
MAKING DECISIONS THAT REDUCE DISCRIMINATORY IMPACT

Matt J. Kusner

Chris Russell

Joshua R. Loftus

Ricardo Silva

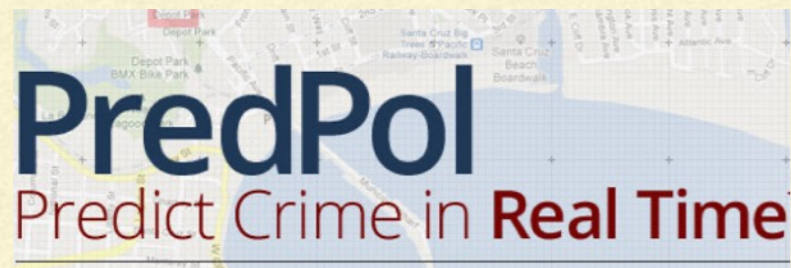


**The
Alan Turing
Institute**

ML IS INVOLVED IN LIFE-CHANGING DECISIONS

Policing

[Ensign et al., 2017]



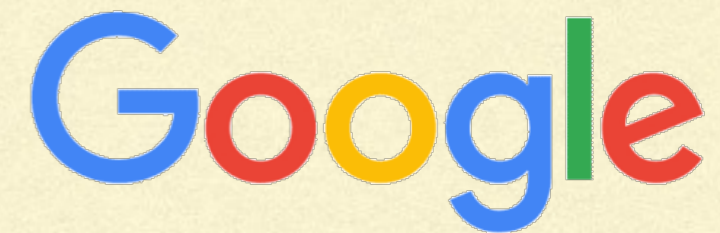
Parole Sentencing

[Larson et al., 2016]



Advertising

[Sweeney, 2013]



Insurance



Lending



**Insert your
favorite
application
here!**

WE HAVE PROBLEMS

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

NEWS

LIVE BBC NEWS AT TEN

Page last updated at 10:35 GMT, Thursday, 24 December 2009

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

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

- News Front Page
- World
- UK
- England
- Northern Ireland
- Scotland
- Wales
- Business
- Politics
- Health
- Education
- Science & Environment
- Technology**
- Entertainment
- Also in the news

ML CAN BE RACIST...

[SWEENEY, 2013]

The Boston Globe

Metro Sports Business & Tech Opinion Politics Lifestyle Arts

Menu **Business** SIGN UP NOW Get Globe.com newsletters delivered to your inbox

Racial bias alleged in Google's ad results

Names associated with blacks prompt link to arrest search

18

Ad related to latanya sweeney

Latanya Sweeney Truth
www.instantcheckmate.com/
Looking for Latanya Sweeney? Check Latanya Sweeney's Arrests.

Ads by Google

Latanya Sweeney, Arrested?
1) Enter Name and State. 2) Access Full Background Checks Instantly.
www.instantcheckmate.com/

Latanya Sweeney

checkmate

LATANYA SWEENEY
1983 Latona Ave
Pittsburgh, PA 15216
DOB: Oct 27, 1983 (30 years old)

Criminal History

This section contains possible citation, arrest, and criminal records for the subject of this report. While our database does contain hundreds of millions of arrest records, different courts have different rules regarding what information they will and will not release.

We strive with you as much information as we possibly can, but a clean state here should not be interpreted as a guarantee that Latanya Sweeney has never been arrested. A empty means that we were not able to locate any existing arrest records in the state that is available to us.

Possible Matching Arrest Records

Name	County and State	Offense	View Details
No matching arrest records were found.			

ML CAN BE SEXIST...

[BOLUKBASI ET AL. 2016]



DEFINITIONS OF FAIRNESS

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]

Counterfactual Fairness
[Kusner et al., 2017]

Path-Specific Fairness
[Shpitser et al., 2017; Chiappa et al., 2018]

PROBLEM #1

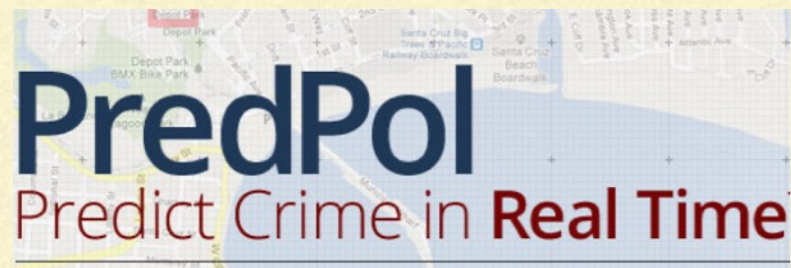
The Discriminatory **Prediction Problem**



ML IS INVOLVED IN LIFE-CHANGING DECISIONS

Policing

[Ensign et al., 2017]



Parole Sentencing

[Larson et al., 2016]



Advertising



Insurance



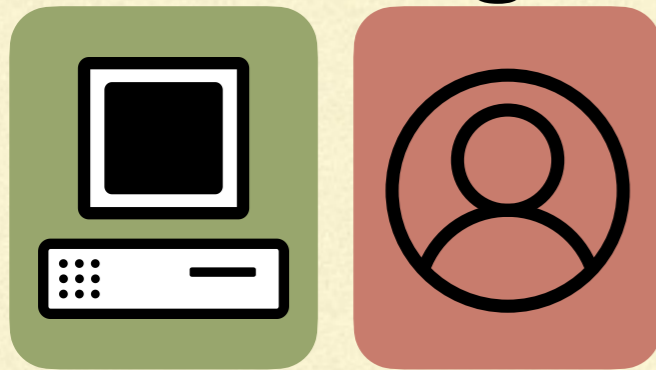
Lending



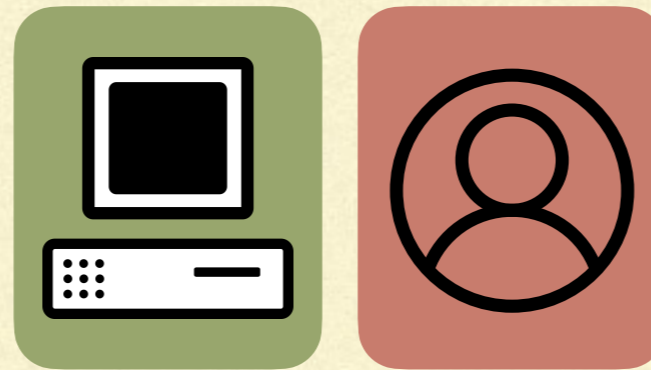
**Insert your
favorite
application
here!**

ML IS INVOLVED IN LIFE-CHANGING DECISIONS

Policing



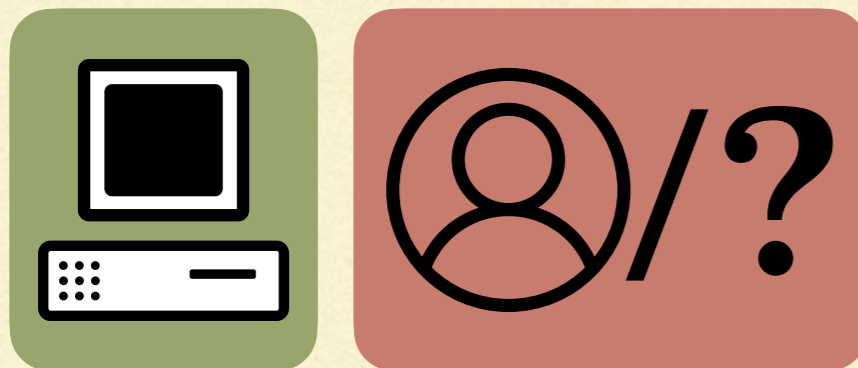
Parole Sentencing



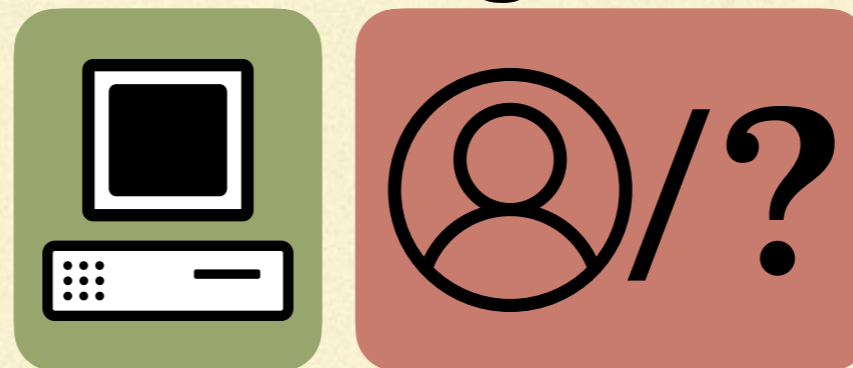
Advertising



Insurance



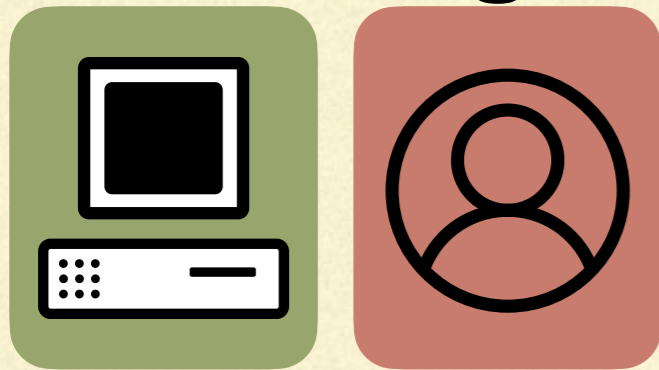
Lending



Insert your favorite application here!

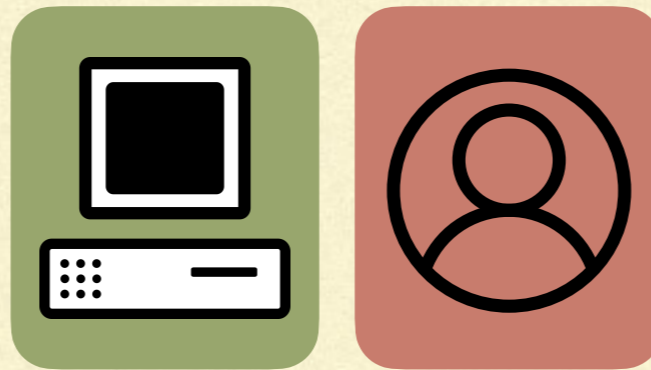
ML IS INVOLVED IN LIFE-CHANGING DECISIONS

Policing



impact: arrest

Parole Sentencing

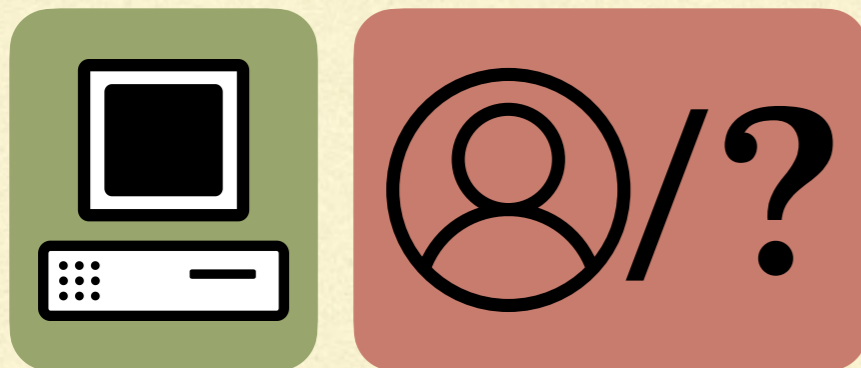


impact: jail-time

Advertising

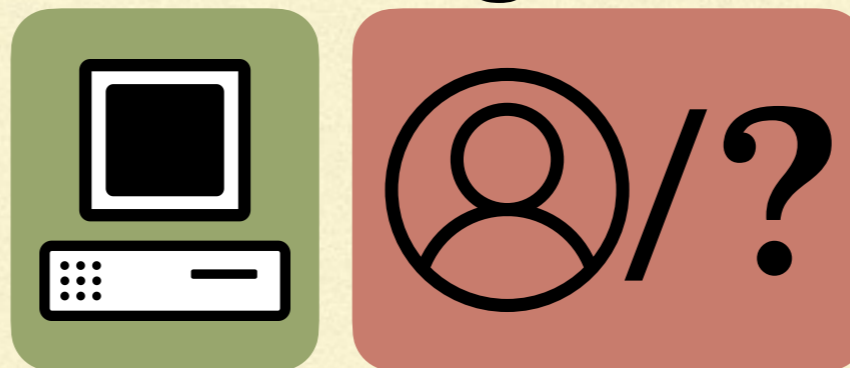


Insurance



impact: improved health

Lending

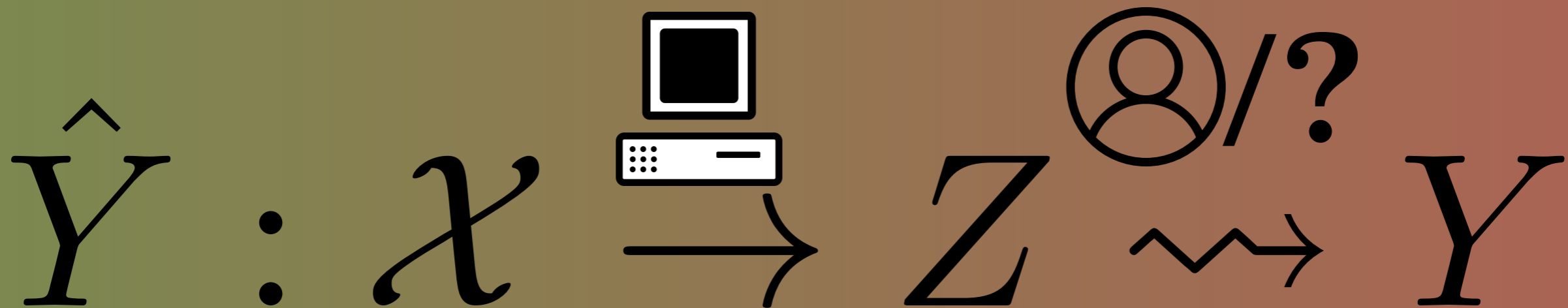


impact: pay off home loans

Insert your favorite application here!

