

Triage Score: A Counterfactual Risk Assessment Instrument

Kosuke Imai

Harvard University

Joint work with Sooahn Shin, D. James Greiner, and Ryan Halen

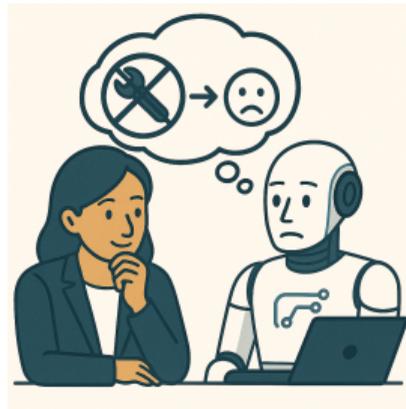
Bridging Prediction and Intervention Problems in Social Systems

Simons Institute for the Theory of Computing

January 12, 2026

Algorithm-assisted Human Decision Making

- AI/machine learning based recommendations
 - Routine decisions made by individuals in daily lives
 - Consequential decisions made by clinicians, judges, etc.
- Risk assessment instruments (“risk scores”)
 - Predicts the likelihood of an undesired outcome if *no intervention* is made
 - High risk → recommend intervention
 - Example: medicine, criminal justice system
- Methodological literature
 - Coston et al. (2020); Rambachan et al. (2022); Ben-Michael et al. (2025)
 - selective labels problem



Background on Triage Scores

- Effective decision-making requires consideration of how decisions affect outcomes
 - risk scores: only consider likelihood of an undesired outcome under no intervention
 - **triage scores**: incorporate counterfactual outcomes under **alternative decisions**
- **Triage**: *trier* or “sort” in French



- Napoleon’s military surgeon: Baron Dominique-Jean Larrey
 - prioritize worst wounds rather than first-come first-serve
 - do not fully consider counterfactual outcomes
-
- Modern triage systems: prioritize those **who would benefit most** to maximize survival given resource constraints

Triage Scores

- **Statistical decision-theoretic framework**
 - Standard decision theory: specify utility for each decision and its consequence
 - **Counterfactual** decision theory: incorporate counterfactual outcomes under alternative decisions when specifying each utility
 - Nuanced notions: regret, difficulty, do-no-harm, fairness etc. (e.g., Imai and Jiang 2023; Mueller and Pearl 2023; Ben-Michael et al. 2024; Christy and Kowalski 2024; Koch and Imai 2025)
- Two axes of comparisons
 1. Risk v.s. Triage scores
 2. Standard v.s. Counterfactual decision theory
- Evaluation of decision making systems
 - Do AI recommendations help humans make better decisions?
 - When should human decision makers ignore or follow AI recommendations?
 - What is the optimal decision rule that maximizes policy maker's expected utility?

Empirical Application: Public Safety Assessment (PSA)

At the first appearance hearing, judges must decide

1. **whether to release** an arrestee pending disposition of charges
2. what conditions (e.g., bail and monitoring) to impose if released

- **Public Safety Assessment (PSA)** provides algorithmic recommendations to judges
 - Inputs: age, criminal history
 - No gender or race is used
 - Outputs: Risk factors \rightsquigarrow recommendation
 - Widely used system across 25 US states
- PSA is a risk score: prediction under release
 - (1) **failure to appear (FTA)**
 - (2) **new criminal activity (NCA)**
 - (3) **new violent criminal activity (NVCA)**

DEVELOPING A NATIONAL MODEL FOR PRETRIAL RISK ASSESSMENT

Every day in America, judges have to answer a critical question again and again: What are the chances that a recently arrested defendant, **if released before trial**, will commit a new crime, a new violent crime, or fail to appear for court?

Does the dataset include instances of defendants who were detained? If so, does the data include outcomes for those people (i.e., did the data account for counterfactual estimation; if so, how)? The approximate 750,000 cases that were considered in the model development process were **only** defendants who had been released at some point in the pretrial process.

