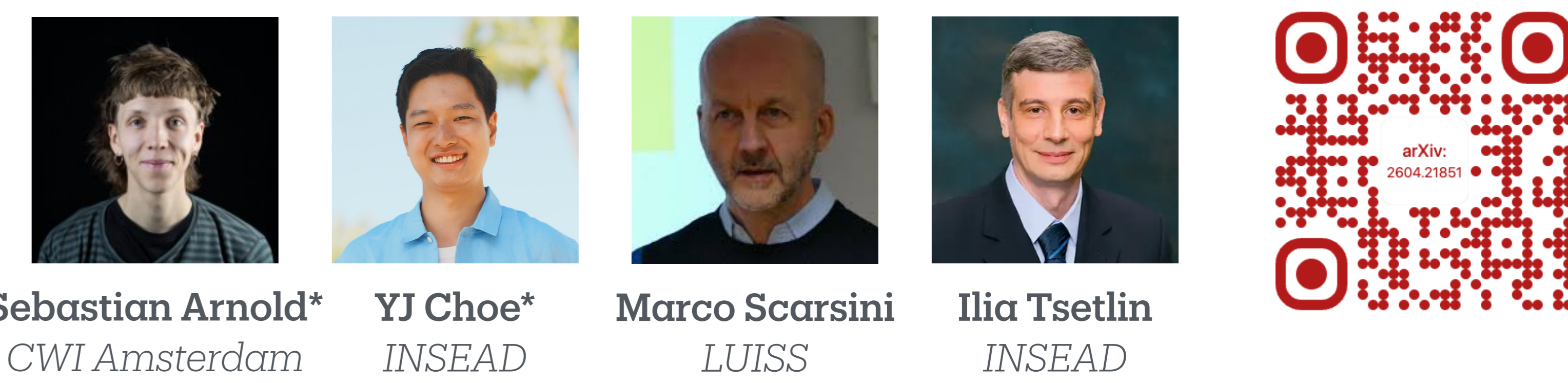


Betting on Bets:

Anytime-Valid Tests for Stochastic Dominance



Motivation

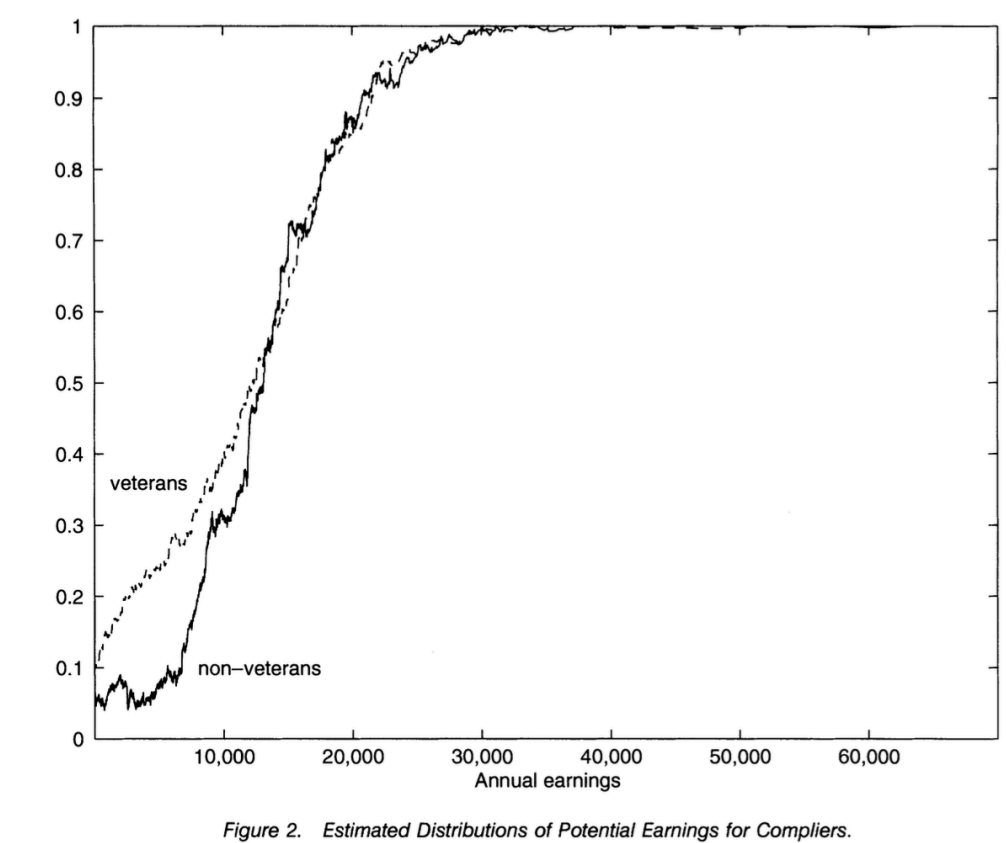
Does **Y** have an **upside** over **X**? Particularly when...