PROBLEM #2

The Discriminatory **Impact** Problem

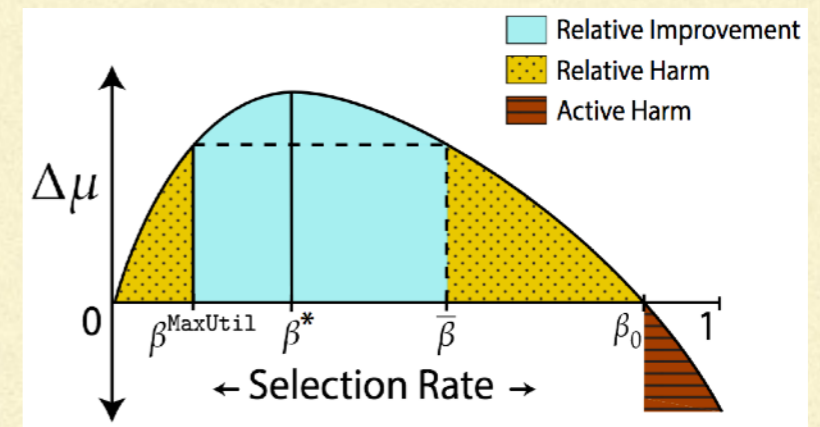


RELATED WORK

Introduction of discriminatory impact problem

[Liu et al., ICML 2018]

[Green & Chen, FAT* 2019]

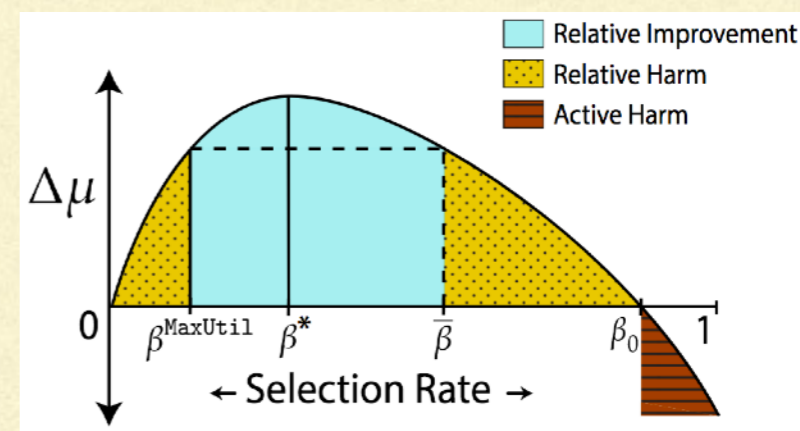


RELATED WORK

Introduction of discriminatory **impact** problem

[Liu et al., ICML 2018]

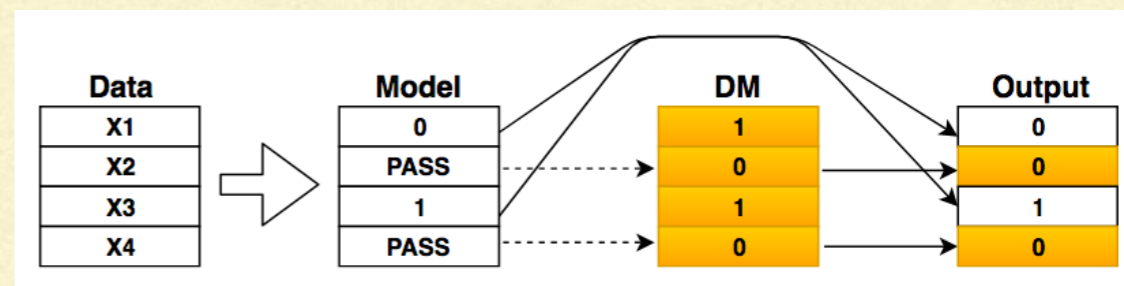
[Green & Chen, FAT* 2019]



Models for special cases

[Madras et al., NeurIPS 2018]

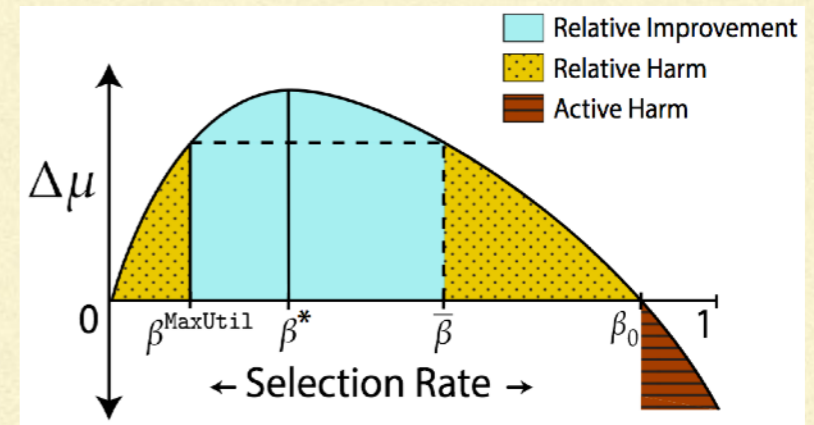
[Kannan et al., FAT* 2019]



RELATED WORK

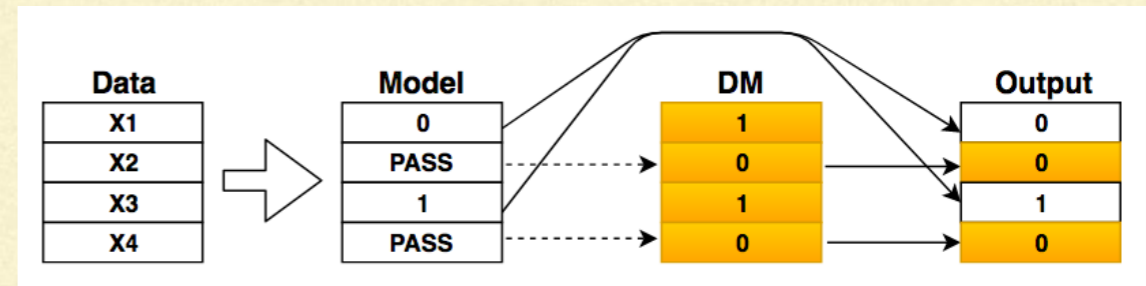
Introduction of discriminatory **impact** problem

[Liu et al., ICML 2018]
[Green & Chen, FAT* 2019]



Models for special cases

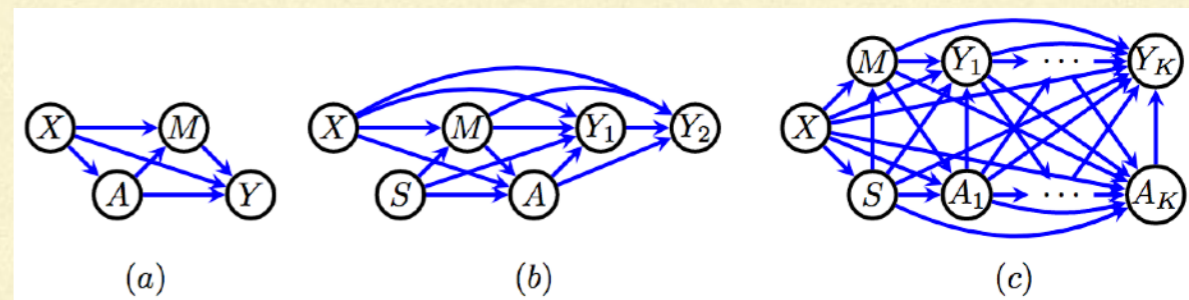
[Madras et al., NeurIPS 2018]
[Kannan et al., FAT* 2019]



Complimentary approaches

(RL, social dynamics)

[Nabi et al., ICML 2019]
[Hedari et al., ICML 2019]



RELATED WORK

our work: a general framework
for *reducing discriminatory impact*
based on causal modeling and MILP

ent

Output

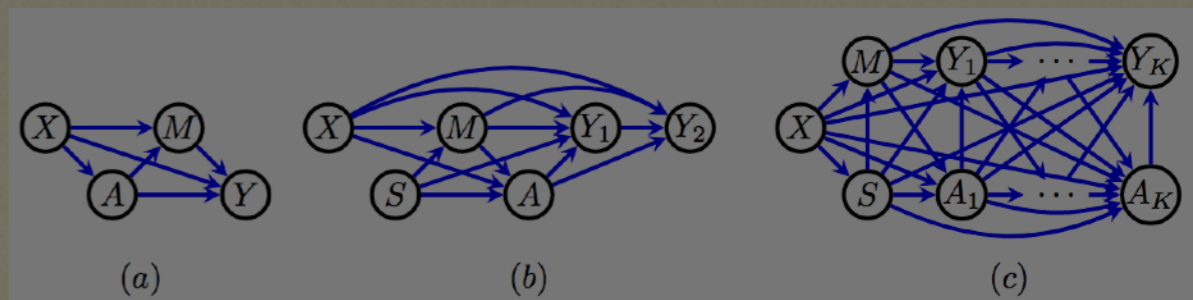
0
0
1
0

Complimentary approaches

(social dynamics, RL)

[Nabi et al., ICML 2019]

[Hedari et al., ICML 2019]



high school 1



high school 2



high school 1



high school 2



intervention: fund advanced classes
impact: to increase college applications
(% SAT/ACT-taking)

high school 1



race
distribution

$$A^{(1)}$$

$$X^{(1)}$$

counselors % SAT/ACT-taking

$$Y^{(1)}$$

high school 2



race
distribution

$$A^{(2)}$$

$$Y^{(2)}$$

% SAT/ACT-taking counselors

$$X^{(2)}$$

high school 1



race distribution intervention:
calculus classes

$A^{(1)}$

$Z^{(1)}$

$X^{(1)}$

$Y^{(1)}$

counselors % SAT/ACT-taking

high school 2



intervention:
calculus classes race
distribution

$Z^{(2)}$

$A^{(2)}$

$Y^{(2)}$

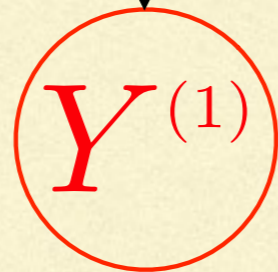
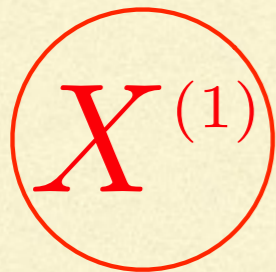
$X^{(2)}$

% SAT/ACT-taking counselors

high school 1



race distribution intervention:
calculus classes

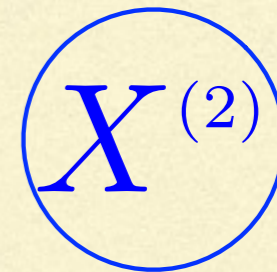
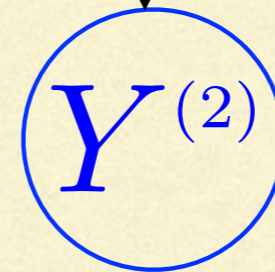
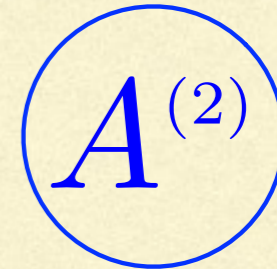
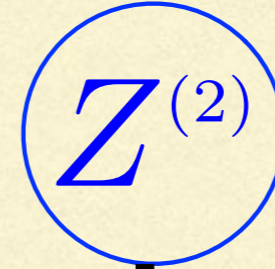


counselors % SAT/ACT-taking

high school 2



intervention:
calculus classes race
distribution

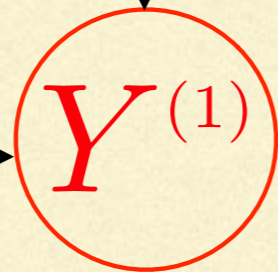
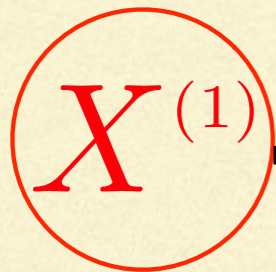


% SAT/ACT-taking counselors

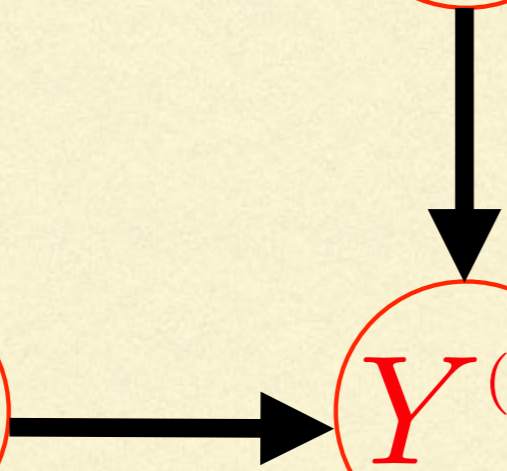
high school 1



race distribution intervention:
calculus classes



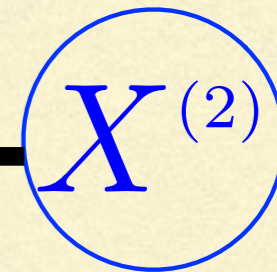
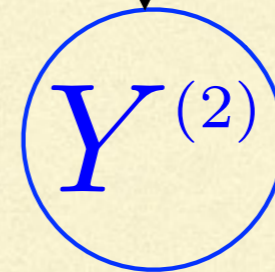
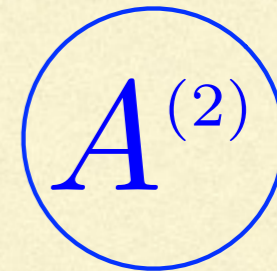
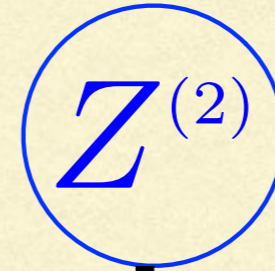
counselors % SAT/ACT-taking