Sources: Arnold Ventures report; Risk Assessment Factsheet (Stanford Law School)



Public Safety Assessment Report

(Date Created: 08/04/2018 03:36:34)

SID:	[REDACTED]	Name:	[REDACTED]	Gender:	M
PC ID:	[REDACTED]	Arrest Date:	08/04/2018	Birth Date:	[REDACTED]

New Violent Criminal Activity Flag: No

New Criminal Activity Scale

1	2	3	4	5	6
---	---	---	---	---	---

Failure to Appear Scale

1	2	3	4	5	6
---	---	---	---	---	---

Charge(s):

58-37A-5(1)(A)	USE OR POSSESSION OF DRUG PARAPHERNALIA
58-37-8(2)(H)(III)	CONTROLLED SUBSTANCE SCHEDULE III, IV, V
76-6-602	RETAIL THEFT (SHOPLIFTING)

Risk Factors:

1. Age at Current Arrest	25
2. Current Violent Offense	No
a. Current Violent Offense & 20 Years Old or Younger	No
3. Pending Charge at the Time of the Offense	No
4. Prior Misdemeanor Conviction	Yes
5. Prior Felony Conviction	Yes
a. Prior Conviction	Yes
6. Prior Violent Offense	4
7. Prior Failure to Appear in Past 2 Years	7
8. Prior Failure to Appear Older Than 2 Years	Yes
9. Prior Sentence to Incarceration	Yes

Recommendations:

1. Release with PRL4 Conditions - Maximum conditions where available: PRL3 conditions, electronic monitoring, home detention, and/or drug testing

Does a Judge Follow PSA?

- We conducted an RCT in Utah (We made the Wisconsin data publicly available)
- Experimental samples: 7K arrest cases
- **Provision** of PSA is **randomized** across cases
- **PSA provision** influenced the judge's decision (more **agreement** between judges and PSA)
 - cash bail
 - release on their own recognizance (ROR)
- Next step: evaluate and improve decision-making systems
 - Does PSA help the judge make better decisions?
 - What is the optimal decision under a given utility function?

Human

		PSA	
		ROR	Cash bail
ROR		13.3% (609)	4.9 (224)
Cash bail		53.3 (2439)	28.4 (1301)

Human+PSA

		PSA	
		ROR	Cash bail
ROR		18.3% (807)	5.4 (240)
Cash bail		46.4 (2046)	29.9 (1320)

Statistical Decision-Theoretic Framework

- Binary decision, binary outcome (generalizable to categorical cases)
- $D^* \in \{0, 1\}$: **decision** made by a decision-maker (e.g., human, human+PSA)
 - $D^* = 1$ if judge imposes cash bail
 - $D^* = 0$ if release on own recognizance
- $Y \in \{0, 1\}$: indicator of an undesirable **outcome**
 - $Y(d)$: potential outcome under decision $D^* = d$
 - $Y = 1$ if rearrested with an NCA, and $Y = 0$ otherwise
- **Principal strata:** $(Y(0), Y(1)) = (y_0, y_1)$

$$(Y(0), Y(1)) = \begin{cases} (0, 0) & \text{"Safe": no new crime under either decision} \\ (0, 1) & \text{"Backlash": new crime only under cash bail} \\ (1, 0) & \text{"Preventable": new crime only if released} \\ (1, 1) & \text{"Risky": new crime regardless of decision} \end{cases}$$

- **Risk score** considers $Y(0)$: two latent groups {safe, backlash} and {preventable, risky}
- **Triage score** considers $(Y(0), Y(1))$: four latent groups

