

# **Lecture 20: Causal article on risk assessment**

**Criminology 250**

**Prof Maria Cuellar**

**University of Pennsylvania**

# Motivation

- **Motivation for class:** Review an article that uses causal inference to study the use of an algorithm for sentencing.

## Experimental Evaluation of Algorithm-Assisted Human Decision-Making: Application to Pretrial Public Safety Assessment\*

Kosuke Imai<sup>†</sup>   Zhichao Jiang<sup>‡</sup>   D. James Greiner<sup>§</sup>   Ryan Halen<sup>¶</sup>   Sooahn Shin<sup>||</sup>

First Draft: July 9, 2020  
This Draft: October 2, 2021

### Abstract

Despite an increasing reliance on fully-automated algorithmic decision-making in our day-to-day lives, human beings still make highly consequential decisions. As frequently seen in business, healthcare, and public policy, recommendations produced by algorithms are provided to human decision-makers to guide their decisions. While there exists a fast-growing literature evaluating

# Motivation

- **Algorithms appeal:** A growing body of literature has suggested the potential superiority of algorithmic decision-making over purely human choices across a variety of tasks.
- **Human appeal:** Many are resistant to have algorithms take over decision-making tasks, especially those that are high-stakes.
- **Hybrid systems appeal:** The desire for a human decision-maker as well as the precision and efficiency of algorithms have led to the adoption of hybrid systems involving both. By far the most popular system uses algorithmic recommendations to inform human decision-making. This has affected us in medicine, hiring, credit lending, investment decisions, and online shopping.
- **Risk assessment in criminal justice:** Of particular interest, algorithmic recommendations are increasingly used in risk assessment in the criminal justice system, with the goal of improving incarceration rulings and other decisions made by judges.

# Contribution of this paper

- **Bias/fairness:** There is a fast-growing literature in computer science/statistics/ML that studies the bias and fairness of algorithms.
- **An overlooked question is:** Whether those algorithms actually help humans make better decisions. In this paper, the authors develop a methodological framework for experimentally evaluating the impacts of algorithmic recommendations on human decision-making.
- **Method:** They conducted the first-ever real-world field experiment by providing, for a randomly selected set of cases, information from a system consisting of Public Safety Assessment (PSA) risk scores, and a recommendation from a Decision Making Framework (DMF) to a judge who makes an initial release decision.

# Outcomes

- **What does it actually test?:** They evaluate whether the PSA-DMF system (\$A\$) helps judges achieve their goal of preventing arrestees from committing a new crime or failing to appear in court while avoiding an unnecessarily harsh decision (\$Y\$).

# Algorithms in court

To date, algorithmic outputs have appeared in

- i) At the "first appearance" or arraignment hearing, during which a judge decides whether to release an arrestee pending disposition of any criminal charges, (in response to an arraignment the accused is expected to enter a plea) and
- ii) At sentencing, in which the judge imposes a punishment on a defendant found guilty.
- This paper deals with point (i).

# Setting

- **What happens at hearing?** We describe a typical first appearance hearing. The key decision the judge must make at a first appearance hearing is whether to release the arrestee pending disposition of any criminal charges, and if the arrestee is to be released, what conditions to impose.
- **Seek to avoid incarceration:** Because arrestees have not yet been adjudicated guilty of any charge at the time of a pretrial hearing, there exists a consensus that pretrial incarceration is to be avoided unless the risks associated with release are sufficiently high.

Judges deciding whether to release arrestees ordinarily consider risk factors among a variety of other concerns:

- 1) **FTA:** Failure to appear - the risk that the arrestee will fail to appear (FTA) at subsequent court dates.
- 2) **NCA:** New criminal activity - the risk that the arrestee will engage in new criminal activity (NCA) before the case is resolved.
- 3) **NVCA:** New violent criminal activity.

# Risk assessment algorithm

- Goal of Public Safety Assessment (PSA) algorithm is to classify arrestees according to FTA and NCA risks.
- They are generally constructed by fitting a statistical model to a training dataset based on past first appearance hearings and the subsequent incidences (or lack thereof) of FTA and NCA.
- The hope is that providing such instruments will improve the assessment of FTA and NCA risks and thereby lead to better decisions.
- The goal of this paper is to develop a general methodological framework for evaluating the impact of providing the PSA to judges at first appearance hearings using an RCT.