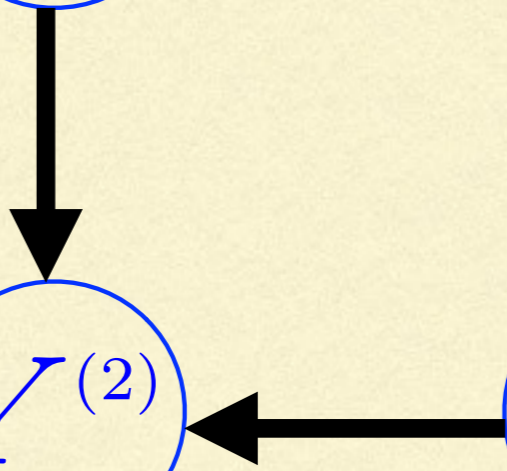
high school 2



intervention:
calculus classes race distribution



% SAT/ACT-taking counselors



high school 1



high school 2



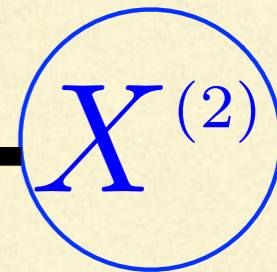
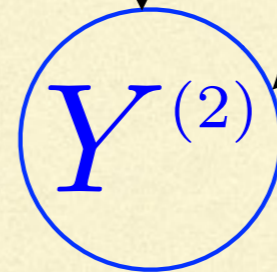
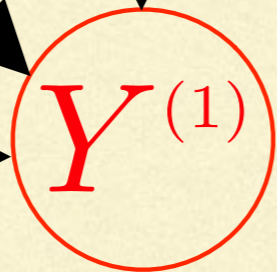
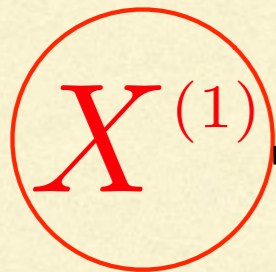
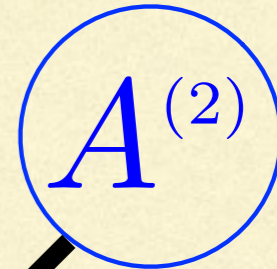
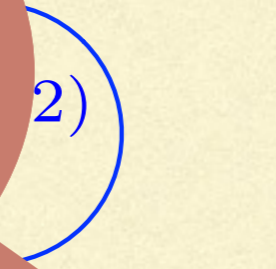
race distribution

intervention: calculus

intervention: classes

race distribution

tied to socio-economic status **due to economic history**
[Bittker, 2003]



counselors % SAT/ACT-taking

% SAT/ACT-taking counselors

high school 1



high school 2



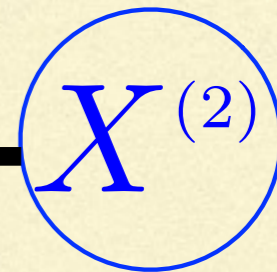
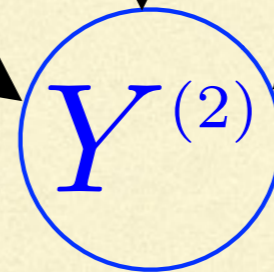
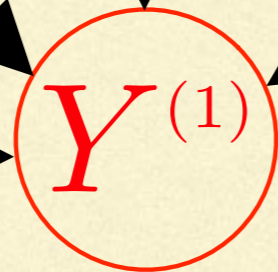
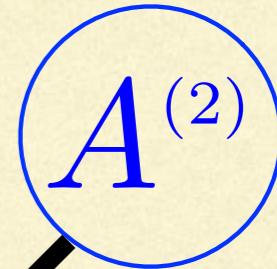
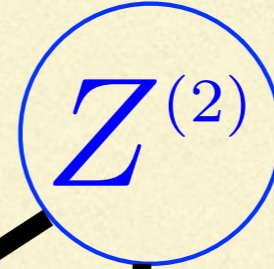
may cause students to **take classes at other schools**

race distribution

intervention: calculus classe

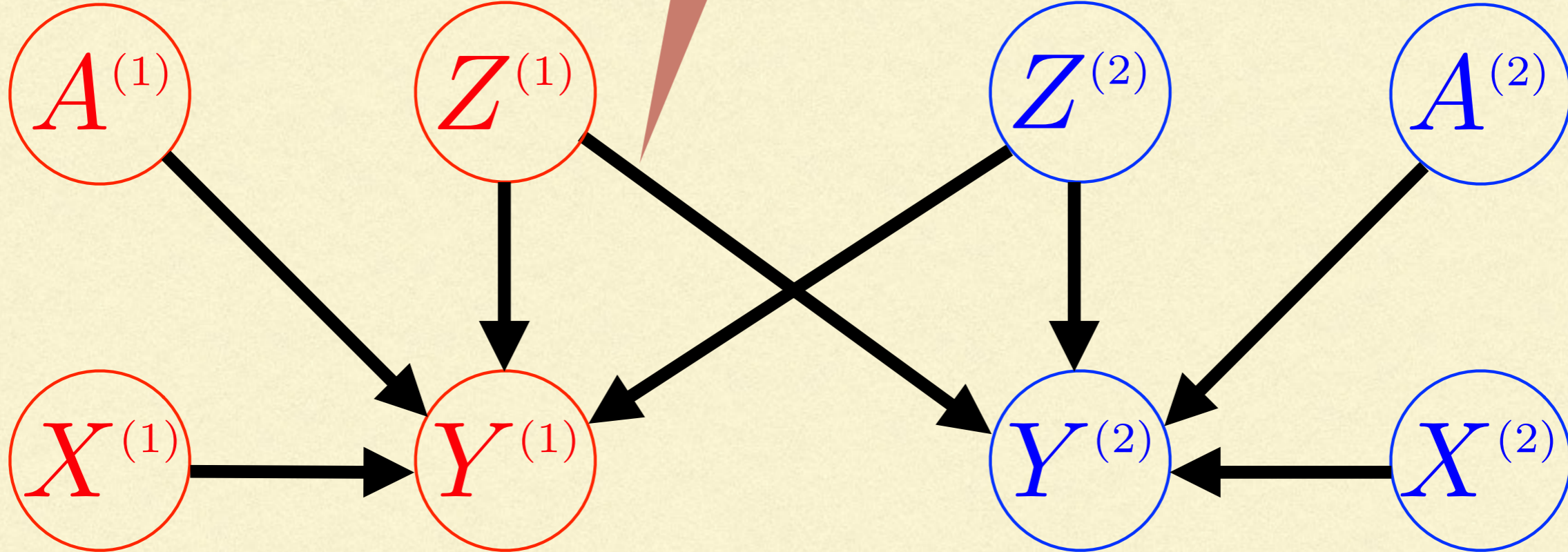
intervention: calculus classes

race distribution



counselors % SAT/ACT-taking

% SAT/ACT-taking counselors

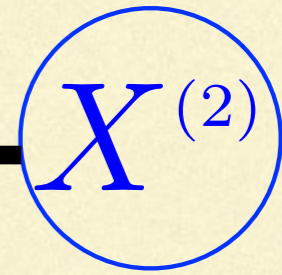
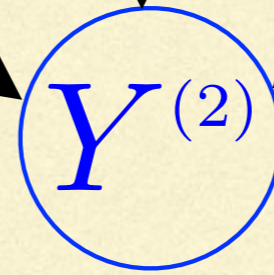
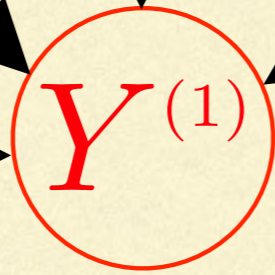
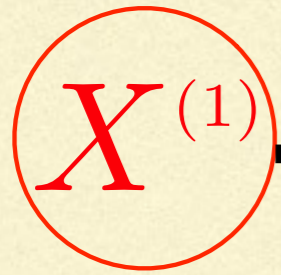
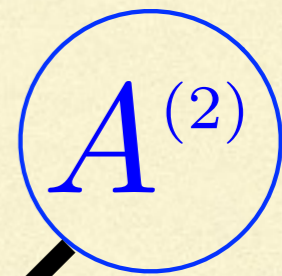
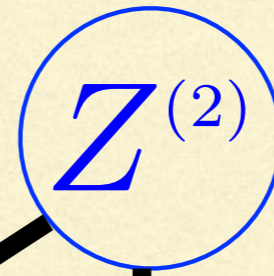
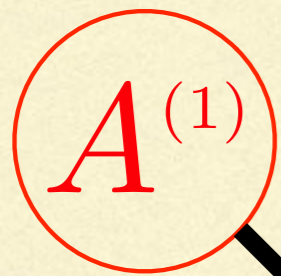


high school 1

high school 2

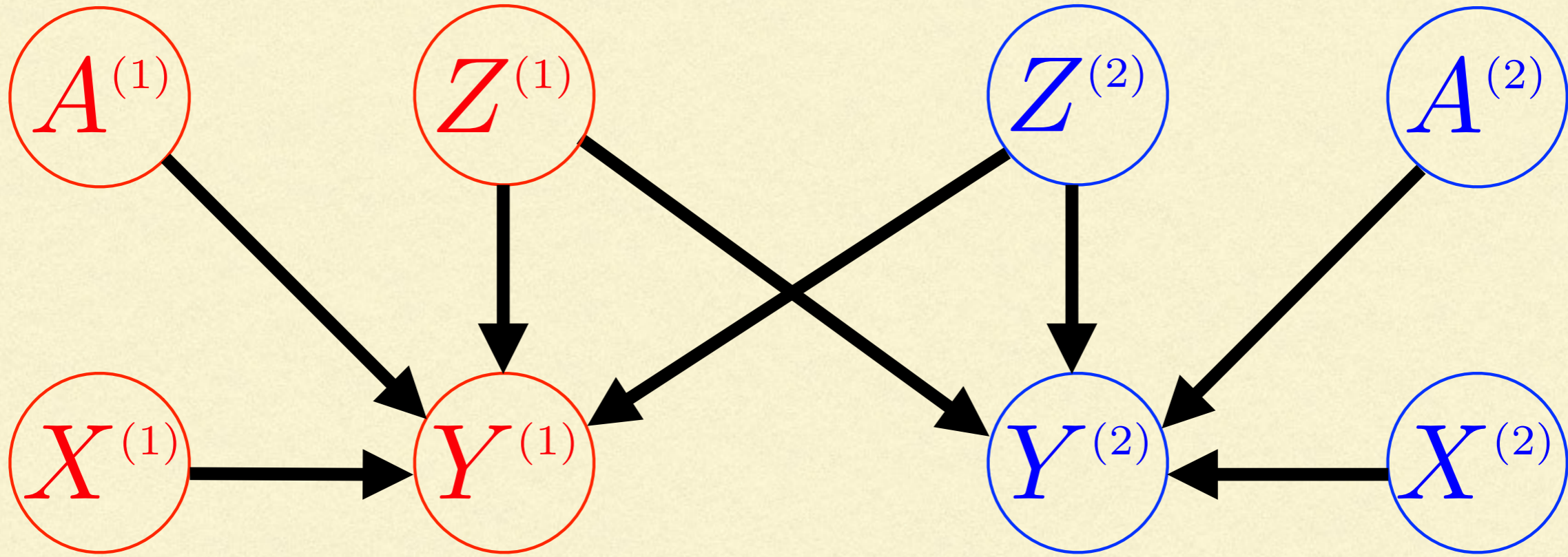
race
distribution intervention:
calculus classes