Evaluating Decision-Making Systems with Triage Score

		Decision	
		Release ($D^* = 0$)	Cash bail ($D^* = 1$)
Principal Strata	Safe ($Y(0) = 0, Y(1) = 0$)	$u_{\neg\text{crime}}^{\text{ror}} + \tilde{u}_{\neg\text{crime}}^{\text{cash}}$	$u_{\neg\text{crime}}^{\text{cash}} + \tilde{u}_{\neg\text{crime}}^{\text{ror}}$
	Backlash ($Y(0) = 0, Y(1) = 1$)	$u_{\neg\text{crime}}^{\text{ror}} + \tilde{u}_{\text{crime}}^{\text{cash}}$	$u_{\text{crime}}^{\text{cash}} + \tilde{u}_{\neg\text{crime}}^{\text{ror}}$
	Preventable ($Y(0) = 1, Y(1) = 0$)	$u_{\text{crime}}^{\text{ror}} + \tilde{u}_{\neg\text{crime}}^{\text{cash}}$	$u_{\neg\text{crime}}^{\text{cash}} + \tilde{u}_{\text{crime}}^{\text{ror}}$
	Risky ($Y(0) = 1, Y(1) = 1$)	$u_{\text{crime}}^{\text{ror}} + \tilde{u}_{\text{crime}}^{\text{cash}}$	$u_{\text{crime}}^{\text{cash}} + \tilde{u}_{\text{crime}}^{\text{ror}}$

- $u_{y_d}^d$: standard utility under decision d with outcome $Y(d) = y_d$
- $\tilde{u}_{y_{1-d}}^{1-d}$: counterfactual utility under alternative decision $1 - d$ with outcome $Y(1 - d) = y_{1-d}$

Evaluating Decision-Making Systems under Risk Score

Decision

Baseline Outcome

		Release ($D^* = 0$)	Cash bail ($D^* = 1$)
$Y(0) = 0$	Safe ($Y(0) = 0, Y(1) = 0$)	$u_{ror, \neg crime}$	$u_{cash, \neg crime}$
	Backlash ($Y(0) = 0, Y(1) = 1$)		
$Y(0) = 1$	Preventable ($Y(0) = 1, Y(1) = 0$)	$u_{ror, crime}$	$u_{cash, crime}$
	Risky ($Y(0) = 1, Y(1) = 1$)		

Decision

		Release ($D^* = 0$)	Cash bail ($D^* = 1$)
$Y(0) = 0$	Safe ($Y(0) = 0, Y(1) = 0$)	$u_{ror, \neg crime}$	$u_{cash, \neg crime}$
	Backlash ($Y(0) = 0, Y(1) = 1$)	$= u_{\neg crime}^{ror} + \tilde{u}_{\neg crime}^{cash}$	$= u_{\neg crime}^{cash} + \tilde{u}_{\neg crime}^{ror}$
$Y(0) = 1$	Preventable ($Y(0) = 1, Y(1) = 0$)	$= u_{crime}^{ror} + \tilde{u}_{crime}^{cash}$	$= u_{crime}^{cash} + \tilde{u}_{crime}^{ror}$
	Risky ($Y(0) = 1, Y(1) = 1$)	$= u_{crime}^{ror} + \tilde{u}_{crime}^{cash}$	$= u_{crime}^{cash} + \tilde{u}_{crime}^{ror}$

Nonparametric Identification of Expected Counterfactual Utility

- Expected counterfactual utility is generally unidentifiable
 - joint potential outcomes are not simultaneously observed

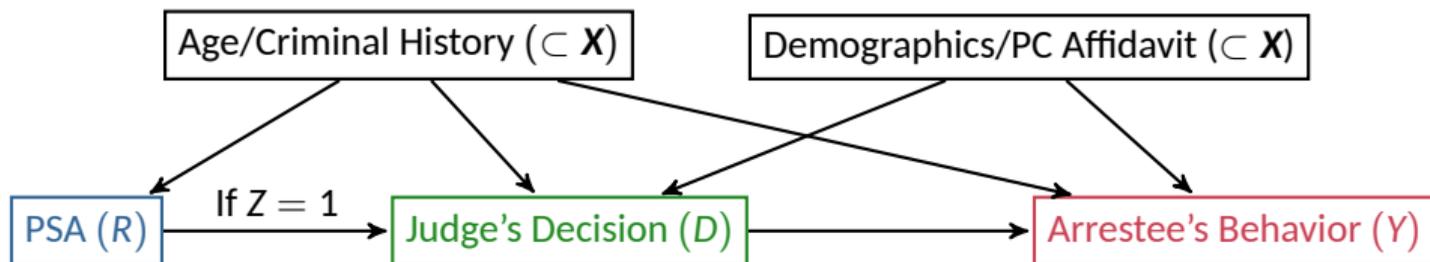
$$\mathbb{E}[U(u; D^*)] = \sum_{d=0}^{K_D-1} \sum_{y_0=0}^{K_Y-1} \dots \sum_{y_{K_D-1}=0}^{K_Y-1} u(d, \{y_d\}_{d=0}^{K_D-1}) \Pr(D^* = d, Y(0) = y_0, \dots, Y(K_D - 1) = y_{K_D-1})$$

