

# Making Decisions that Reduce Discriminatory Impact

Matt J. Kusner<sup>1,2</sup>

Chris Russell<sup>1,3</sup>

**Joshua R. Loftus**<sup>4</sup>

Ricardo Silva<sup>1,5</sup>

<sup>1</sup>Alan Turing Institute, <sup>2</sup>Oxford, <sup>3</sup>Surrey, <sup>4</sup>NYU, <sup>5</sup>UCL

6/13/2019

## HIGH-CRIME AREA MONITORING

### PROBLEM

- Parking lot thefts
- Frequent drug trafficking
- Looting threats
- Dark alleys
- Multiple locations
- Limited force
- Suspect identification
- Evidence gathering

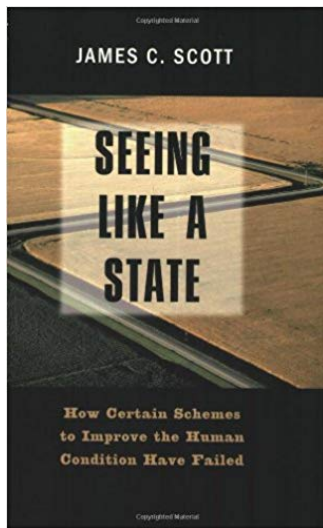
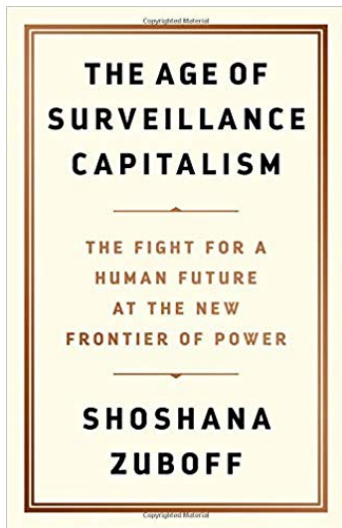
### SOLUTION

- Visible deterrent
- Man or unmanned
- Towable portability & rapid deployment
- Single officer
- Generator power
- Thermal & visible surveillance cameras and recording
- Extra spotlights



Source: flir.com “SkyWatch”

Maybe closer to the opposite of fair by default. . .



- ▶ Propose to formalize the **impact problem**
- ▶ Design fair(er) interventions under **causal interference**

### Defining impact

An **impact** is an event caused jointly by the decisions under our control and other real world factors. *Decisions about one individual can impact another individual.*

See also Liu et al. (ICML 2018), Green & Chen (FAT\* 2019)

*Fair predictions/decisions do not imply fair impacts, since other downstream factors can make the impact unfair (possibly to different individuals than the subjects of the original prediction/decision)*



- ▶ Propose to formalize the **impact problem**
- ▶ Design fair(er) interventions under **causal interference**

### Defining impact

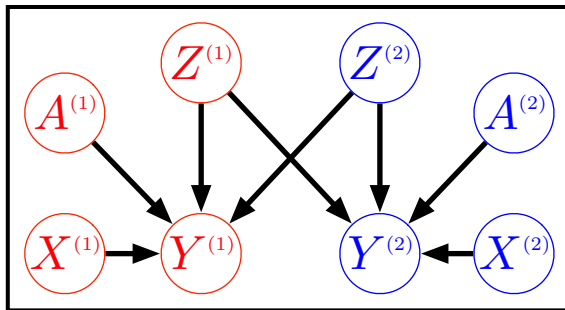
An **impact** is an event caused jointly by the decisions under our control and other real world factors. *Decisions about one individual can impact another individual.*

See also Liu et al. (ICML 2018), Green & Chen (FAT\* 2019)

*Fair predictions/decisions do not imply fair impacts*, since other downstream factors can make the impact unfair (possibly to different individuals than the subjects of the original prediction/decision)

## Causal interference: decisions affect multiple individuals

We use the structural causal model (SCM) framework



$Z$  is the intervention or policy we want to optimize,  $A$  the protected attribute,  $X$  other predictors, and  $Y$  the outcome (higher values are desirable), superscript for observation index

## School example

- ▶ Budget to pay for calculus classes in highschools (that do not already have them)
- ▶ Intervention:  $\mathbf{Z}^{(i)} = 1$  if school  $i$  receives funding for a class and 0 otherwise
- ▶ Outcome:  $\mathbf{Y}^{(i)}$  percent of students at school  $i$  taking the SAT (planning to go to college)
- ▶ Protected attribute:  $\mathbf{A}^{(i)}$  encodes whether school  $i$  is majority black, Hispanic, or white
- ▶ Interference: students at school  $i$  *may be able to take a calculus class at nearby schools*

Given causal model and data, design the best **fair** intervention  $\mathbf{Z}$

## Fair? What does that mean?

Predictions or decisions should be the same in the actual world and in a counterfactual world where the value of the protected attribute had been different

