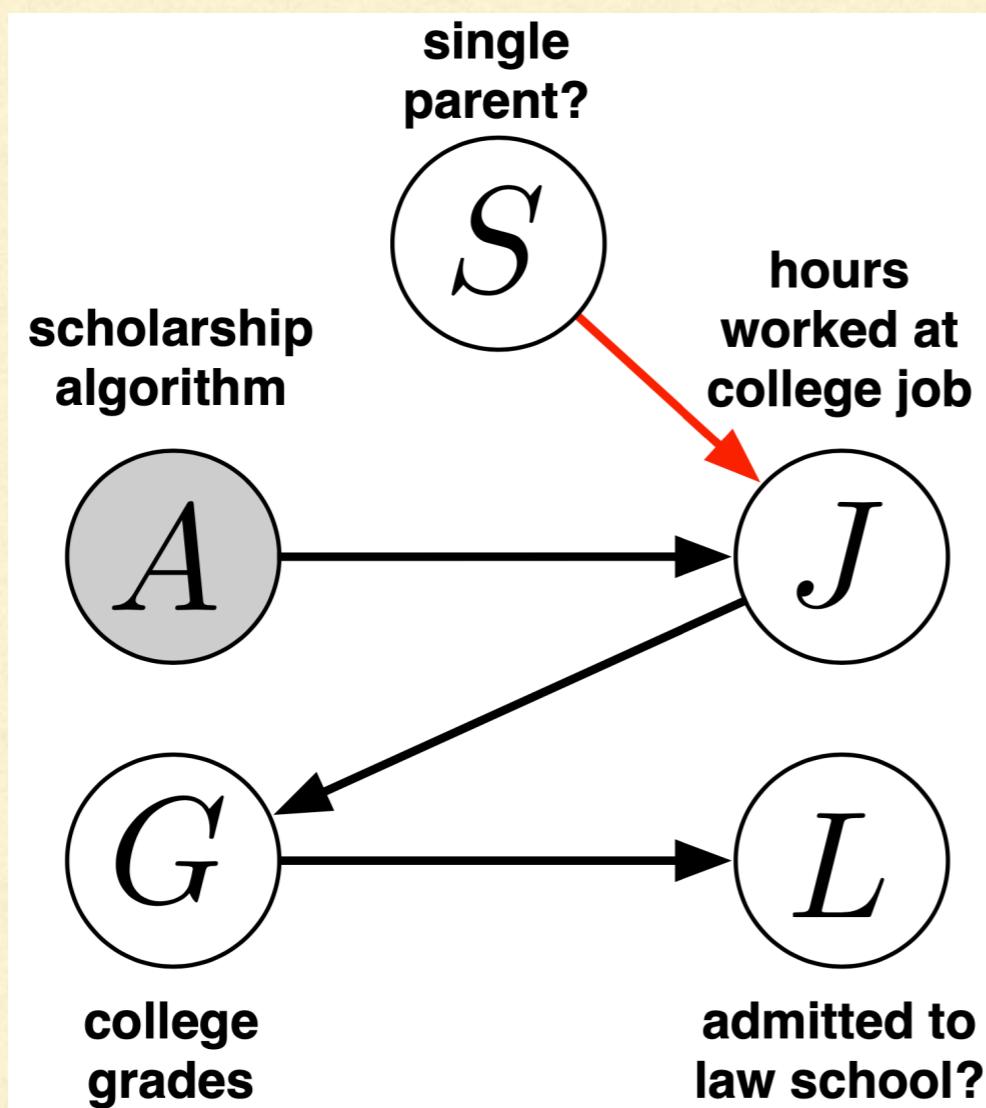


“How wrong could
a fairness method
be if **racism** is
incorrectly
estimated?”