- partial identification approaches (Li and Pearl 2019; Ben-Michael et al. 2024)
- **Additive counterfactual utility** (Koch and Imai, 2025):

$$u(d, \{y_d\}_{d=0}^{K_D-1}) = \underbrace{u_{y_d}^d}_{\text{standard utility}} + \sum_{d' \neq d} \underbrace{\tilde{u}_{y_{d'}}^{d'}}_{\text{counterfactual utility}} .$$

- standard utility as a special case
- can be a function of covariates \mathbf{X}
- Additivity as a **necessary and sufficient condition** for nonparametric identification under unconfoundedness

Assumptions in Our Empirical Application



- **Assumption 1: Single-blinded and unconfounded treatment assignment**

1. Single-blinded treatment: $Y_i(0, D_i(0)) = Y_i(1, D_i(1))$ if $D_i(0) = D_i(1)$
 \rightsquigarrow PSA can only affect the **outcome** through the **decision**
2. Unconfounded treatment assignment: $Z_i \perp\!\!\!\perp \{R_i, \{D_i(z), Y_i(d)\}_{z \in \{0,1\}, d \in \mathcal{D}}\} \mid \mathbf{X}_i$
3. Overlap: $e(\mathbf{x}) := \Pr(Z_i = 1 \mid \mathbf{X}_i = \mathbf{x}) \in [\lambda_Z, 1 - \lambda_Z]$ for $0 < \lambda_Z \leq \frac{1}{2}$

- **Assumption 2: Unconfounded decision and overlap**

$$\{Y_i(d)\}_{d \in \mathcal{D}} \perp\!\!\!\perp D_i \mid \mathbf{X}_i, Z_i, R_i, \quad \Pr(Y_i = y \mid D_i = d, Z_i = z, R_i = r, \mathbf{X}_i) \geq \lambda_D$$

- holds if we have all the information that the judge takes into account (*and we do!*)
- R can be confounded: $R \leftarrow U_D \rightarrow D$ and $R \leftarrow U_Y \rightarrow Y$ are allowed

Affidavit of Probable Cause and GenAI Powered Inference (GPI)

IN THE
[REDACTED] JUSTICE COURT
COUNTY OF [REDACTED] STATE OF UTAH

State of Utah vs. [REDACTED] Date of Birth: [REDACTED]	Arrestee	Affidavit of Probable Cause
---	----------	-----------------------------

On [REDACTED] 2020 [REDACTED] the defendant was arrested for the offense(s) of:

	Offense Date	Offense Description	Statute	Gov Code	Severity	DV
1	[REDACTED] 2020	INTOXICATION	76-9-701	UT	MC	No

I believe there is probable cause to charge the defendant with these charges because:

On [REDACTED] 2020 I responded to a public intoxication in the area of [REDACTED] Security in the area so a male later identified as [REDACTED] punching windows of local businesses and walking into the streets while yelling at pedestrians. When I arrived on scene I saw [REDACTED] stumbling down the sidewalk screaming into the sky and waving his arms around his body. I stopped out with [REDACTED] and he would not answer my questions without yelling incoherently and jumping around. A flashlight was used on [REDACTED] eyes and I noticed his pupils did not constrict and remained pinpoint. Based on his behavior and his eyes, I believed that [REDACTED] was under the influence of an unknown substance. I did not believe [REDACTED] could safely make his way to shelter and would be a danger to himself and possibly others out in the area with him. [REDACTED] was taken to [REDACTED] County Jail on a public Intoxication charge.