- ▶ Changing  $a$  to  $a'$  also changes descendents of  $\mathbf{A}$  in the SCM graph (*model-based counterfactuals*)
- ▶ **Counterfactual fairness** (Kusner et al, NeuRIPs 2017) is the property of invariance to those specific changes
- ▶ In this paper we instead bound **counterfactual privilege**

$$\mathbb{E}[\hat{Y}(a, \mathbf{Z})] - \mathbb{E}[\hat{Y}(a', \mathbf{Z})] < \tau$$

- ▶ In practice these *asymmetric constraints* will only be active for privileged values of  $a$  (actual, left term), and inactive otherwise

## Fair? What does that mean?

Predictions or decisions should be the same in the actual world and in a counterfactual world where the value of the protected attribute had been different

- ▶ Changing  $a$  to  $a'$  also changes descendants of  $\mathbf{A}$  in the SCM graph (*model-based counterfactuals*)
- ▶ **Counterfactual fairness** (Kusner et al, NeuRIPs 2017) is the property of invariance to those specific changes
- ▶ In this paper we instead bound **counterfactual privilege**

$$\mathbb{E}[\widehat{\mathbf{Y}}(a, \mathbf{Z})] - \mathbb{E}[\widehat{\mathbf{Y}}(a', \mathbf{Z})] < \tau$$

- ▶ In practice these *asymmetric constraints* will only be active for privileged values of  $a$  (actual, left term), and inactive otherwise

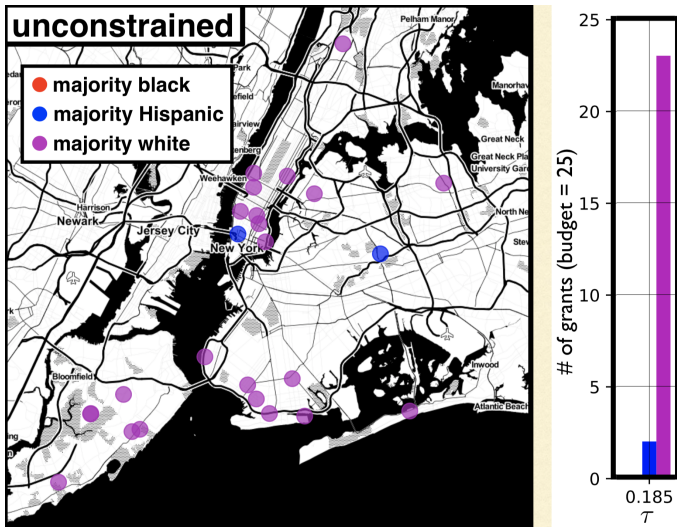
## Optimal intervention under interference

- ▶ Our goal is to design optimal interventions or policies  $\mathbf{Z}$  subject to a budget constraint, e.g.

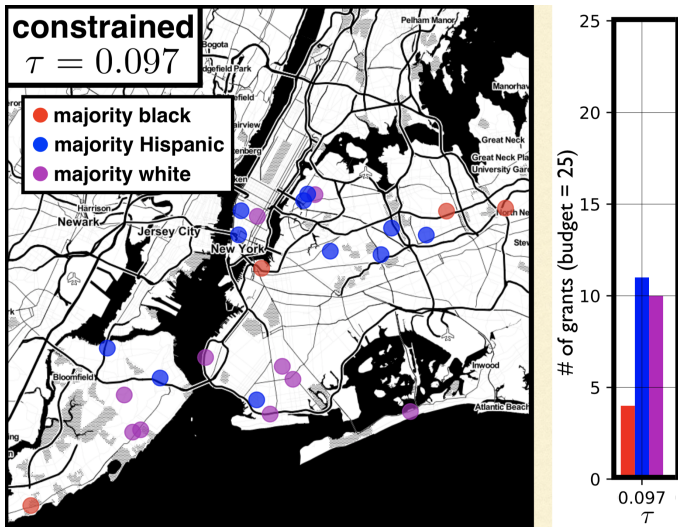
$$\mathbf{Z} = \arg \max \sum_i \mathbb{E} \left[ \widehat{\mathbf{Y}}^{(i)}(a^{(i)}, \mathbf{Z}) | \mathbf{A}^{(i)}, \mathbf{X}^{(i)} \right] \quad s.t. \quad \sum_i \mathbf{Z}^{(i)} \leq b$$

- ▶ Interference means  $\mathbf{Y}^{(i)}$  is potentially a function of all of  $\mathbf{Z}$  and not just  $\mathbf{Z}^{(i)}$
- ▶ Next two slides: optimal interventions with and without counterfactual privilege constraint

# School resource allocation without fairness constraint



# School resource allocation bounded counterfactual privilege





Thanks for your attention! See paper/poster(138) for more details

Matt Kusner



Chris Russell



Ricardo Silva