intervention:
calculus classes race
distribution



counselors % SAT/ACT-taking

% SAT/ACT-taking counselors



high school 1

high school 2

race
distribution

intervention:
calculus classes

intervention:
calculus classes

race
distribution

$A^{(1)}$

$A^{(1)} \neq A^{(2)}$

$A^{(2)}$

$X^{(1)} = X^{(2)}$

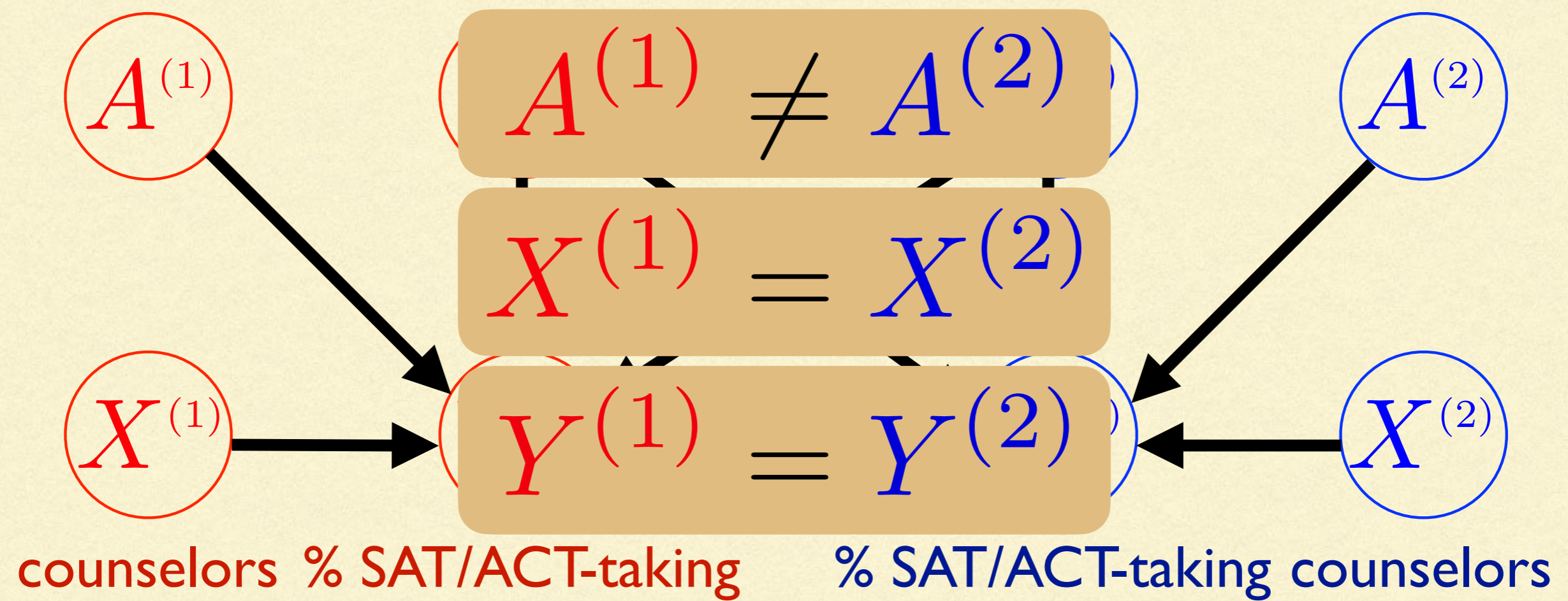
$X^{(1)}$

$Y^{(1)} = Y^{(2)}$

$X^{(2)}$

counselors % SAT/ACT-taking

% SAT/ACT-taking counselors

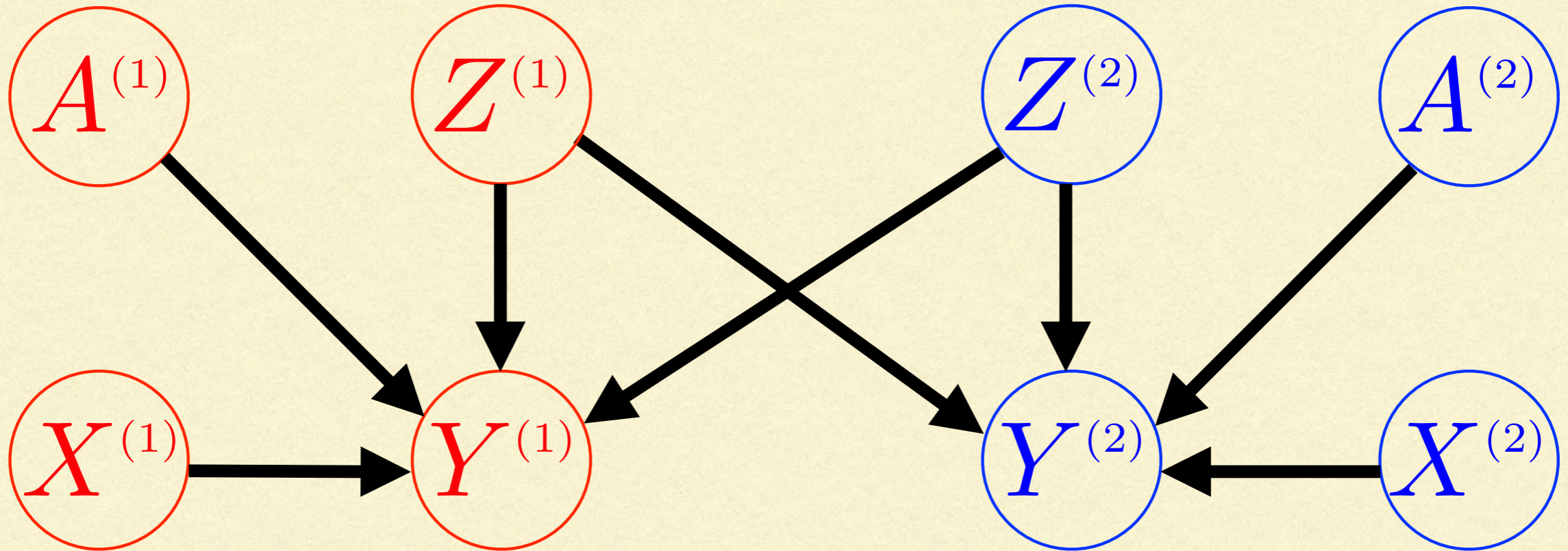


high school 1

high school 2

race distribution intervention:
calculus classes

intervention:
calculus classes race
distribution



counselors % SAT/ACT-taking

% SAT/ACT-taking counselors

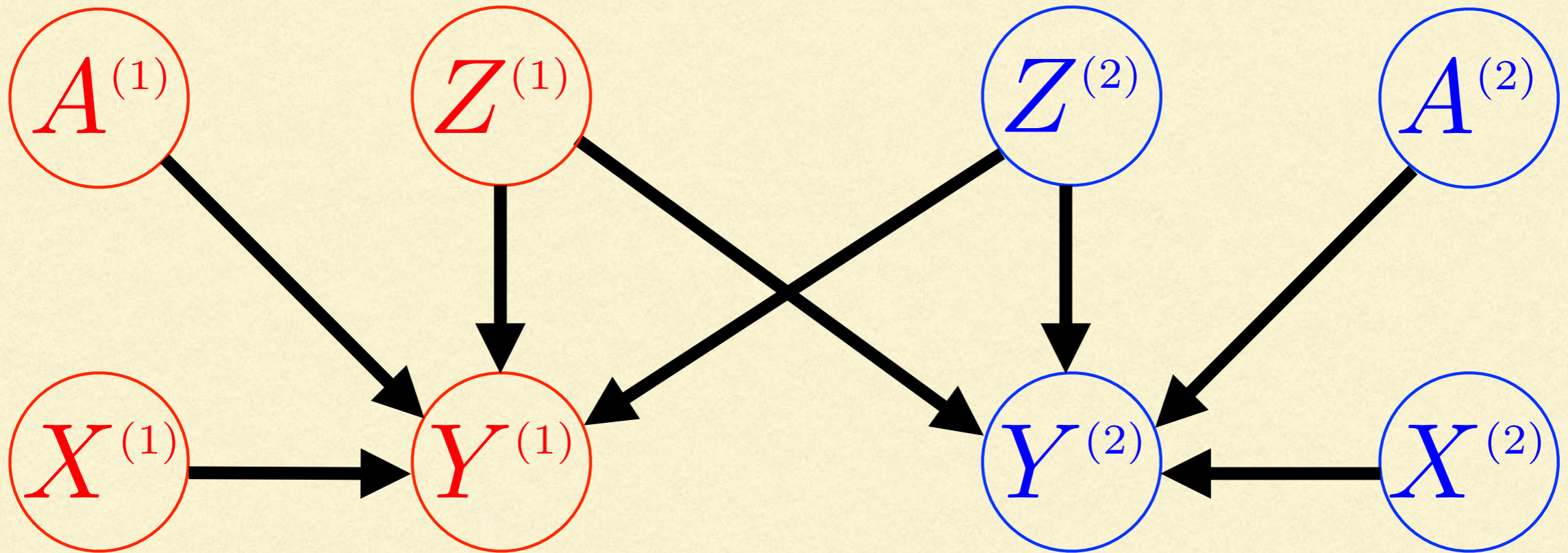
$$Y^{(1)}([z^{(1)} = 1, z^{(2)} = 0]) = 1.0$$

high school 1

high school 2

race distribution intervention:
calculus classes

intervention:
calculus classes race
distribution



counselors % SAT/ACT-taking

% SAT/ACT-taking counselors

$$Y^{(1)}([z^{(1)} = 1, z^{(2)} = 0]) = 1.0$$

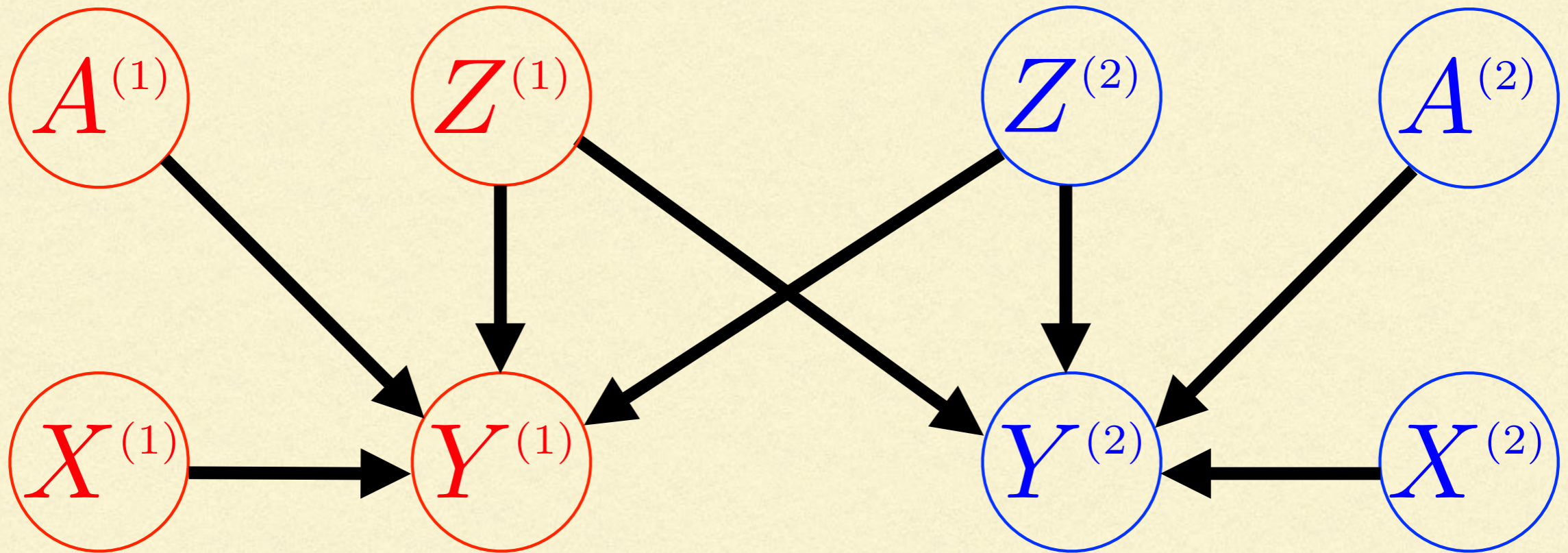
$$Y^{(2)}([z^{(1)} = 1, z^{(2)} = 0]) = 0.1$$

high school 1

high school 2

race distribution intervention:
calculus classes

intervention:
calculus classes race
distribution



counselors % SAT/ACT-taking

% SAT/ACT-taking counselors

$$Y^{(1)}([z^{(1)} = 1, z^{(2)} = 0]) = 1.0$$

$$Y^{(2)}([z^{(1)} = 1, z^{(2)} = 0]) = 0.1$$

$$Y^{(1)}([z^{(1)} = 0, z^{(2)} = 1]) = 0.6$$

$$Y^{(2)}([z^{(1)} = 0, z^{(2)} = 1]) = 0.5$$

high school 1

high school 2

race
distribution

decision:
calculus classes

decision:
calculus classes

race
distribution

$A^{(1)}$

$Z^{(1)}$

$Z^{(2)}$

$A^{(2)}$

either intervention causes the same
average overall impact

$X^{(1)}$

but it seems unfair to give classes to
school 1 as they have better impact
solely due to race

$X^{(2)}$

counselors % SAT/ACT-taking

% SAT/ACT-taking counselors

$$Y^{(1)}([z^{(1)} = 1, z^{(2)} = 0]) = 1.0$$

$$Y^{(2)}([z^{(1)} = 1, z^{(2)} = 0]) = 0.1$$

$$Y^{(1)}([z^{(1)} = 0, z^{(2)} = 1]) = 0.6$$

$$Y^{(2)}([z^{(1)} = 0, z^{(2)} = 1]) = 0.5$$

MORE FORMALLY

Maximizing overall impact

$$\begin{aligned} \max_{\mathbf{z} \in \{0,1\}^n} & \sum_{i=1}^n \mathbb{E}[Y^{(i)}(a^{(i)}, \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}], \\ \text{s.t.}, & \sum_{i=1}^n z^{(i)} \leq b \end{aligned}$$

MORE FORMALLY

Maximizing overall impact

whether to grant advanced classes

of schools

$$\max_{\mathbf{z} \in \{0,1\}^n} \sum_{i=1}^n \mathbb{E}[Y^{(i)}(a^{(i)}, \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}]$$

$$s.t., \sum_{i=1}^n z^{(i)} \leq b$$

government budget

QUANTIFYING PRIVILEGE

$$\mathbb{E}[Y^{(i)}(a^{(i)}, \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}]$$
$$\mathbb{E}[Y^{(i)}(a', \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}]$$

QUANTIFYING PRIVILEGE

$$\mathbb{E}[Y^{(i)}(a^{(i)}, \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}]$$
$$\mathbb{E}[Y^{(i)}(a', \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}]$$

a counterfactual:
impact school i *would*
have gotten for
interventions \mathbf{z} if race
distribution was a'

QUANTIFYING PRIVILEGE

$$\mathbb{E}[Y^{(i)}(a^{(i)}, \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] - \mathbb{E}[Y^{(i)}(a', \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] < \tau$$

a counterfactual:
impact school i *would*
have gotten for
interventions \mathbf{z} if race
distribution was a'

QUANTIFYING PRIVILEGE

constraint on **counterfactual privilege**

$$\mathbb{E}[Y^{(i)}(a^{(i)}, \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] - \mathbb{E}[Y^{(i)}(a', \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] < \tau$$

a counterfactual:
impact school i *would*
have gotten for
interventions \mathbf{z} if race
distribution was a'

MORE FORMALLY

Maximizing impact with privilege constraints

$$\max_{\mathbf{z} \in \{0,1\}^n} \sum_{i=1}^n \mathbb{E}[Y^{(i)}(a^{(i)}, \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}]$$

$$s.t., \sum_{i=1}^n z^{(i)} \leq b$$

$$C_{ia'} \leq \tau \quad \forall a' \in \mathcal{A}, i \in \{1, \dots, n\},$$

$$\begin{aligned} & \mathbb{E}[Y^{(i)}(a^{(i)}, \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] \\ & - \mathbb{E}[Y^{(i)}(a', \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] \end{aligned}$$

$$\mathbb{E}[Y^{(i)}(a^{(i)}, \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] - \mathbb{E}[Y^{(i)}(a', \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] < \tau$$

$$\tau = 0$$

$$Y^{(1)}([z^{(1)} = 1, z^{(2)} = 0]) = 1.0$$

$$Y^{(2)}([z^{(1)} = 1, z^{(2)} = 0]) = 0.1$$

$$Y^{(1)}([z^{(1)} = 0, z^{(2)} = 1]) = 0.6$$

$$Y^{(2)}([z^{(1)} = 0, z^{(2)} = 1]) = 0.5$$

high school 1

high school 2

$$\mathbb{E}[Y^{(i)}(a^{(i)}, \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] - \mathbb{E}[Y^{(i)}(a', \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] < \tau$$

$$\tau = 0$$

essentially counterfactuals!

$$Y^{(1)}([z^{(1)} = 1, z^{(2)} = 0]) = 1.0$$

$$Y^{(2)}([z^{(1)} = 1, z^{(2)} = 0]) = 0.1$$

$$Y^{(1)}([z^{(1)} = 0, z^{(2)} = 1]) = 0.6$$

$$Y^{(2)}([z^{(1)} = 0, z^{(2)} = 1]) = 0.5$$

high school 1

high school 2

$$\mathbb{E}[Y^{(i)}(a^{(i)}, \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] - \mathbb{E}[Y^{(i)}(a', \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] < \tau$$

$$\tau = 0$$

school 1
gets classes

$$Y^{(1)}([z^{(1)} = 1, z^{(2)} = 0]) - Y^{(2)}([z^{(1)} = 1, z^{(2)} = 0]) = 0.9$$

$$Y^{(1)}([z^{(1)} = 1, z^{(2)} = 0]) = 1.0$$

$$Y^{(2)}([z^{(1)} = 1, z^{(2)} = 0]) = 0.1$$

$$Y^{(1)}([z^{(1)} = 0, z^{(2)} = 1]) = 0.6$$

$$Y^{(2)}([z^{(1)} = 0, z^{(2)} = 1]) = 0.5$$

high school 1

high school 2

$$\mathbb{E}[Y^{(i)}(a^{(i)}, \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] - \mathbb{E}[Y^{(i)}(a', \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] < \tau$$

$$\tau = 0$$

school 1
gets classes

$$Y^{(1)}([z^{(1)} = 1, z^{(2)} = 0]) - Y^{(2)}([z^{(1)} = 1, z^{(2)} = 0]) = 0.9$$

school 2
gets classes

$$Y^{(2)}([z^{(1)} = 0, z^{(2)} = 1]) - Y^{(1)}([z^{(1)} = 0, z^{(2)} = 1]) = -0.1$$

$$Y^{(1)}([z^{(1)} = 1, z^{(2)} = 0]) = 1.0$$

$$Y^{(2)}([z^{(1)} = 1, z^{(2)} = 0]) = 0.1$$

$$Y^{(1)}([z^{(1)} = 0, z^{(2)} = 1]) = 0.6$$

$$Y^{(2)}([z^{(1)} = 0, z^{(2)} = 1]) = 0.5$$

high school 1

high school 2

$$\mathbb{E}[Y^{(i)}(a^{(i)}, \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] - \mathbb{E}[Y^{(i)}(a', \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] < \tau$$

$$\tau = 0$$

school 1
gets classes

$$Y^{(1)}([z^{(1)} = 1, z^{(2)} = 0]) - Y^{(2)}([z^{(1)} = 1, z^{(2)} = 0]) = 0.9$$



school 2
gets classes

$$Y^{(2)}([z^{(1)} = 0, z^{(2)} = 1]) - Y^{(1)}([z^{(1)} = 0, z^{(2)} = 1]) = -0.1$$



$$Y^{(1)}([z^{(1)} = 1, z^{(2)} = 0]) = 1.0$$

$$Y^{(2)}([z^{(1)} = 1, z^{(2)} = 0]) = 0.1$$

$$Y^{(1)}([z^{(1)} = 0, z^{(2)} = 1]) = 0.6$$

$$Y^{(2)}([z^{(1)} = 0, z^{(2)} = 1]) = 0.5$$

high school 1

high school 2

VS. COUNTERFACTUAL FAIRNESS

constraint on **counterfactual privilege**

$$\mathbb{E}[Y^{(i)}(a^{(i)}, \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] - \mathbb{E}[Y^{(i)}(a', \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] < \tau$$

VS.

$$P(\hat{Y}^{(i)}(a^{(i)}) = y \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}) = P(\hat{Y}^{(i)}(a') = y \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)})$$

counterfactual fairness

[Kusner et al., 2017]