Officer Name: [REDACTED]	Badge ID: [REDACTED]
I am a sworn officer with: [REDACTED] Police Dept	
Arresting agency case number: [REDACTED]	Associated citation number: [REDACTED]

I declare under penalty of perjury and under the laws of the State of Utah that the foregoing is true and correct.
/s/ [REDACTED]

SUBMISSION IDENTIFICATION INFORMATION

Booking agency: [REDACTED]	Booking agency ORI: [REDACTED]	
Booking agency case number: [REDACTED]	SID: [REDACTED]	UTN: [REDACTED]
Booking UserID: [REDACTED]	Booking date/time: [REDACTED]	Submission ID: [REDACTED] (Version 1)

- In Utah, PC Affidavit (along with PSA) contains all the information a judge has
- structured information: types of offense
- unstructured information: probable cause statement

- GPI (Imai and Nakamura 2025)
 - use open-source LLM to regenerate texts
 - obtain their internal representation
 - guaranteed to contain all the information
 - leverage LLM's internal structure

Policy Evaluation and Learning

- **Nonparametric identification** of expected additive counterfactual utilities
 - We consider a general setting (Assumptions 1 and 2), which are plausible in our application
 - Under binary decision, Assumption 2 (unconfoundedness of decision) is not required for the comparison between human-alone and human+PSA
- **AIPW estimators** given the decision and outcome models
 - GPI to learn the “deconfounder” through the DragonNet neural network architecture
 - asymptotic normality, double-robustness, rate-robustness
- **Statistical hypothesis test**
 1. **When do we prefer Human+PSA?** Invert the hypothesis test

$$H_0 : \bar{U}(u; D(1)) \leq \bar{U}(u; D(0)) \text{ vs. } H_1 : \bar{U}(u; D(1)) > \bar{U}(u; D(0))$$

for a given utility function u

2. **What is an optimal decision rule?** Consider a covariate-dependent policy $\pi : \mathcal{X} \rightarrow \mathcal{D}$ and estimate optimal policy by solving the empirical utility maximization problem

Utility Specification

- Focus on “loss” aspects: cost and regret
 - Cost of decision: c^{cash}
 - Cost of outcome: $c_{\text{crime}}^{\text{ror}}$ and $c_{\text{crime}}^{\text{cash}}$
 - Regret from counterfactual outcome: $r_{\neg\text{crime}}^{\text{ror}}$ and $r_{\neg\text{crime}}^{\text{cash}}$

		Decision	
		Release ($D^* = 0$)	Cash bail ($D^* = 1$)
Principal Strata	Safe ($Y(0) = 0, Y(1) = 0$)	$-r_{\neg\text{crime}}^{\text{cash}}$	$-c^{\text{cash}} - r_{\neg\text{crime}}^{\text{ror}}$
	Backlash ($Y(0) = 0, Y(1) = 1$)	0	$-c^{\text{cash}} - c_{\text{crime}}^{\text{cash}} - r_{\neg\text{crime}}^{\text{ror}}$
	Preventable ($Y(0) = 1, Y(1) = 0$)	$-c_{\text{crime}}^{\text{ror}} - r_{\neg\text{crime}}^{\text{cash}}$	$-c^{\text{cash}}$
	Risky ($Y(0) = 1, Y(1) = 1$)	$-c_{\text{crime}}^{\text{ror}}$	$-c^{\text{cash}} - c_{\text{crime}}^{\text{cash}}$

Utility Specification (cont.)