1. their means are not meaningfully different, or
2. "means" aren't even well-defined (e.g., ordinal data).

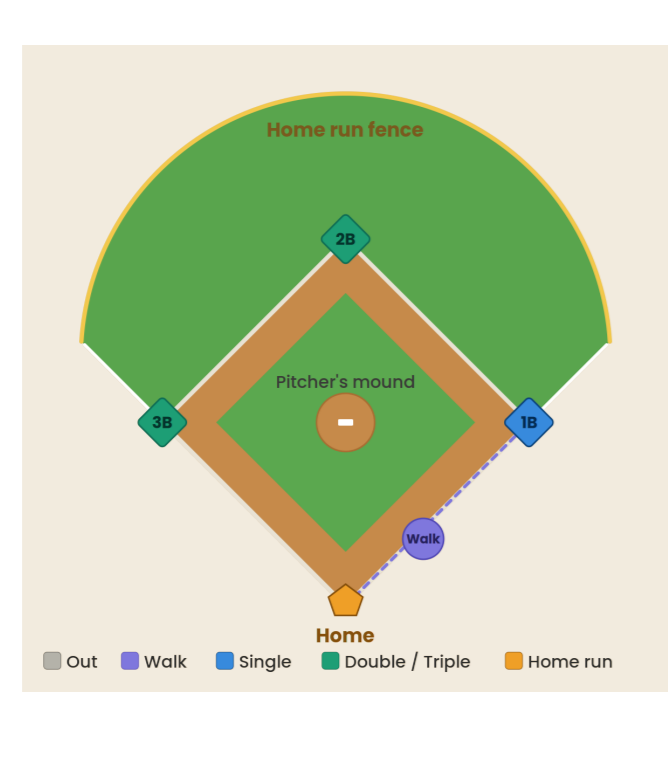
Can we test for upsides **in real time**? We want...

1. validity under continuous monitoring, and
2. comparable power to classical (e.g., bootstrap) tests.

Effects of veteran status on future incomes (Abadie, 2002)



"3rd time effect" in baseball



The Manager, the Ace and a Decision That Will Haunt the Rays

Blake Snell was dominating the Dodgers in a must-win World Series game. Kevin Cash still pulled him, raising painful questions of what might have been.



Building block: Growth-rate optimal e-values

Def. An **e-process** $(E_t)_{t \in \mathbb{N}}$ for a composite hypothesis H is a sequence of nonnegative variables s.t.

$$\mathbb{E}_{\mathbb{P}}[E_\tau] \leq 1 \text{ for every } \mathbb{P} \in H \text{ and stopping time } \tau.$$

At any given time t , we refer to E_t as an **e-value** for H .

Example: likelihood ratio process of any alternative \mathbb{Q} over \mathbb{P} .
An e-process is a nonparametric & composite generalization.

Interpretation as statistical evidence: how much wealth we multiply by betting against the odds put forth by \mathbb{P} .

Lemma. Fix any test threshold $z \in \mathbb{R}$. For any $\lambda \in [0, 1]$,

$$S(\lambda, z) = 1 + \lambda[\mathbf{1}(X \leq z) - \mathbf{1}(Y \leq z)]$$

is an **e-value** for $H_0(z)$. (λ is the "betting parameter").

Proposition. Given any alternative $\mathbb{Q} \notin H_0(z)$, there is a **growth-rate optimal (GRO)** bet:

$$\lambda^{\text{GRO}}(z) = \frac{\mathbb{Q}(X \leq z < Y) - \mathbb{Q}(Y \leq z < X)}{\mathbb{Q}(X \leq z < Y) + \mathbb{Q}(Y \leq z < X)}$$

Multiplying these e-values yields an e-process: $E_t = \prod_{\ell=1}^t S_\ell$.