OPTIMIZATION: MILP

$$\begin{aligned} & \max_{\substack{\mathbf{z} \in \{0,1\}^n \\ \mathbf{H} \in [0,1]^{(n,2^K)}}} \sum_{i=1}^n \sum_{j=1}^{2^K} h_{ij} \xi^{ij}(a^{(i)}) \\ & \text{s.t.}, \sum_{j=1}^{2^K} h_{i,j} \left[\xi_{\prec}^{ij}(a^{(i)}) - \xi_{\prec}^{ij}(a') \right] < \tau, \quad \forall a', i \\ & \mathbb{I}[\mathbf{e}_j = 1] h_{ij} \leq \mathbf{z}^{N(i)}, \quad \forall i, j \\ & \mathbb{I}[\mathbf{e}_j = 0] h_{ij} \leq 1 - \mathbf{z}^{N(i)}, \quad \forall i, j \\ & \sum_{j=1}^{2^K} h_{ij} = 1, \quad \forall i \\ & \sum_{i=1}^n z^{(i)} \leq b. \end{aligned}$$

OPTIMIZATION: MILP

can accommodate any formulation of impact

$$\max_{\substack{\mathbf{z} \in \{0,1\}^n \\ \mathbf{H} \in [0,1]^{(n,2^K)}}} \sum_{i=1}^n \sum_{j=1}^{2^K} h_{ij} \xi^{ij}(a^{(i)})$$

$$\xi^{ij}(a^{(i)}) := \mathbb{E}[Y^{(i)}(a^{(i)}, \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}]$$

$$s.t., \sum_{j=1}^{2^K} h_{i,j} \left[\xi_{\prec}^{ij}(a^{(i)}) - \xi_{\prec}^{ij}(a') \right] < \tau, \quad \forall a', i$$

$$\mathbb{I}[\mathbf{e}_j = 1] h_{ij} \leq \mathbf{z}^{N(i)}, \quad \forall i, j$$

$$\mathbb{I}[\mathbf{e}_j = 0] h_{ij} \leq 1 - \mathbf{z}^{N(i)}, \quad \forall i, j$$

$$\sum_{j=1}^{2^K} h_{ij} = 1, \quad \forall i$$

$$\sum_{i=1}^n z^{(i)} \leq b.$$

NYC PUBLIC SCHOOL FUNDING

[CRDC, [HTTPS://OCRDATA.ED.GOV/](https://ocrdata.ed.gov/)]

345 schools



school 1

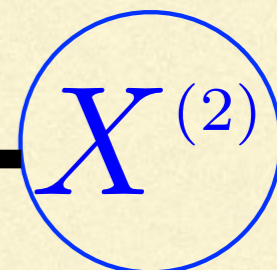
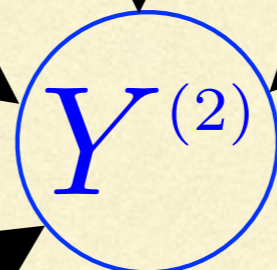
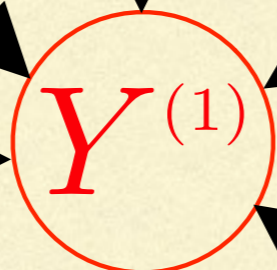
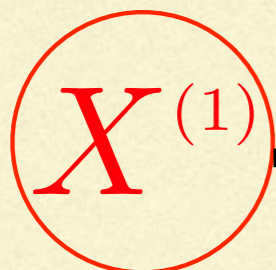
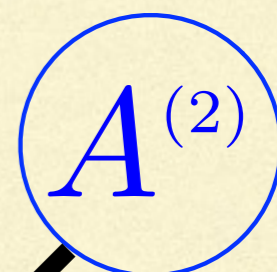
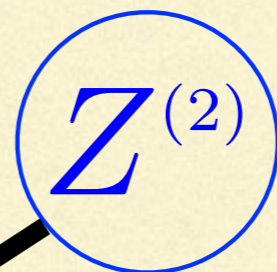
school 2

race
distribution

intervention:
calculus classes

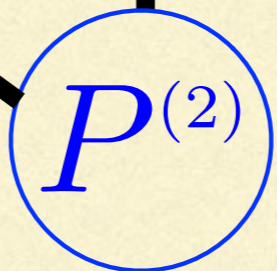
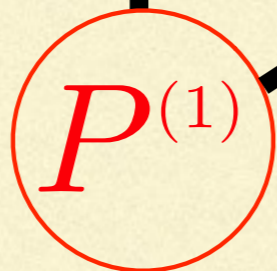
intervention:
calculus classes

race
distribution



counselors

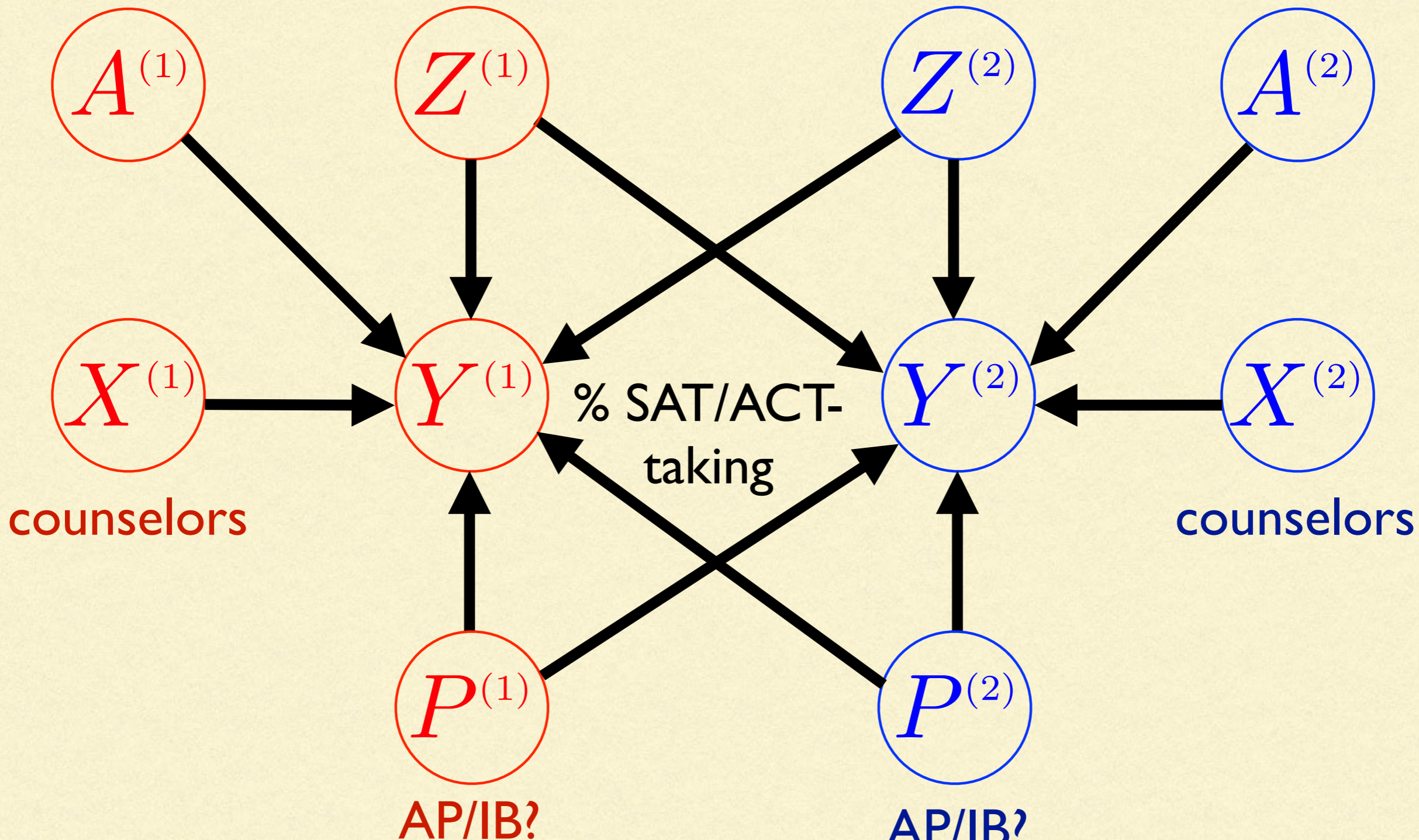
counselors



AP/IB?

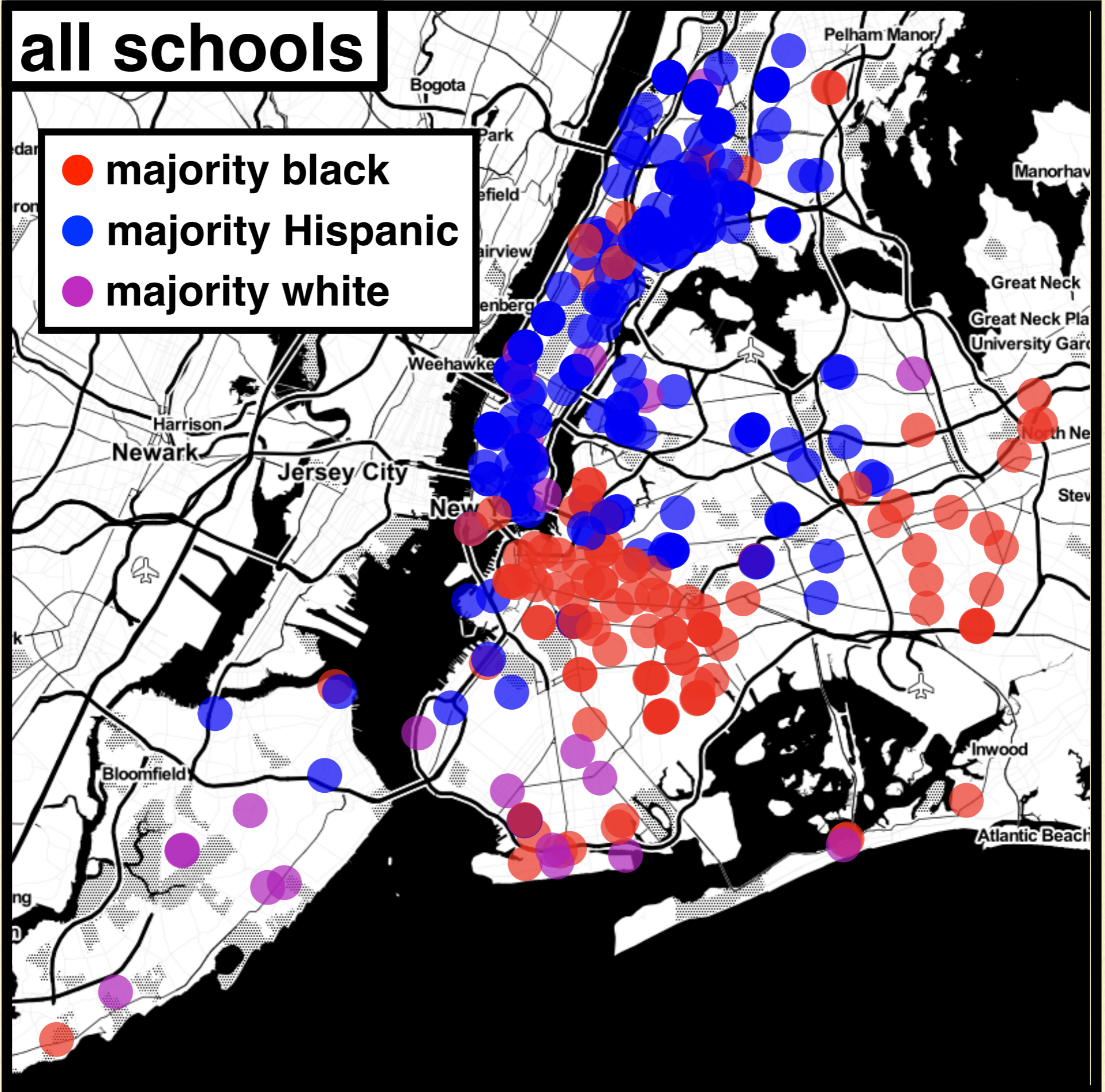
AP/IB?