OTHER CAUSAL TOOLS

Long-term Impacts

[Liu et al., 2018; Kannan et al., 2018; Kusner et al., 2019]



“What are the
**long-term
impacts of
algorithmic
decisions?**”

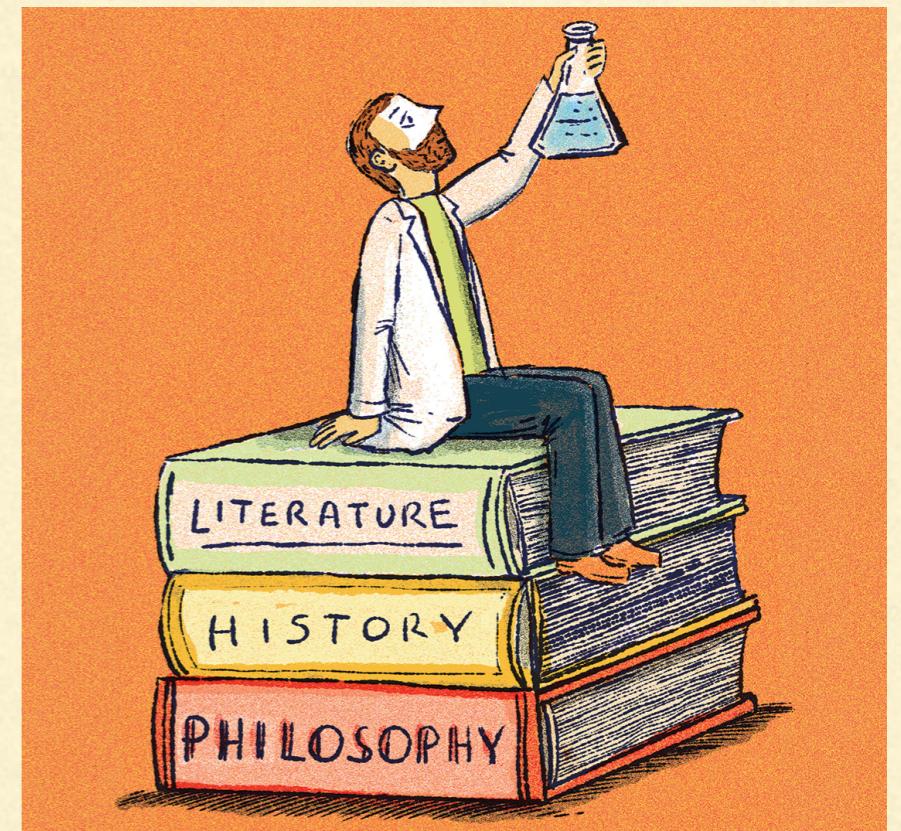
NEXT STEP #1: TAKE DIRECTIONS FROM THE HUMANITIES

Machine Learning



+

Humanities



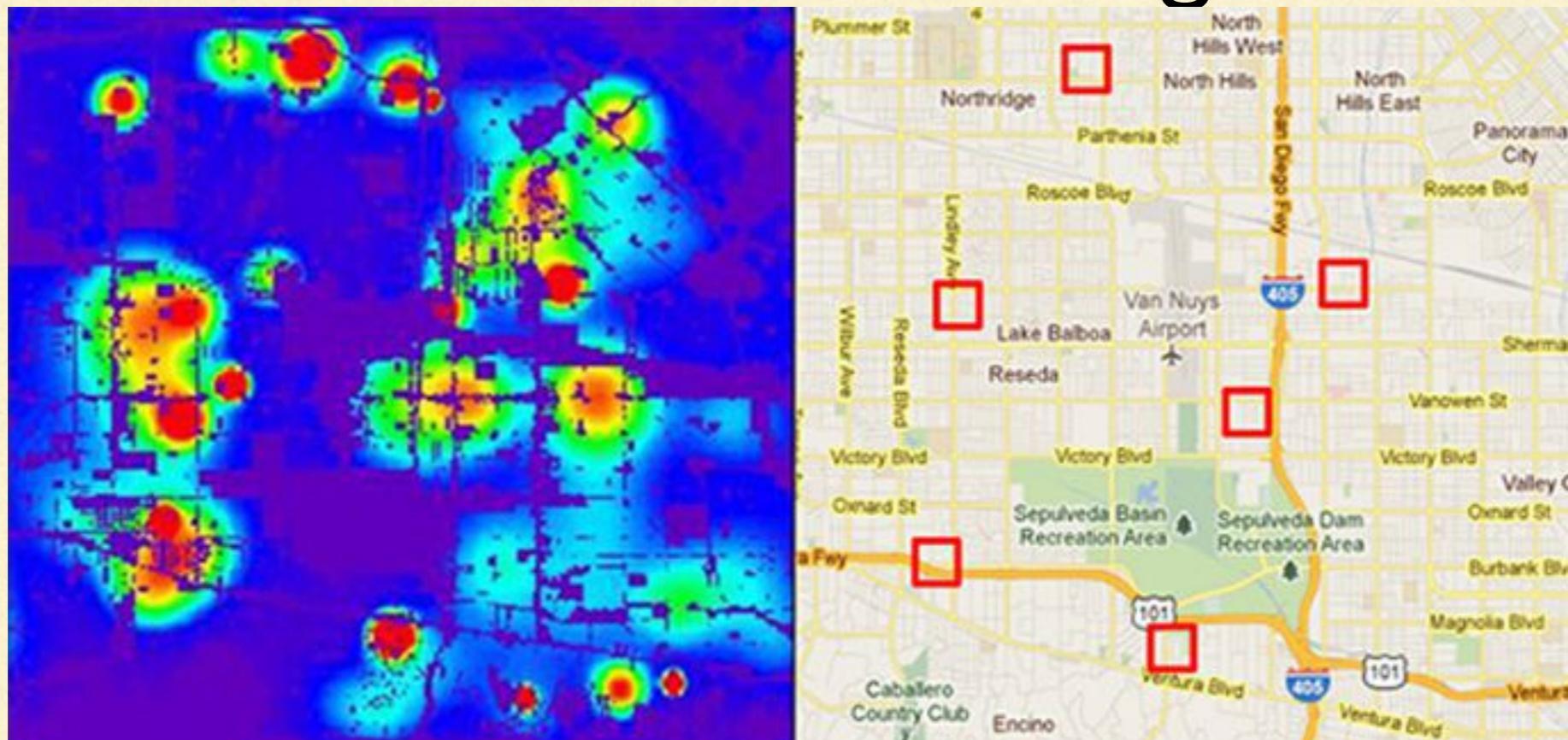
NEXT STEP #2: WORK ALONGSIDE HARMED INDIVIDUALS

