

# Improved Offline Contextual Bandits with Second-Order Bounds: Betting and Freezing

Jongha (Jon) Ryu

MIT EECS

RL Theory Seminars

December 16, 2025



Jeongyeol Kwon  
UW-Madison/Meta



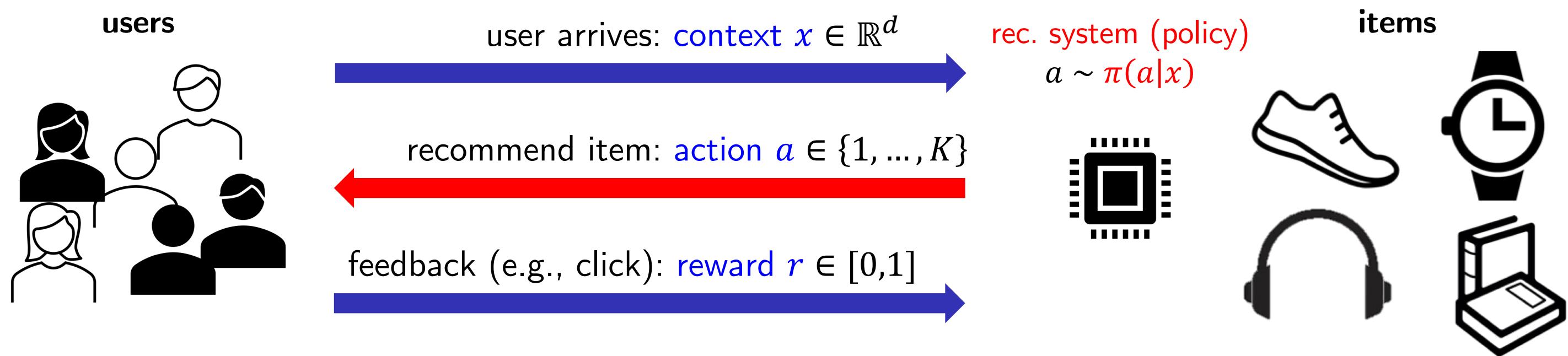
Benjamin Koppe  
Cornell



Kwang-Sung Jun  
U. of Arizona

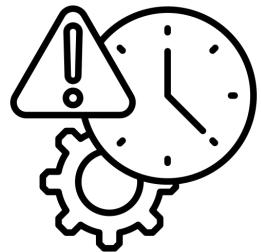
# Contextual Bandits

- **Example:** recommendation systems



- **Goal:** find  $\pi$  that maximizes cumulative reward!
- Ideally, we can **try** different policies  $\pi \in \Pi$  to find the best (via **online interaction**)

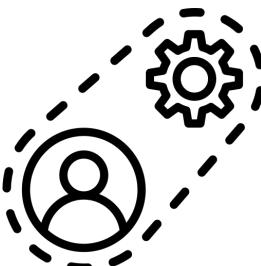
# Online Interaction is Expensive!



- System constraints



- Organizational policies



- Real-time feedback is expensive

Let's go **offline!**



# Problem: Off-Policy Contextual Bandit

- Behavior (logging) policy  $\pi_0(a|x)$
- Offline data  $D_n = \{(x_t, a_t, r_t)\}_{t=1}^n \sim p(x)\pi_0(a|x)p(r|a, x)$ 
  - For  $t = 1, \dots, n$ 
    - Observe context  $x_t \sim p(x)$
    - Choose action  $a_t \sim \pi_0(a|x_t)$
    - Observe reward  $r_t \sim p(r|a_t, x_t)$
- **Goal:** Given a policy class  $\Pi$ , find the best policy  $\hat{\pi} \in \arg \max_{\pi \in \Pi} v(\pi)$
- **Criterion:** Value  $v(\pi) \stackrel{\text{def}}{=} \mathbb{E}_{p(x)\pi(a|x)p(r|a,x)}[r]$  (expected reward)