Main result: Mixture GRO e-process is powerful

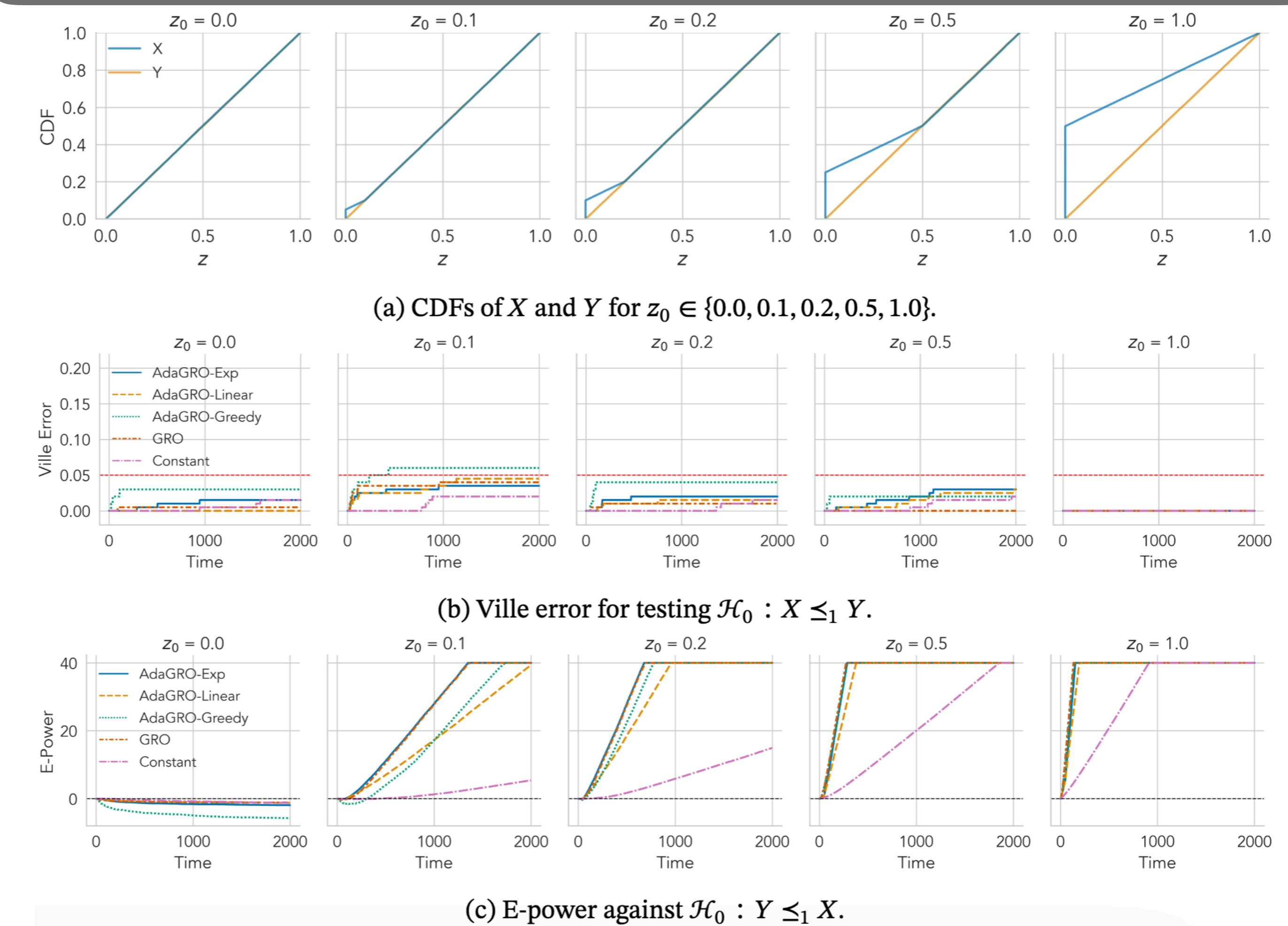
To test the *intersection* null H_0 , we can take mixtures over multiple test thresholds $Z_t \subset \mathbb{R}$ at each time t .