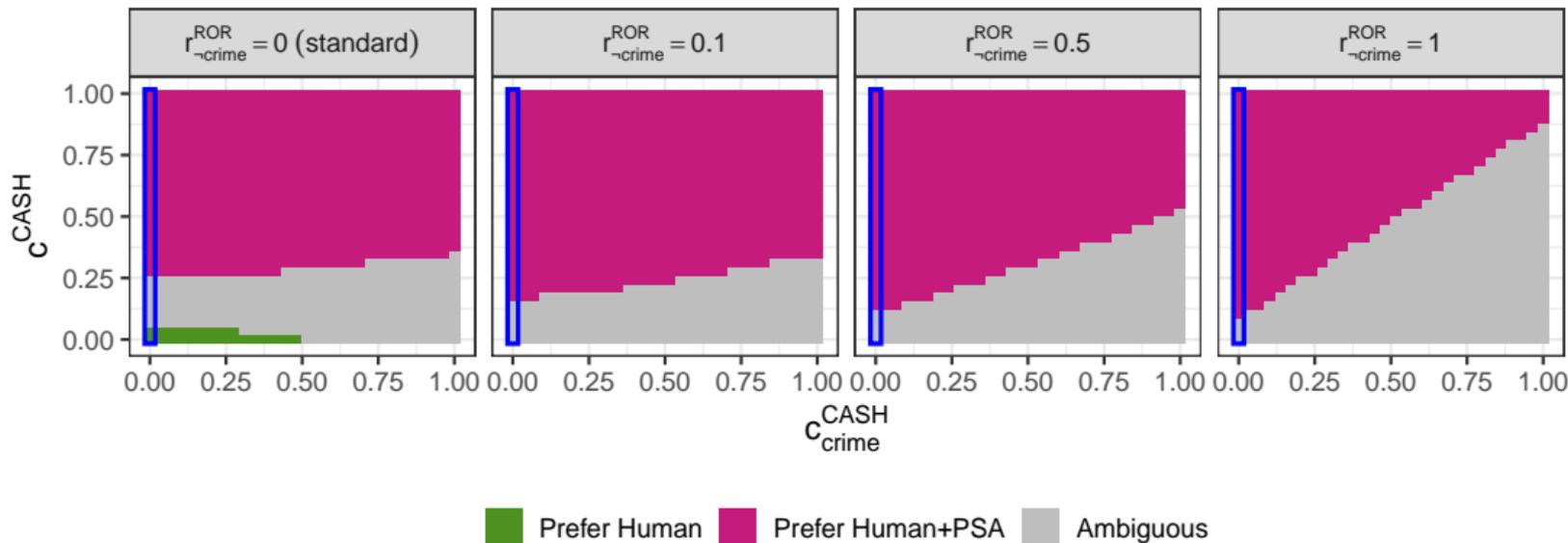
- We standardize $c_{\text{crime}}^{\text{ror}} = 1$ (cost of new crime under ROR)
 1. **Standard** decision theory with the **risk** score framework: c^{cash}
 2. **Standard** decision theory with the **triage** score framework: $c^{\text{cash}}, c_{\text{crime}}^{\text{cash}}$
 3. **Counterfactual** decision theory with the **risk** score framework: $c^{\text{cash}}, r_{\text{-crime}}^{\text{ror}}$
 4. **Counterfactual** decision theory with the **triage** score framework: $c^{\text{cash}}, c_{\text{crime}}^{\text{cash}}, r_{\text{-crime}}^{\text{ror}}, r_{\text{-crime}}^{\text{cash}}$
- To further reduce the number of parameters for our analysis:

$$\frac{\text{Cost of new crime under cash bail } (c_{\text{crime}}^{\text{cash}})}{\text{Cost of new crime under ROR } (c_{\text{crime}}^{\text{ror}})}$$

$$= \frac{\text{Regret of no new crime under cash bail } (r_{\text{-crime}}^{\text{cash}})}{\text{Regret of no new crime under ROR } (r_{\text{-crime}}^{\text{ror}})}$$