% SAT/ACT-
taking



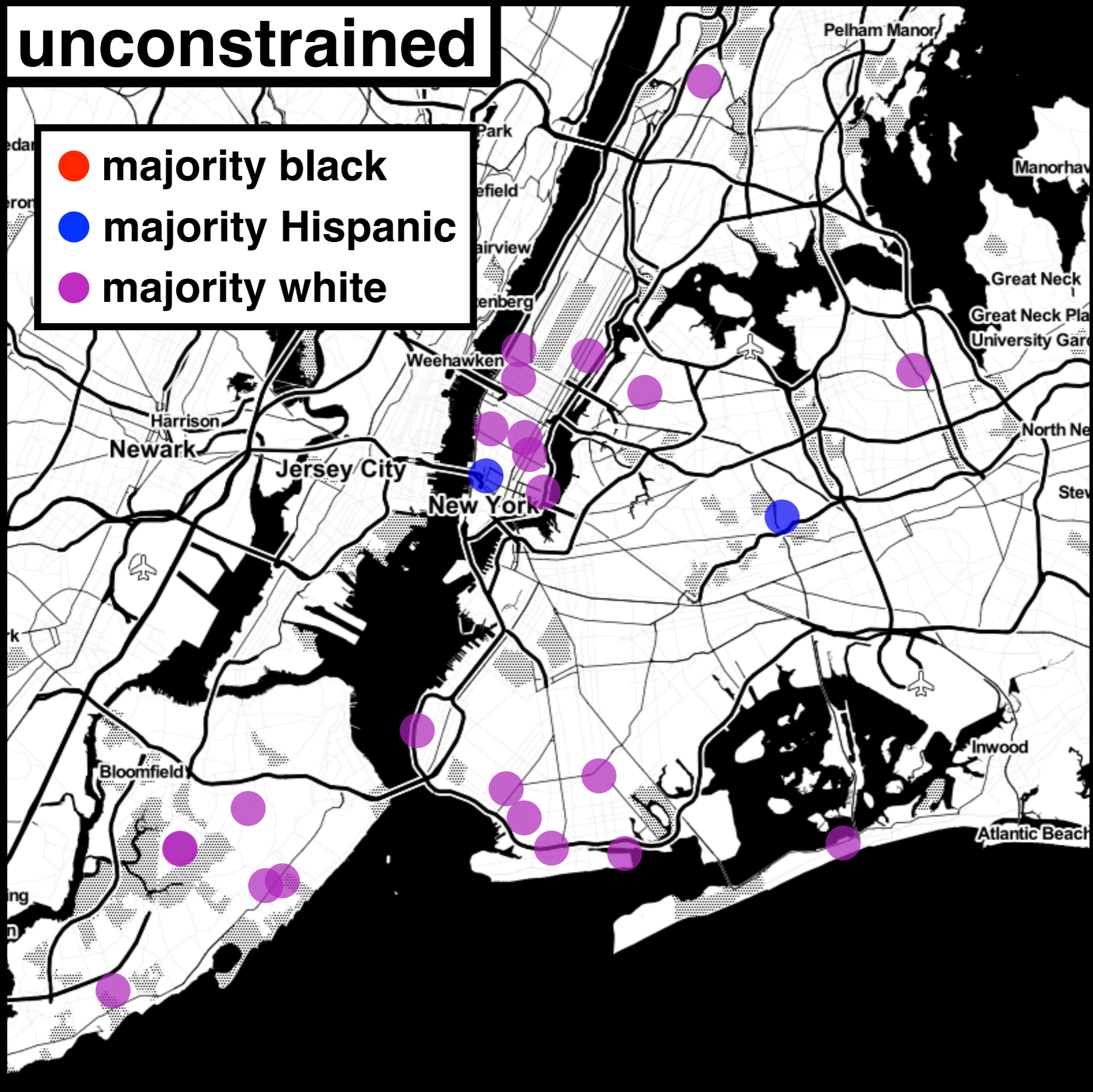
all schools

- majority black
- majority Hispanic
- majority white

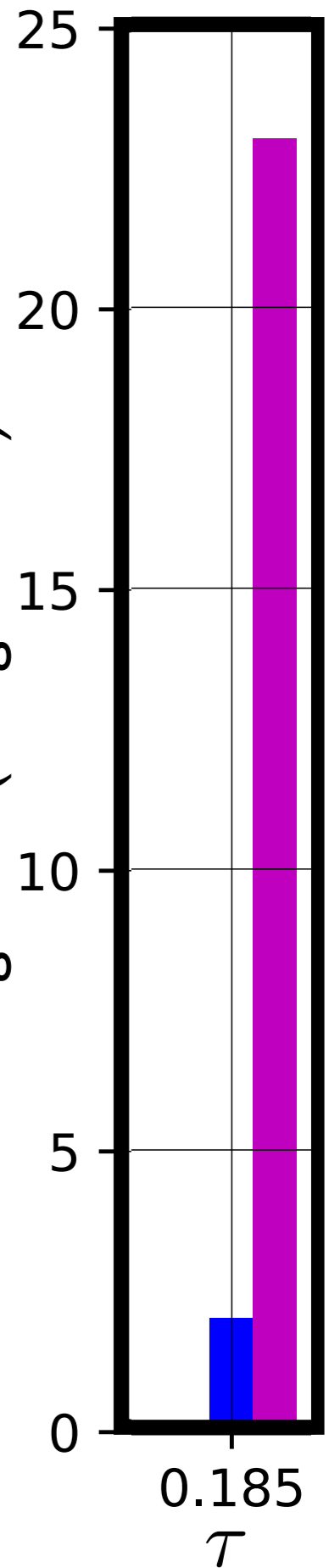


unconstrained

- majority black
- majority Hispanic
- majority white



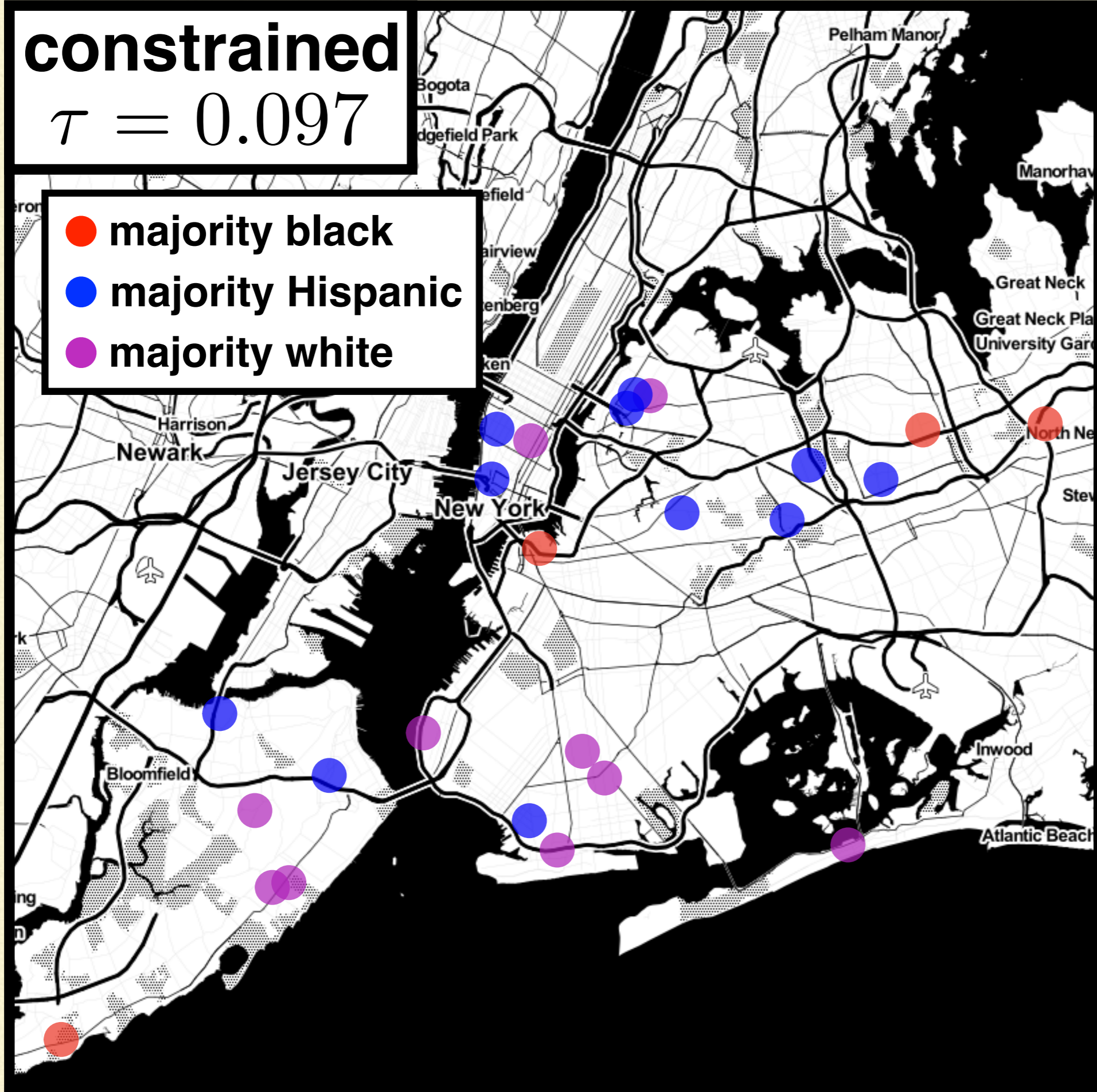
of grants (budget = 25)



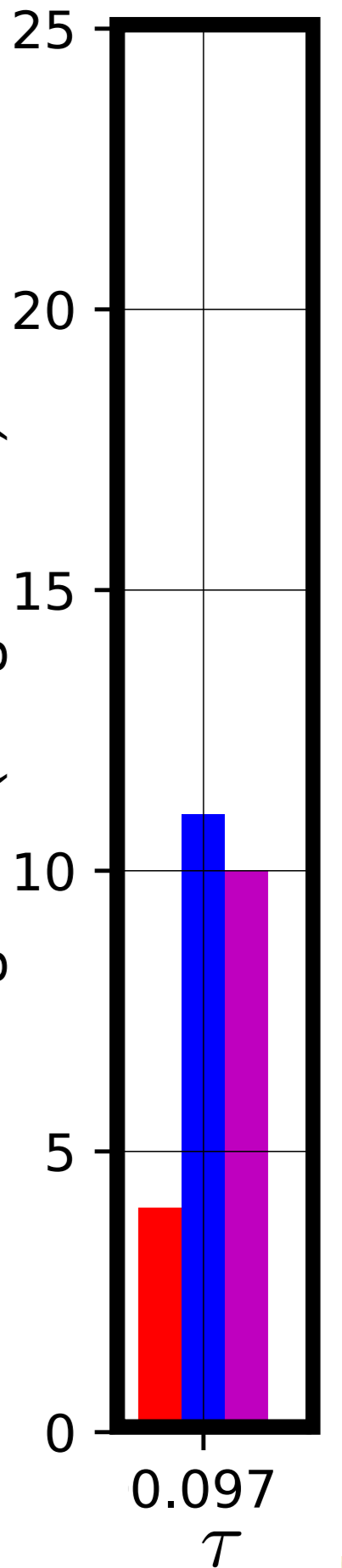
constrained

$$\tau = 0.097$$

- majority black
- majority Hispanic
- majority white



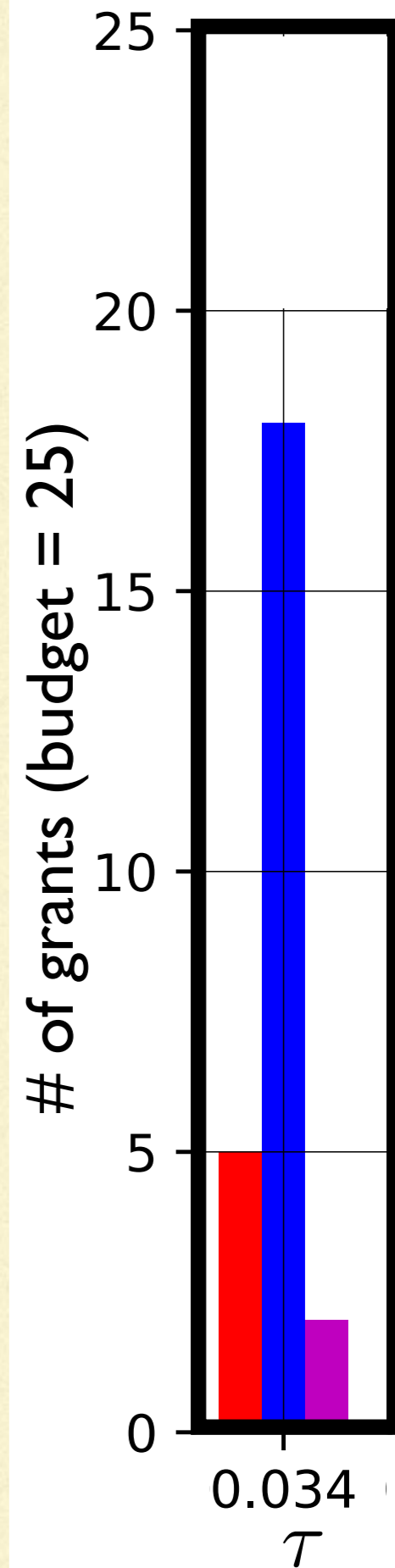
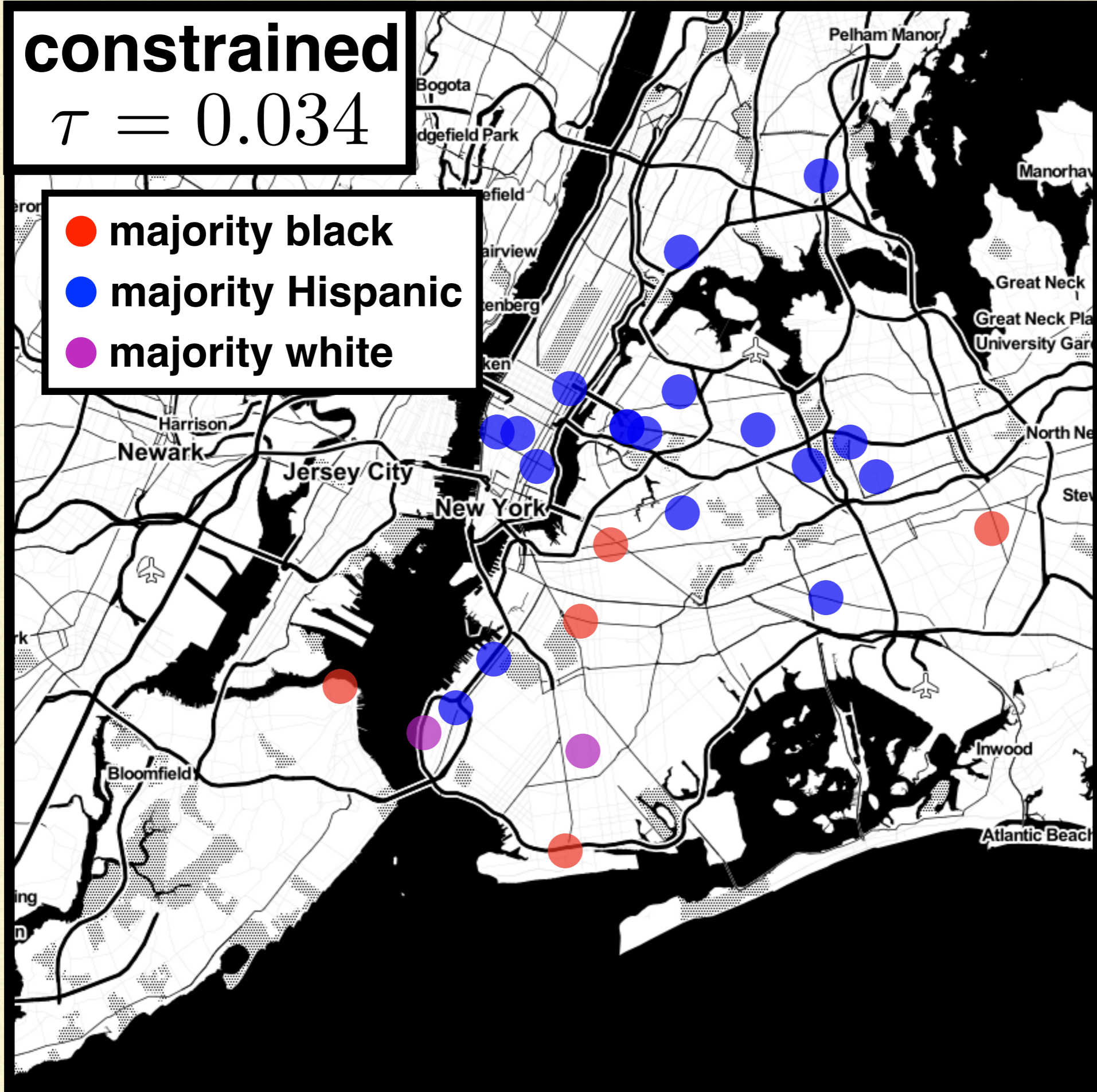
of grants (budget = 25)