Theorem (GRO e-process is anytime-valid & powerful).

- (a) $(E_t)_{t \in \mathbb{N}}$ is an e-process for the 1-SD null H_0 .
- (b) For "reasonable" predictable mixtures $(\psi_t)_{t \in \mathbb{N}}$, $(E_t)_{t \in \mathbb{N}}$ is **powerful** against any non-1-SD alternative \mathbb{Q} :

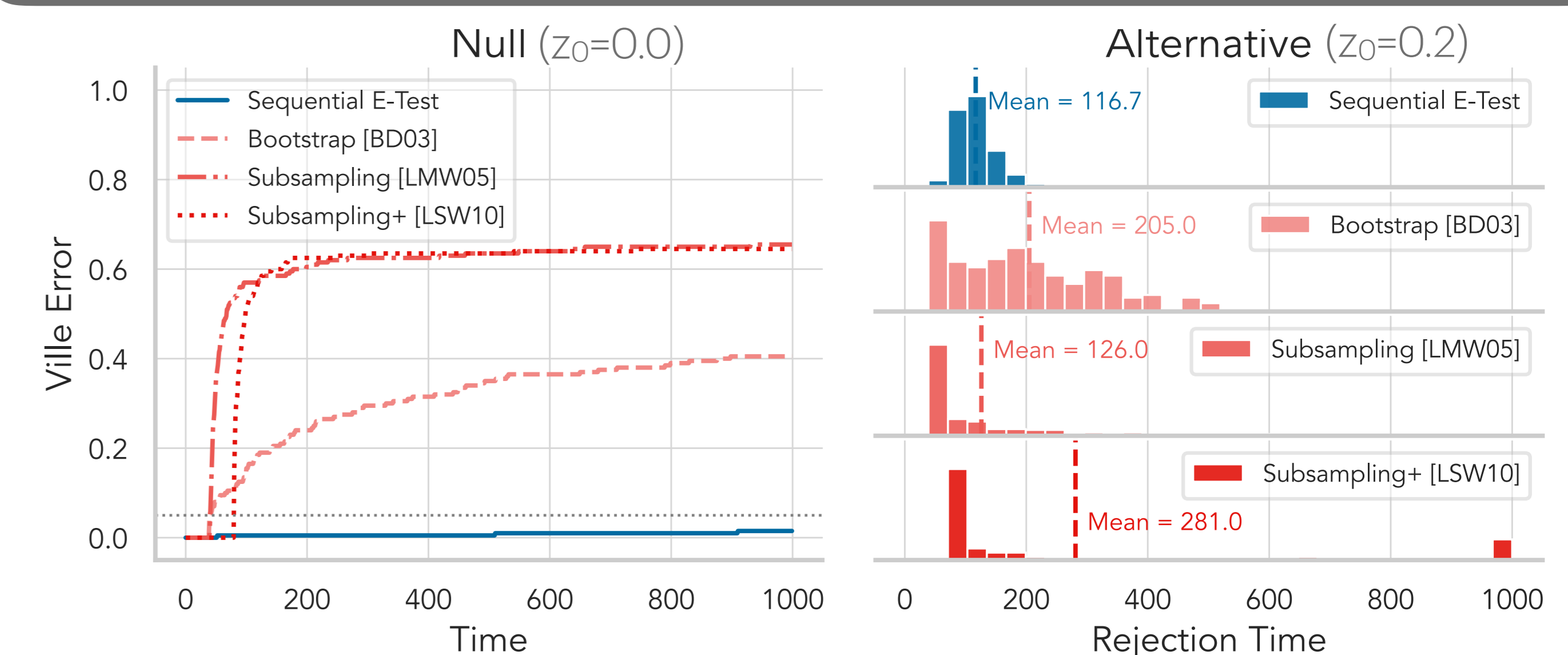
$$\mathbb{Q} \left(\liminf_{t \rightarrow \infty} E_t = \infty \right) = 1. \text{ (yields a test of power one)}$$