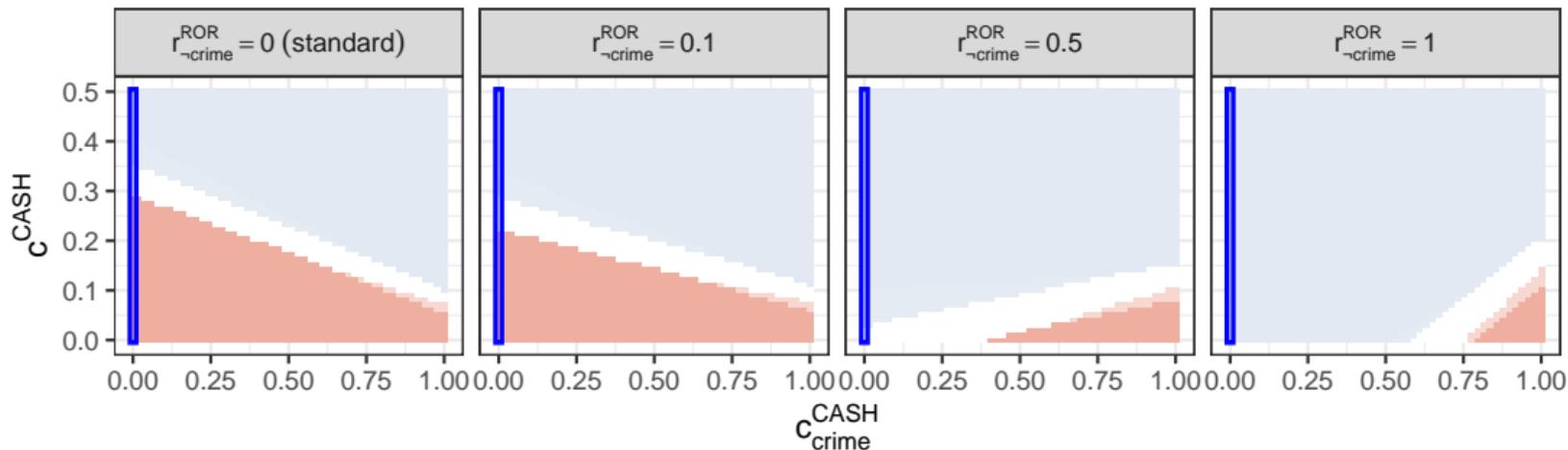
$$\rightsquigarrow r_{\text{-crime}}^{\text{cash}} = r_{\text{-crime}}^{\text{ror}} c_{\text{crime}}^{\text{cash}}$$

When do we prefer Human+PSA over Human? (Preliminary Results)



- Outcome: new criminal activity
- Blue boxes: risk score system (when $c_{crime}^{cash} = 0$)
- First panel: standard decision framework (when $r_{-crime}^{ror} = r_{crime}^{ror} = 0$)
- More ambiguous under triage score than risk score

Optimal Decision Rules (Preliminary Results)



Change in Cash Bail Proportion (relative to current) 31% 12% 0% -5% -6%

- Outcome: failure to appear
- Policy class: current PSA decision rule with varying NVCA flag threshold (θ)
- PSA recommends cash bail if $NVCA(\mathbf{X}) > \theta$ where $NVCA(\mathbf{X}) \in [0, 7]$
- triage scores lead to more lenient optimal decisions than risk scores

Conclusion

We propose **triage score** as a comprehensive framework for decision-making

- Unlike risk scores, triage scores incorporate counterfactual outcomes and utilities
- It enables decision-makers to incorporate a wide range of ethical and practical factors
- It can be used to assess **Human**, **Human+PSA**, and **PSA** decision-making systems
- We apply our framework to a RCT to evaluate **a pre-trial risk assessment instrument**
 - Triage scores capture richer utility structures than risk scores and yield distinct results
 - Our preliminary results show that triage score framework yields higher expected utility for Human + PSA and more lenient optimal recommendations than risk score framework
- Next steps:
 - RCT with the optimal triage score (in planning!)
 - Evaluate Human+Triage system