# Experiment

- Dane County, Wisconsin.
- PSA scores are based on the weighted indices of nine factors drawn from criminal history information, primarily prior convictions and FTA, and a single demographic factor, age.
- Notably, gender and race are not used to compute the PSA.
- The weights are calculated using past data.
- More information here: <https://advancingpretrial.org/psa/factors/>

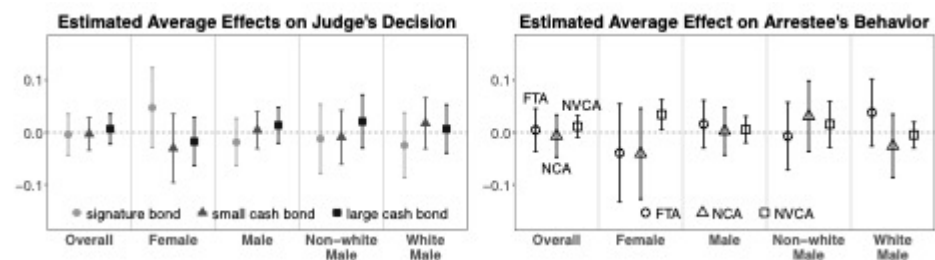
# Data

- 30 month treatment assignment period (2017-2019), followed by the collection of data on FTA, NCA, NCVA and other outcomes
- N = 1891 cases, 40% are white male arrestees and 13% are white female arrestees (non-white male is 39% and non-white female is 8%).

	<i>no</i> PSA (Control Group)			PSA (Treatment Group)			Total (%)
	Signature bond	Cash bond ≤\$1000	Cash bond >\$1000	Signature bond	Cash bond ≤\$1000	Cash bond >\$1000	
Non-white Female	64 (3.4)	11 (0.6)	6 (0.3)	67 (3.5)	6 (0.3)	0 (0.0)	154 (8.1)
White Female	91 (4.8)	17 (0.9)	7 (0.4)	104 (5.5)	17 (0.9)	10 (0.5)	246 (13.0)
Non-white Male	261 (13.8)	56 (3.0)	49 (2.6)	258 (13.6)	53 (2.8)	57 (3.0)	734 (38.8)
White Male	289 (15.3)	48 (2.5)	44 (2.3)	276 (14.6)	54 (2.9)	46 (2.4)	757 (40.0)
FTA committed	218 (11.5)	42 (2.2)	16 (0.8)	221 (11.7)	45 (2.4)	16 (0.8)	558 (29.4)
not committed	487 (25.8)	90 (4.8)	90 (4.8)	484 (25.6)	85 (4.5)	97 (5.1)	1333 (70.6)
NCA committed	211 (11.2)	39 (2.1)	14 (0.7)	202 (10.7)	40 (2.1)	17 (0.9)	523 (27.7)
not committed	494 (26.1)	93 (4.9)	92 (4.9)	503 (26.6)	90 (4.8)	96 (5.1)	1368 (72.4)
NCVA committed	36 (1.9)	10 (0.5)	3 (0.2)	44 (2.3)	10 (0.5)	6 (0.3)	109 (5.7)
not committed	669 (35.4)	122 (6.5)	103 (5.4)	661 (35.0)	120 (6.3)	107 (5.7)	1782 (94.3)
Total	705 (37.3)	132 (7.0)	106 (5.6)	705 (37.3)	130 (6.9)	113 (6.0)	1891 (100)

# Findings

- In general, we observe a positive association between the PSA scores and judge's decisions, implying that a higher PSA score is associated with a harsher decision.
- This figure presents the estimated average causal effect of PSA provisions on the judge's decisions (left plot) and three outcomes of interest (right plot).



- The vertical bars are the 95% confidence intervals.
- PSA provision appears to have little overall effect on the judge's decision and arrestee's behavior, on average, although it may slightly increase NVCA among female arrestees.
- **Results:** The results imply that PSA provision, on average, has little effect on the judge's decisions.

# Causal analysis about subgroups

$Z_i$ : Binary treatment variable indicating whether the PSA is presented to the judge of case  $i = 1, 2, \dots, n$ .

$D_i$ : Binary detention decision made by the judge to either detain ( $D_i=1$ ) or release ( $D_i=0$ ) the arrestee prior to the trial.  $D_i(z)$  is potential outcome under treatment  $z$ .

$Y_i$ : Binary outcome - they code all their outcomes (NCA, NVCA, FTA) as binary variables. For example,  $Y_i = 1$  implies that the arrestee of case  $i$  commits an NCA.

$\mathbf{X}_i$ : A vector of observed pre-treatment covariates for case  $i$ . They include age, gender, race, and prior criminal history.

Preventable cases:  $APCE_p = E(D_i(1) - D_i(0) | Y_i(1) = 0, Y_i(0) = 1)$

Risky cases:  $APCE_r = E(D_i(1) - D_i(0) | Y_i(1) = 1, Y_i(0) = 1)$

Safe cases:  $APCE_s = E(D_i(1) - D_i(0) | Y_i(1) = 0, Y_i(0) = 0)$

# Assumptions

- Randomization of the treatment assignment
- Exclusion restriction.
- Monotonicity.
- They also use nonparametric estimation.

# Paper conclusions

- First, we find that PSA provision has little overall impact on the judge's decisions.
- Second, we find potentially suggestive evidence PSA provision may encourage the judge to make more lenient decisions for female arrestees regardless of their risk levels while leading to more stringent decisions for males who are classified as risky.
- Third, PSA provision appears to widen the existing gender difference of the judge's decision against male arrestees whereas it does not seem to alter decision-making across race among male arrestees.
- Our results suggest that the PSA's recommendations may be harsher than necessary.