An International Algorithm Regulatory Institute



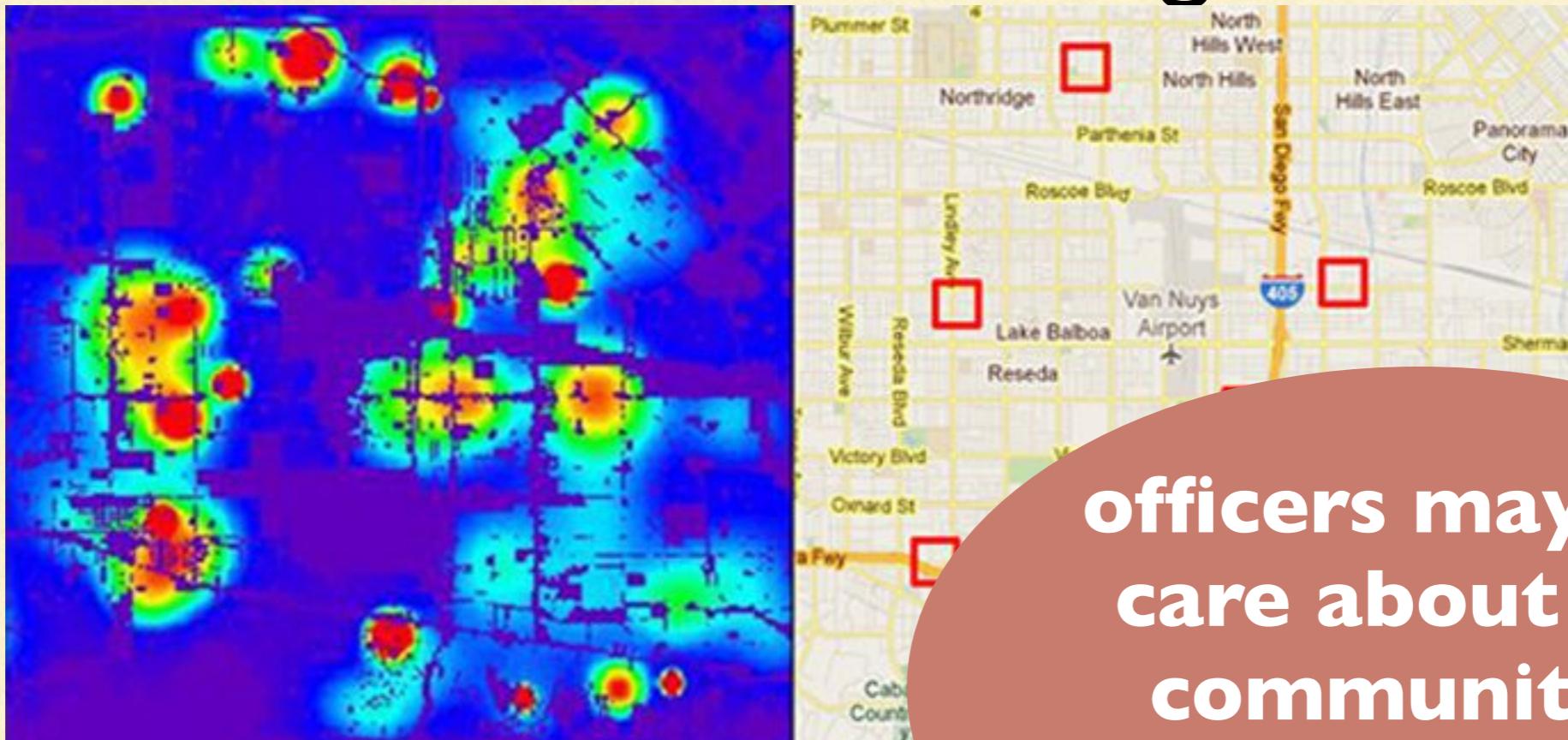
NEXT STEP #3: IDENTIFY WHEN ALGORITHMS ARE INAPPROPRIATE

Predictive Policing