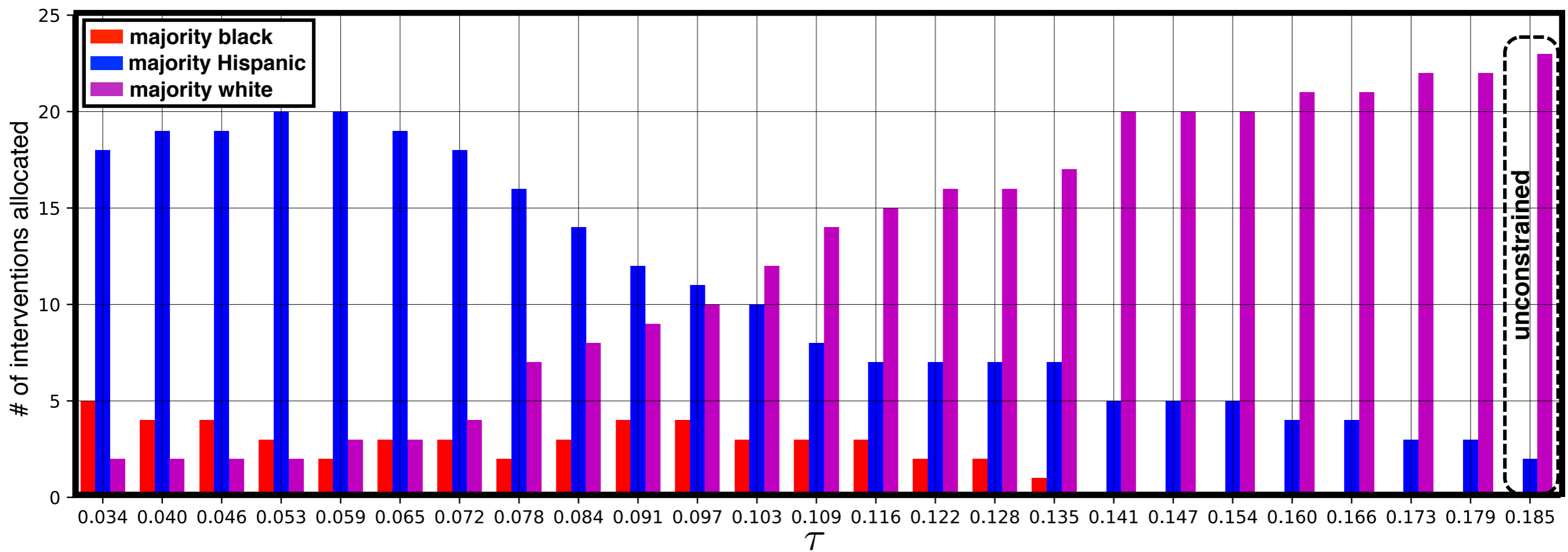


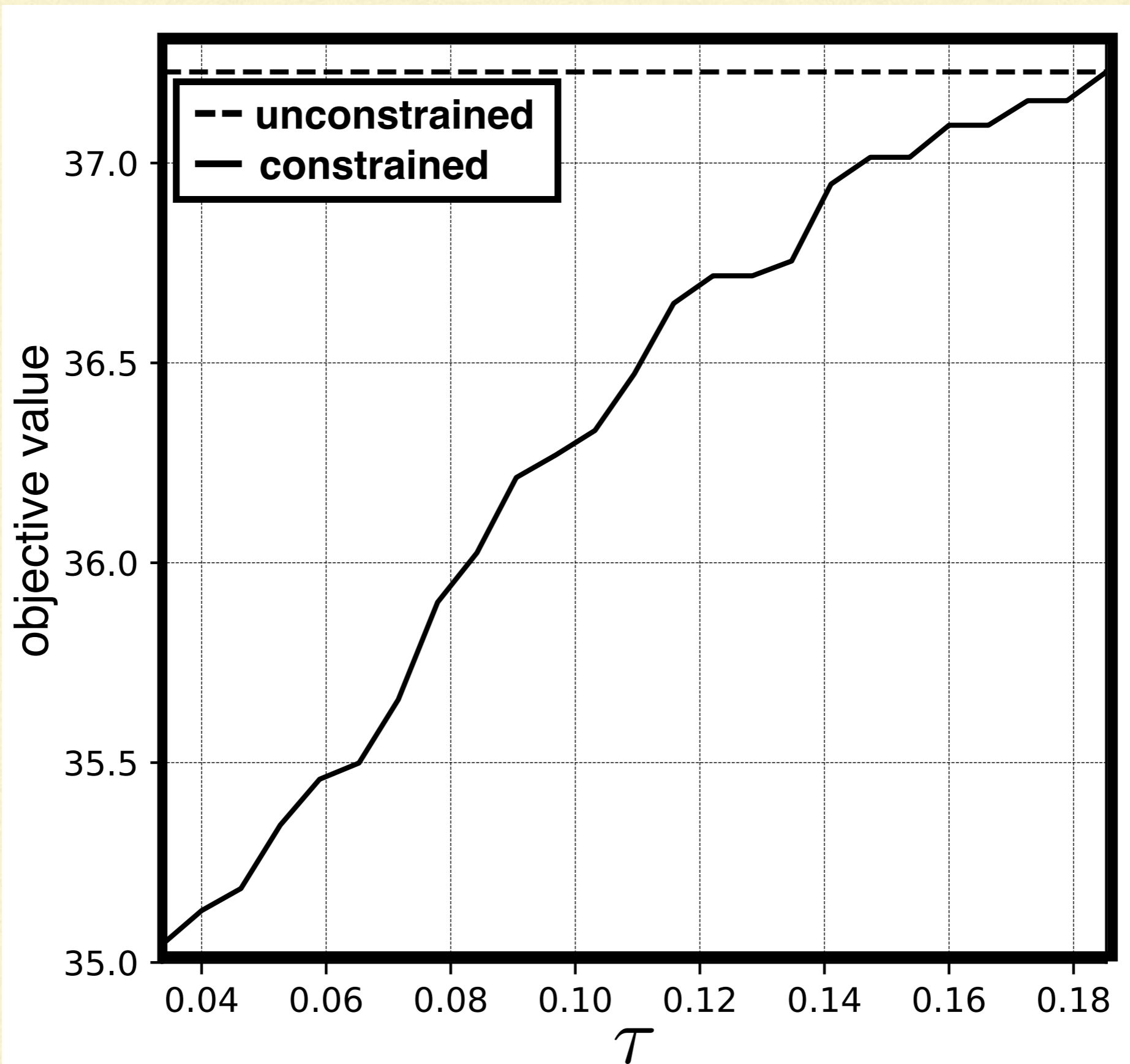
constrained

$$\tau = 0.034$$

- majority black
- majority Hispanic
- majority white



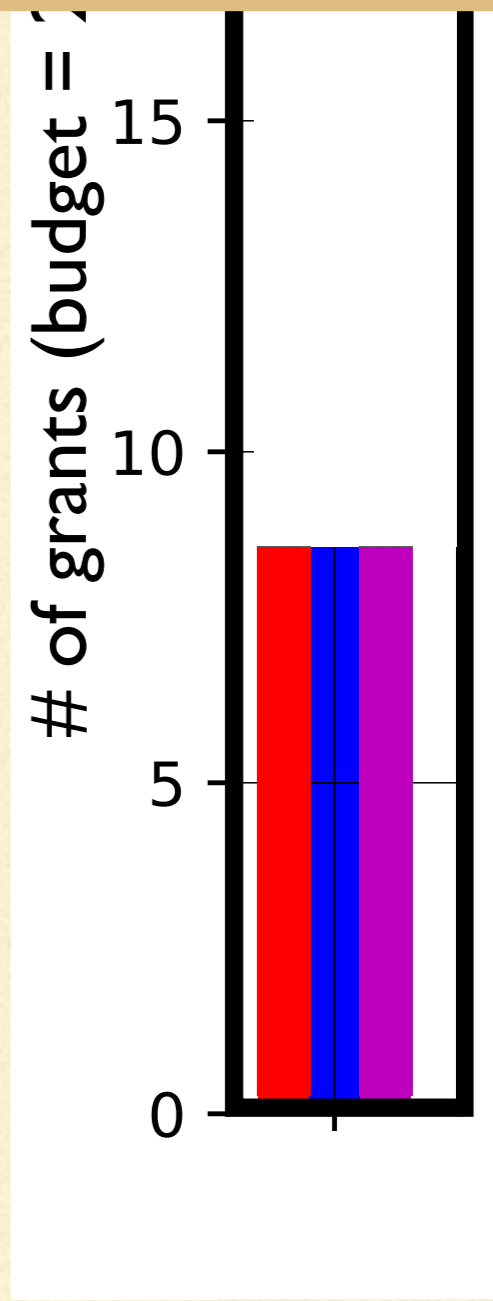
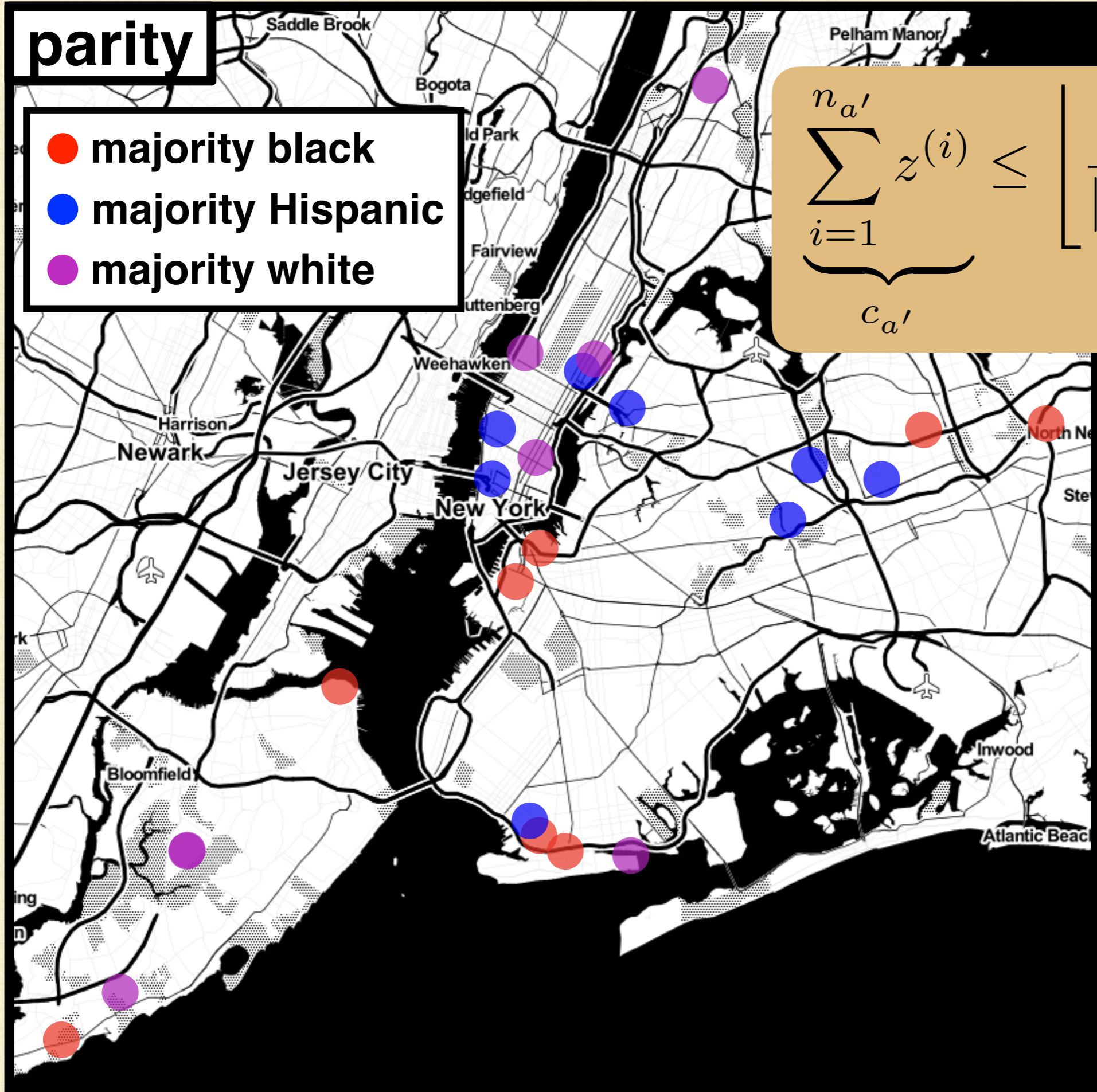




parity

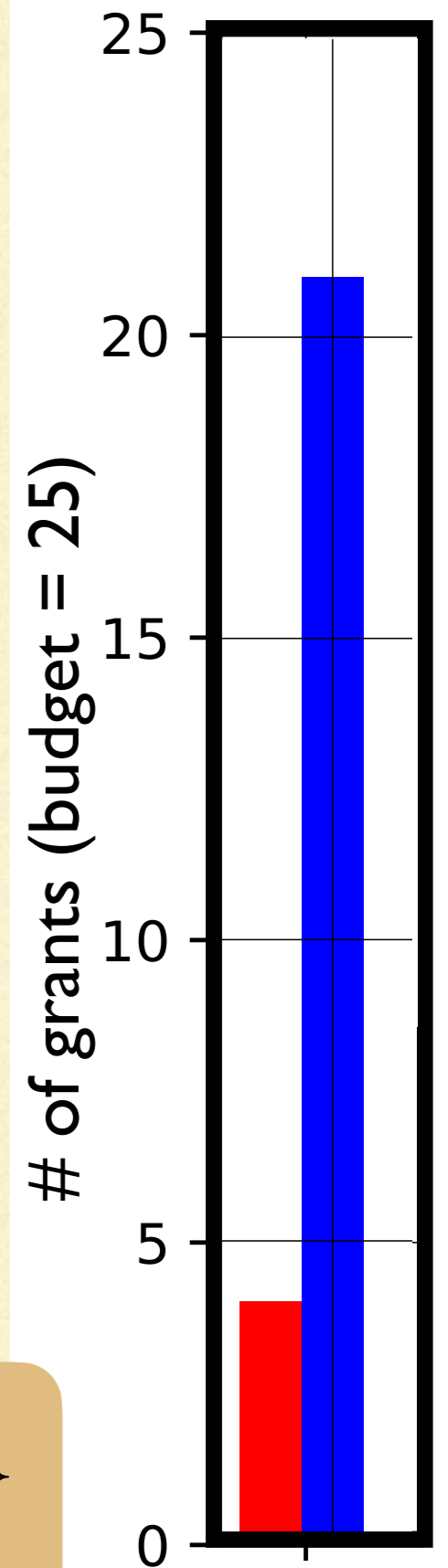
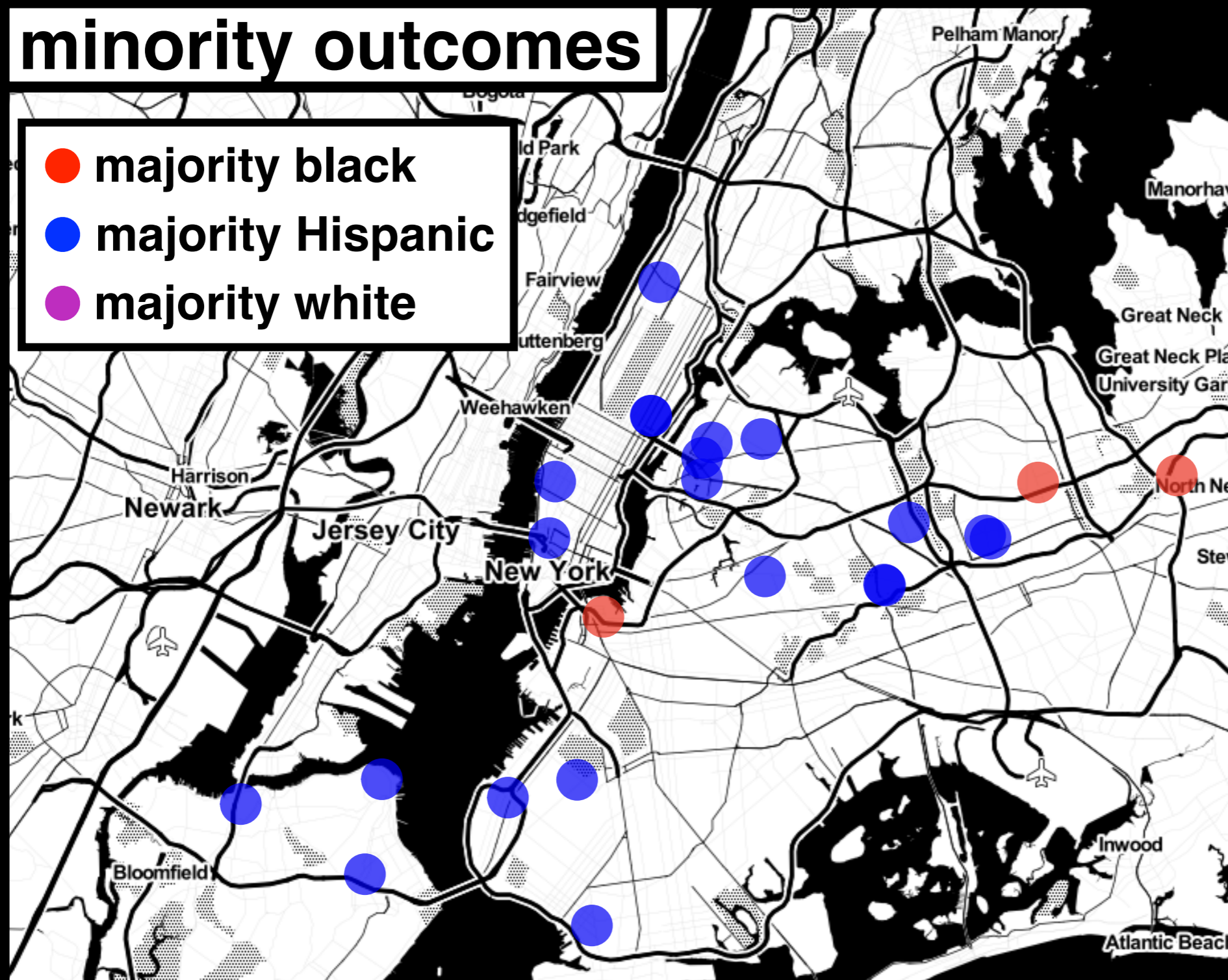
- majority black
- majority Hispanic
- majority white