Simulation Experiments

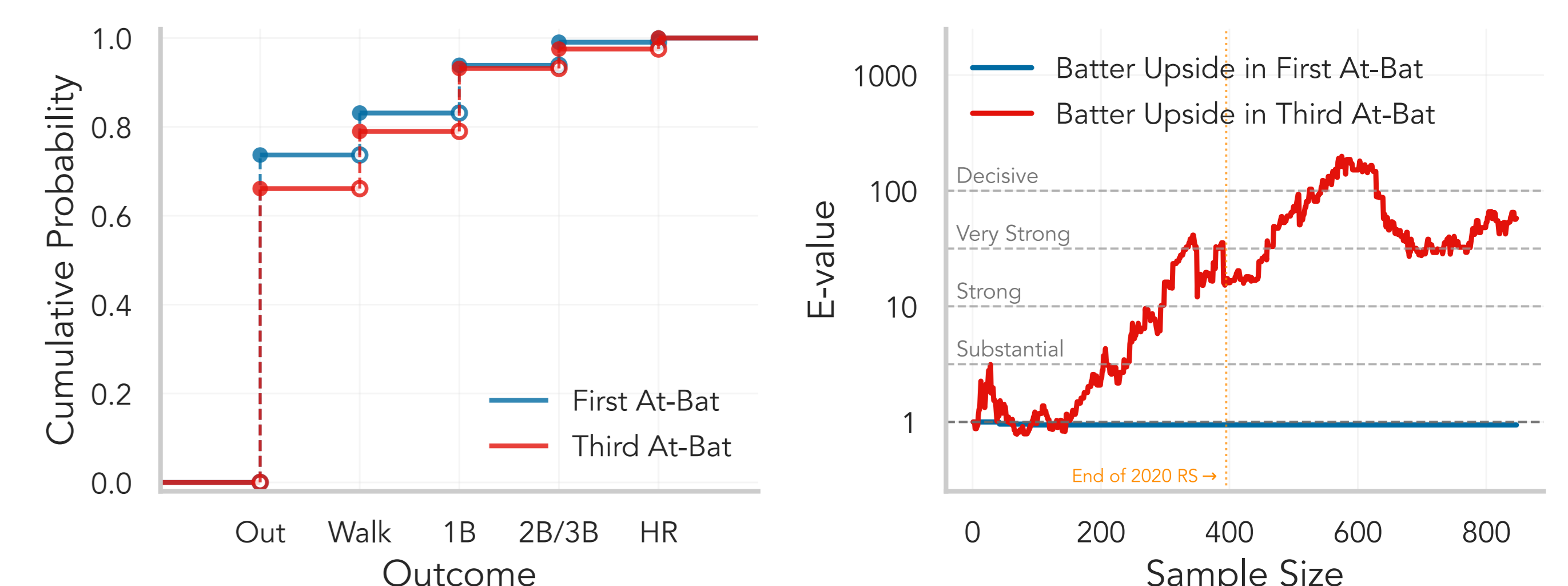


- Ville Error (cumulative type I error): $\hat{\mathbb{P}}(\exists \ell \leq t : E_\ell \geq 1/\alpha)$
- E-Power: $\hat{\mathbb{E}}[\log E_t]$ • Rejection Time: $\tau_\alpha = \inf \{t : E_t \geq 1/\alpha\}$

Comparison with non-SAVI methods



Application: Monitoring the 3TTO effect



Data: Every paired at-bat outcome against MLB pitcher Blake Snell, from 2016 to 2025 (regular seasons).

The sequential SD testing problem

X, Y : RVs with marginal CDFs F_X, F_Y .

Def. X **stochastically dominates** Y

if F_X falls entirely below F_Y :

$$Y \leq_1 X \iff F_X(z) \leq F_Y(z), \forall z.$$

"Y has no upside": At *every* threshold z , the prob. that Y exceeds z is no more than the prob. that X exceeds z .

The sequential SD testing problem:

With a stream of data $(X_1, Y_1), (X_2, Y_2), \dots \sim \mathbb{P}_{XY}$, test:

$$H_0 : Y \leq_1 X \text{ vs. } H_1 : Y \not\leq_1 X.$$

Fact. The null H_0 is an **intersection**: $H_0 = \cap_z H_0(z)$, where $H_0(z) = \{\mathbb{P} : \mathbb{P}(X \leq z) \leq \mathbb{P}(Y \leq z)\}$.