NEXT STEP #3: IDENTIFY WHEN ALGORITHMS ARE INAPPROPRIATE

Predictive Policing

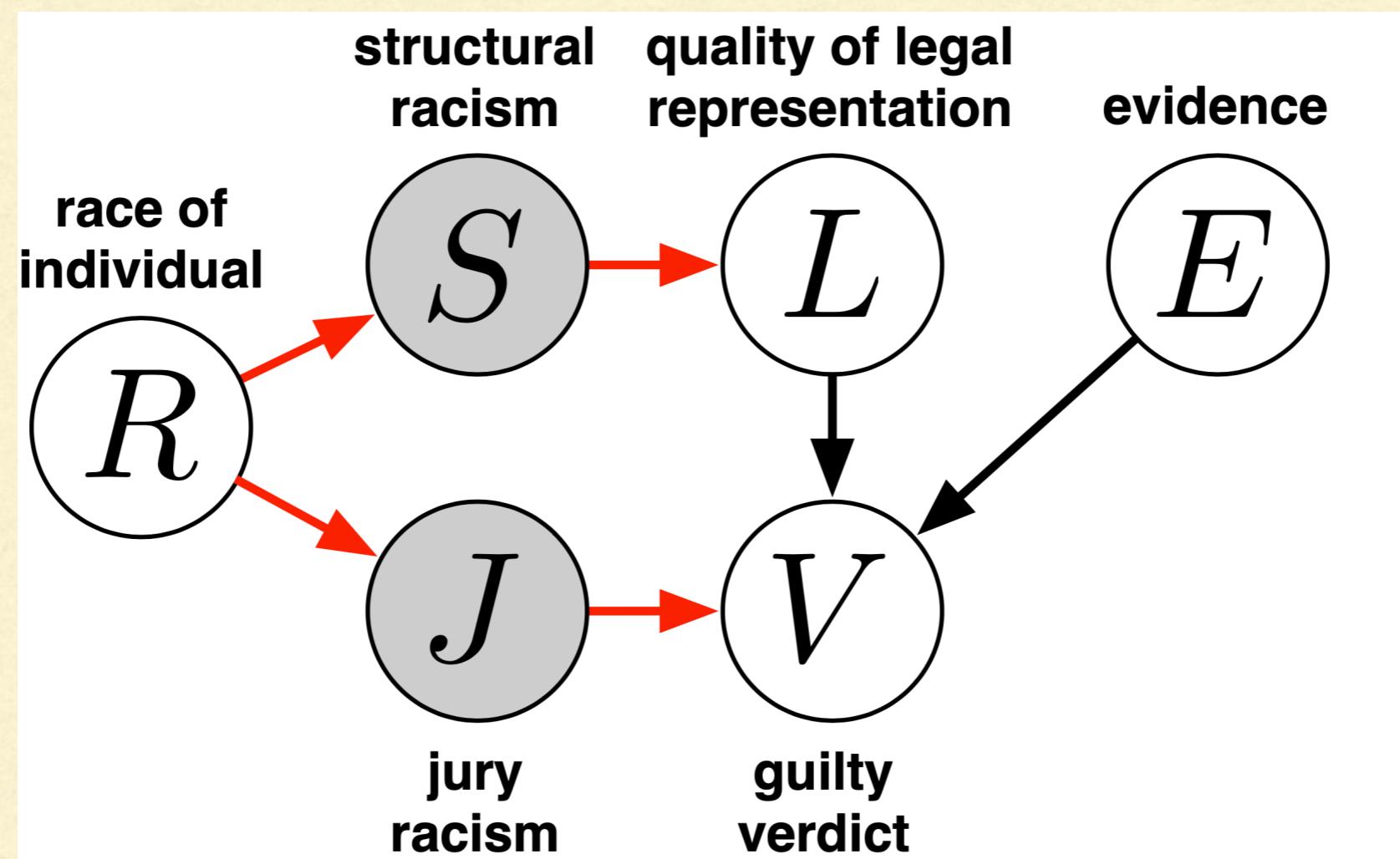


officers may not care about the communities they police

[McManus et al., 2019]