# Off-Policy Contextual Bandit: Challenge

- **Goal:** Given a policy class  $\Pi$ , find the best policy  $\hat{\pi} \in \arg \max_{\pi \in \Pi} v(\pi)$

- Difficulty? logging  $\pi_0 \neq$  target  $\pi$ !
- At time  $t$ , given context  $x_t$

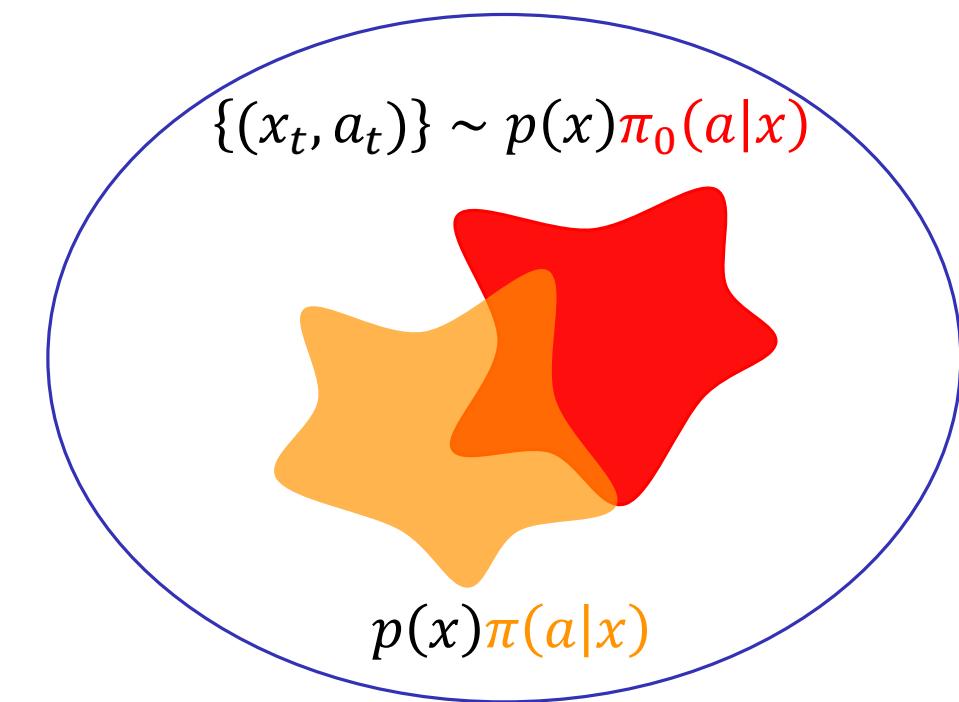
behavior policy  $\pi_0$  chose action 2

action	reward
1	1.0
2	0.1
3	0.5
...	
K	0.3

target policy  $\pi$  chooses action 3

observed  
unobserved!

data space  $\mathcal{X} \times \mathcal{A}$



Risk of **over-optimization**:  
a low-value policy  $\pi$  might  
seem good on  $\mathcal{D}$

# Off-Policy Problems

~35 min

~10 min

	evaluation	selection	learning (=optimization)
$ \Pi $ (size of policy class)	1	$< \infty$	$\infty$ (e.g., neural nets)
<b>Goal</b>	estimate $v(\pi)$	find $\operatorname{argmax}_{\pi \in \Pi} v(\pi)$	find $\operatorname{argmax}_{\pi \in \Pi} v(\pi)$
<b>Data split used</b>	test	validation	train
<b>Requirement</b>	none	none	"argmax-oracle-efficiency"