$$\underbrace{\sum_{i=1}^{n_{a'}} z^{(i)}}_{c_{a'}} \leq \left\lfloor \frac{b}{|\mathcal{A}|} \right\rfloor \quad \forall a' \in \mathcal{A}$$



minority outcomes

- majority black
- majority Hispanic
- majority white



$$\underbrace{z^{(i)}}_{c_{ia'}} \leq 0 \quad \forall a' \in \mathcal{A}_{\text{maj}} \subseteq \mathcal{A}, i \in \{1, \dots, n\}$$

bank 1



bank 2

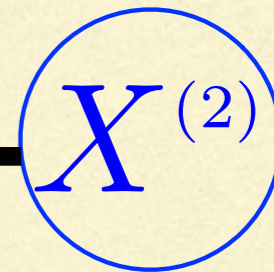
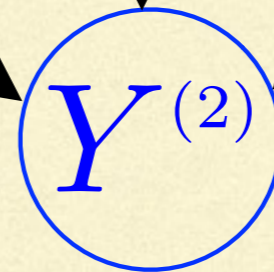
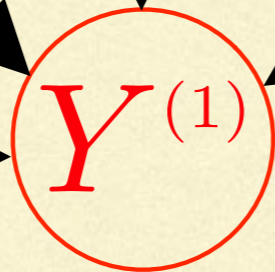
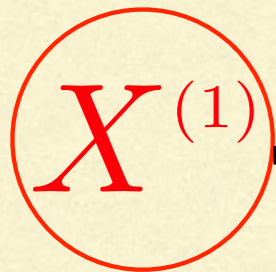
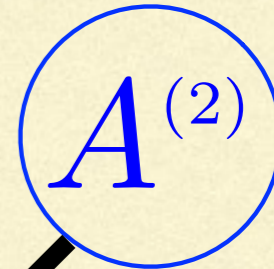
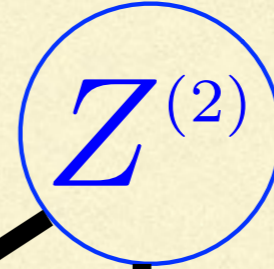


race
distribution

intervention:
audit

intervention:
audit

race
distribution

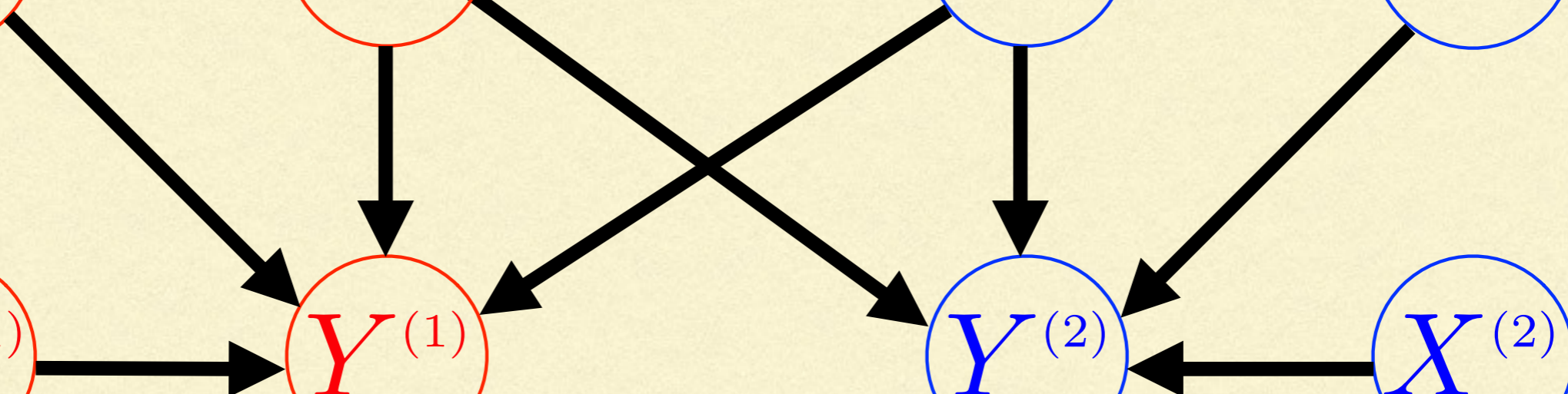


capital

% successful loans

% successful loans

capital



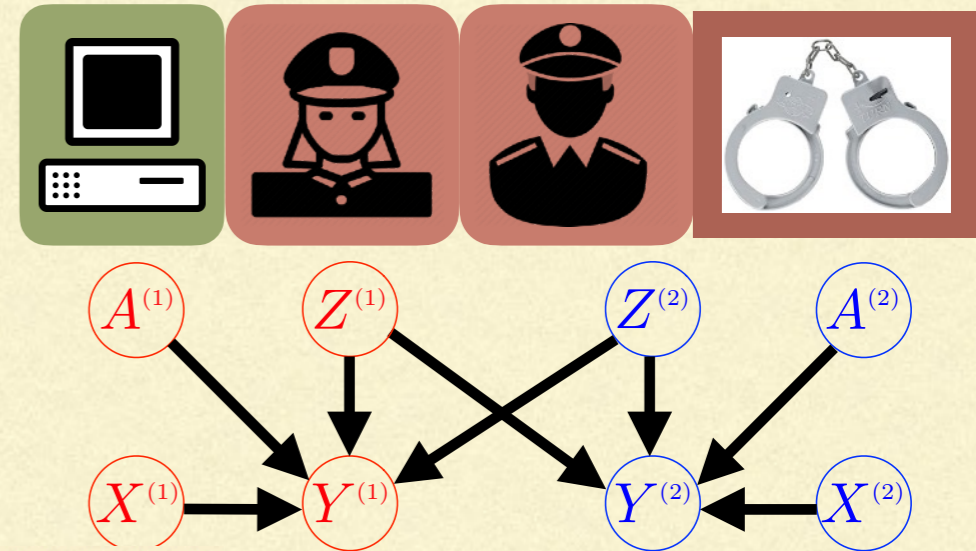
TAKE-AWAYS

- Many cases where ML algorithms decide only **part of an impact**



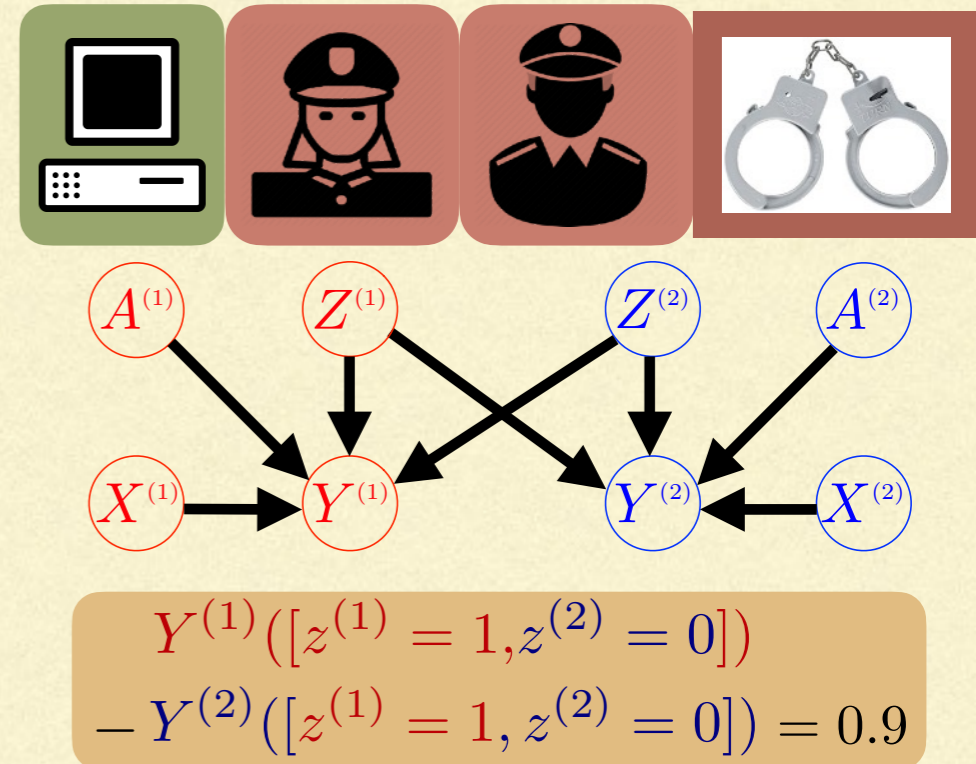
TAKE-AWAYS

- Many cases where ML algorithms decide only **part of an impact**
- Idea: formalize algorithmic decisions within society **using causal models**



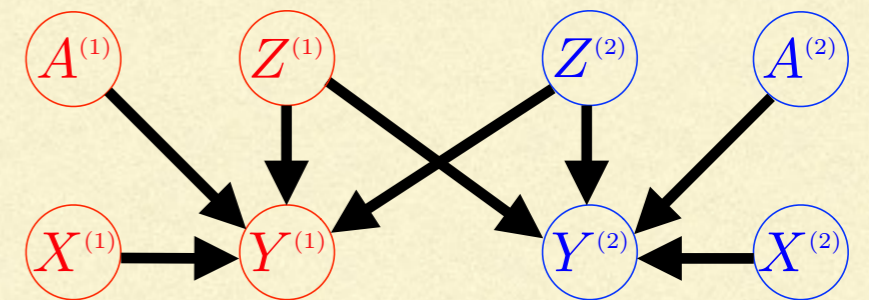
TAKE-AWAYS

- Many cases where ML algorithms decide only **part of an impact**
- Idea: formalize algorithmic decisions within society **using causal models**
- **Counterfactuals** allow us to formulate **discriminatory privilege**



TAKE-AWAYS

- Many cases where ML algorithms decide only **part of an impact**
- Idea: formalize algorithmic decisions within society **using causal models**
- **Counterfactuals** allow us to formulate **discriminatory privilege**
- We propose a constrained optimization problem that maximizes overall impact while **reducing privileged impact**

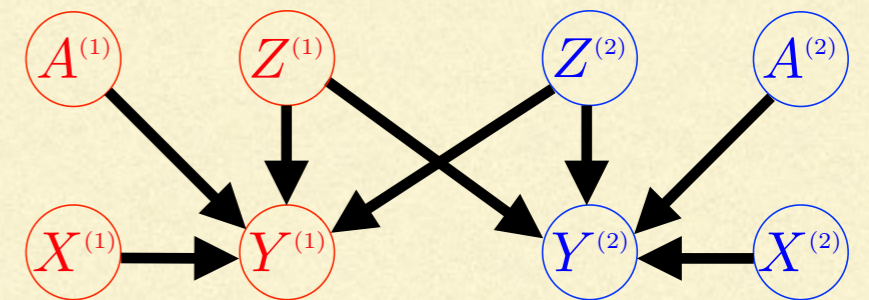


$$Y^{(1)}([z^{(1)} = 1, z^{(2)} = 0]) \\ - Y^{(2)}([z^{(1)} = 1, z^{(2)} = 0]) = 0.9$$

$$\max_{\mathbf{z} \in \{0,1\}^n} \sum_{i=1}^n \mathbb{E}[Y^{(i)}(a^{(i)}, \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] \\ s.t., \sum_{i=1}^n z^{(i)} \leq b \\ c_{ia'} \leq \tau \quad \forall a' \in \mathcal{A}, i \in \{1, \dots, n\},$$

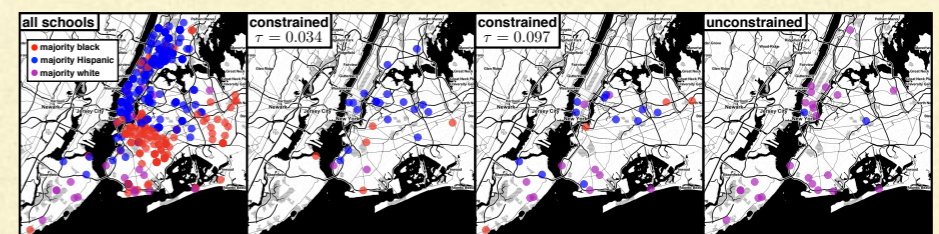
TAKE-AWAYS

- Many cases where ML algorithms decide only **part of an impact**
- Idea: formalize algorithmic decisions within society **using causal models**
- **Counterfactuals** allow us to formulate **discriminatory privilege**
- We propose a constrained optimization problem that maximizes overall impact while **reducing privileged impact**
- Allows one to make less discriminatory policy decisions for school funding



$$Y^{(1)}([z^{(1)} = 1, z^{(2)} = 0]) \\ - Y^{(2)}([z^{(1)} = 1, z^{(2)} = 0]) = 0.9$$

$$\max_{z \in \{0,1\}^n} \sum_{i=1}^n \mathbb{E}[Y^{(i)}(a^{(i)}, \mathbf{z}) \mid A^{(i)} = a^{(i)}, X^{(i)} = \mathbf{x}^{(i)}] \\ s.t., \sum_{i=1}^n z^{(i)} \leq b \\ c_{ia'} \leq \tau \quad \forall a' \in \mathcal{A}, i \in \{1, \dots, n\},$$



MY COAUTHORS

Chris Russell



Joshua R Loftus



Ricardo Silva



**The
Alan Turing
Institute**



**The
Alan Turing
Institute**