NEXT STEP #4: STAY CRITICAL

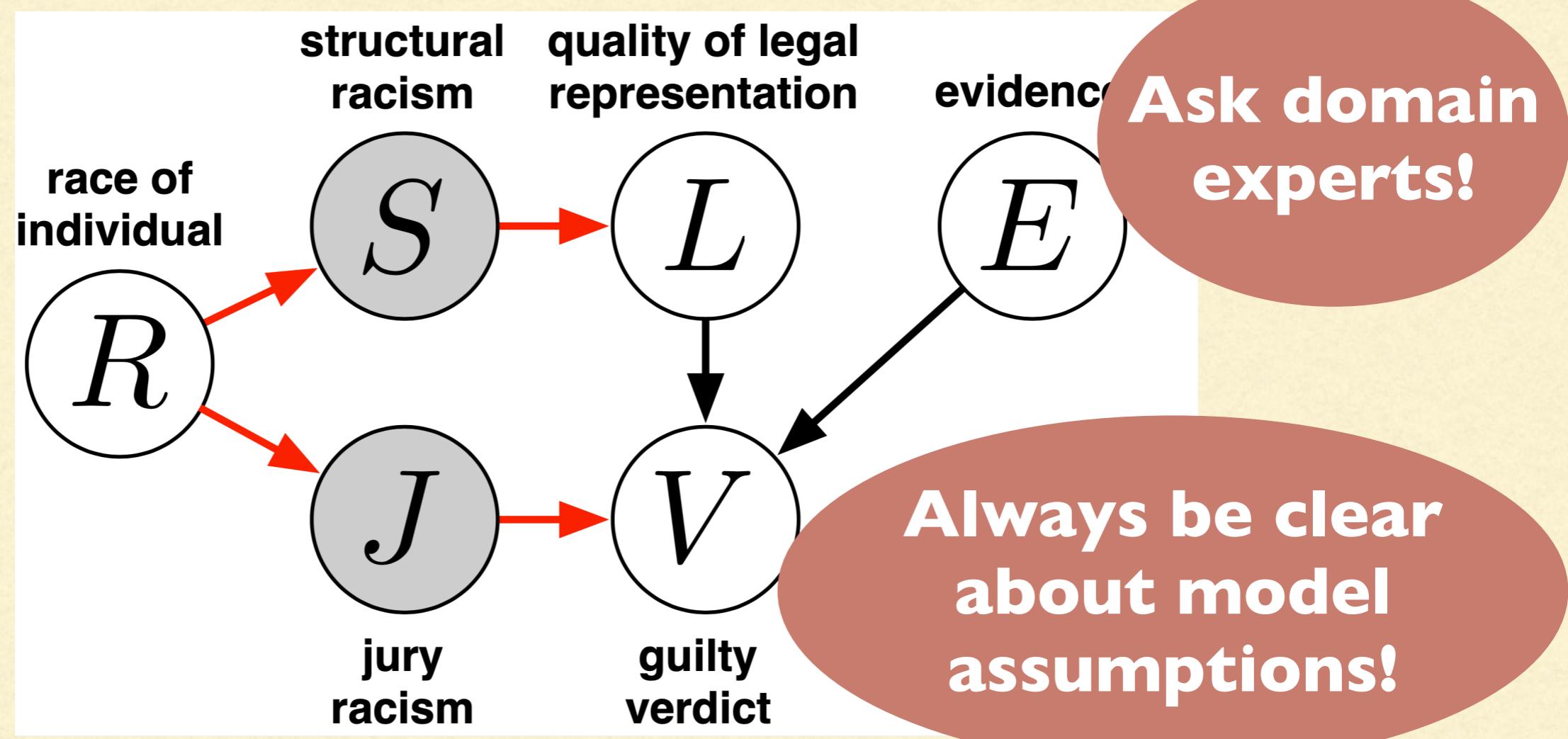
Causal Models Can Also Be Wrong



Counterfactuals Are Impossible To Verify

NEXT STEP #4: STAY CRITICAL

Causal Models Can Also Be Wrong



Counterfactuals Are Impossible To Verify

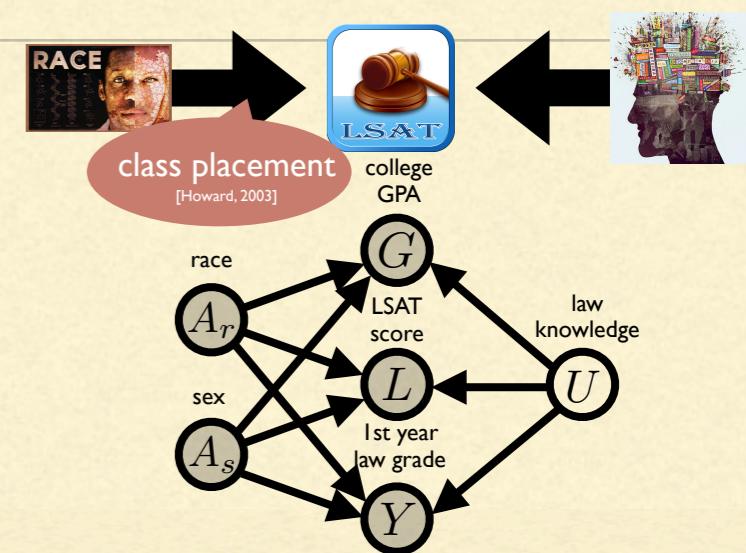
TAKE-AWAYS

- Race/Gender/Sexual Orientation/etc. could cause AI decisions to change unfairly



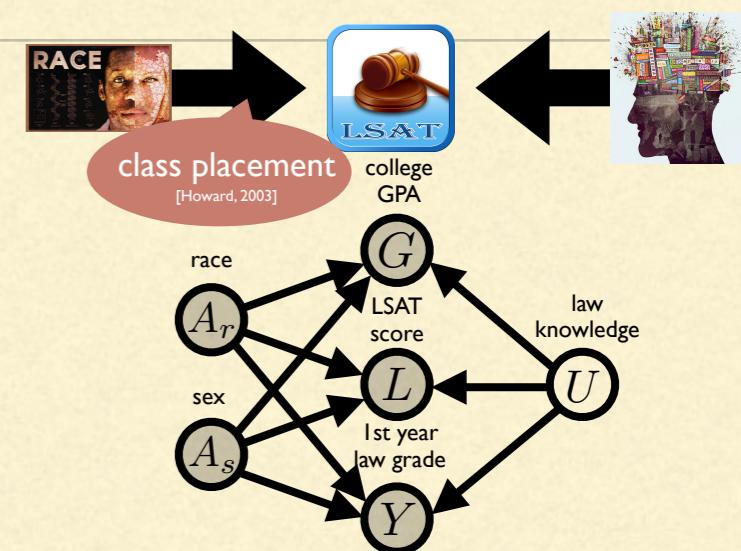
TAKE-AWAYS

- Race/Gender/Sexual Orientation/etc. could cause AI decisions to change unfairly
- Idea: Model how these attributes *cause* unfair decisions via *causal models*



TAKE-AWAYS

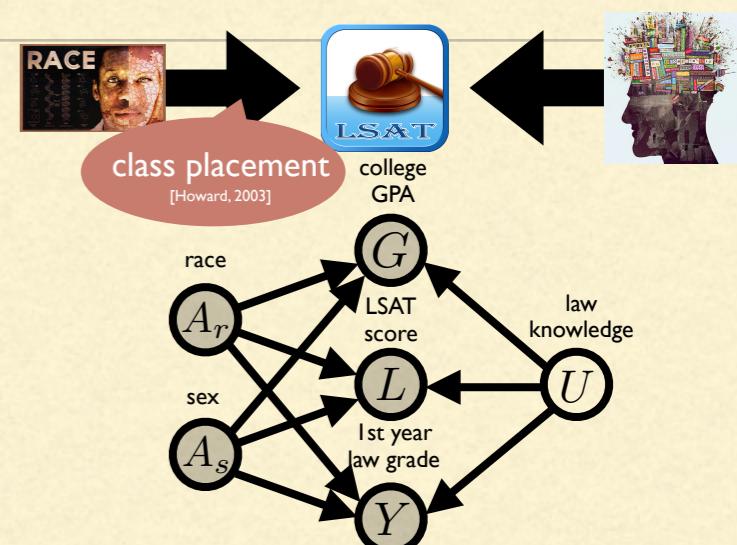
- Race/Gender/Sexual Orientation/etc. could cause AI decisions to change unfairly
- Idea: Model how these attributes *cause* unfair decisions via *causal models*
- **Counterfactual Fairness:** A predictor is fair if it would give the same prediction in a world where you were different



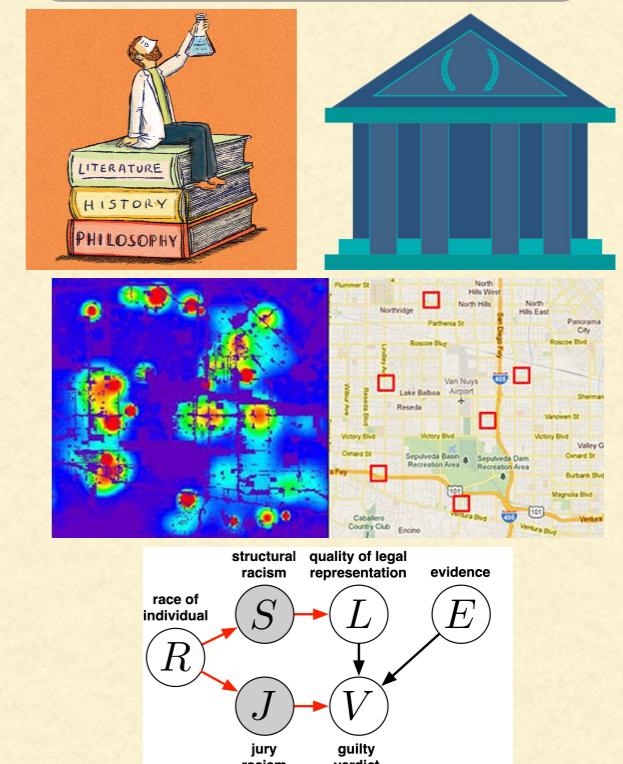
$$\begin{aligned}\hat{Y}(G, L) = \\ \hat{Y}(G_{A_r \leftarrow a'}, L_{A_r \leftarrow a'})\end{aligned}$$

TAKE-AWAYS

- Race/Gender/Sexual Orientation/etc. could cause AI decisions to change unfairly
- Idea: Model how these attributes *cause* unfair decisions via *causal models*
- **Counterfactual Fairness:** A predictor is fair if it would give the same prediction in a world where you were different
- Much still to do in the future:
 - #1 Take direction from the humanities
 - #2 Work alongside harmed individuals
 - #3 Don't always use algorithms
 - #4 Critique all models (even causal ones)



$$\hat{Y}(G, L) = \hat{Y}(G_{A_r \leftarrow a'}, L_{A_r \leftarrow a'})$$



THANKS TO:

Ricardo Silva



Chris Russell



Josh Loftus



Niki Kilbertus



Phillip Ball