$$v(\pi) \stackrel{\text{def}}{=} \mathbb{E}_{p(x)\pi(a|x)p(r|a,x)}[r]$$

# Part 1. Off-Policy Selection

---

# Off-Policy Selection

# Off-Policy Selection

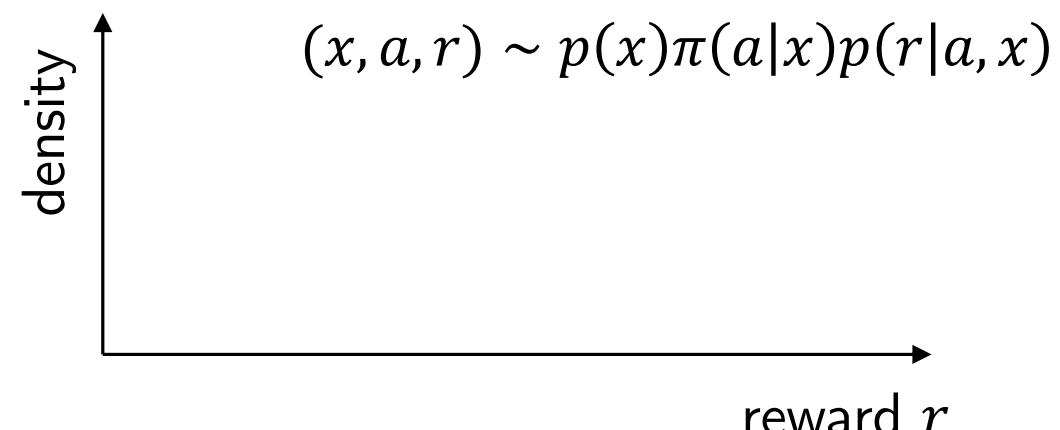
- Straightforward if we estimate **the value function**  $\pi \mapsto v(\pi)$ , using  $\mathcal{D}$

$$v(\pi) \stackrel{\text{def}}{=} \mathbb{E}_{p(x)\pi(a|x)p(r|a,x)}[r]$$

# Off-Policy Selection

- Straightforward if we estimate **the value function**  $\pi \mapsto v(\pi)$ , using  $\mathcal{D}$

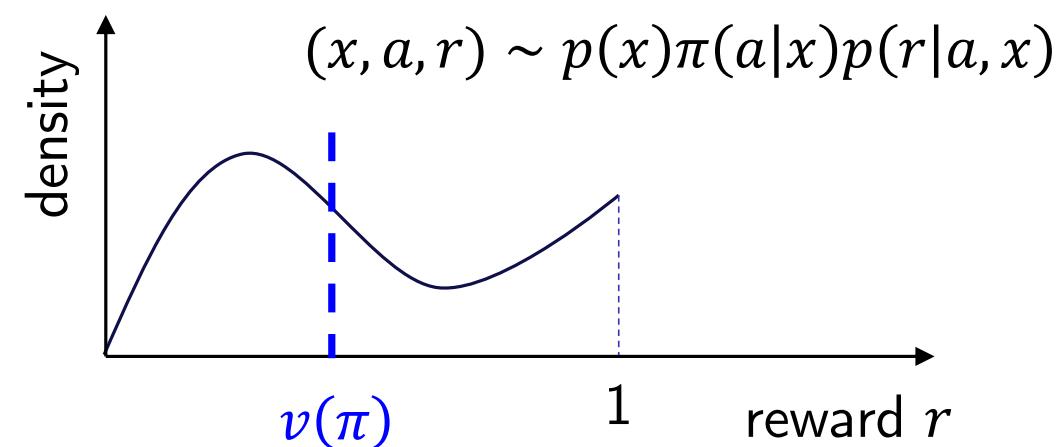
$$v(\pi) \stackrel{\text{def}}{=} \mathbb{E}_{p(x)\pi(a|x)p(r|a,x)}[r]$$



# Off-Policy Selection

- Straightforward if we estimate **the value function**  $\pi \mapsto v(\pi)$ , using  $\mathcal{D}$

$$v(\pi) \stackrel{\text{def}}{=} \mathbb{E}_{p(x)\pi(a|x)p(r|a,x)}[r]$$

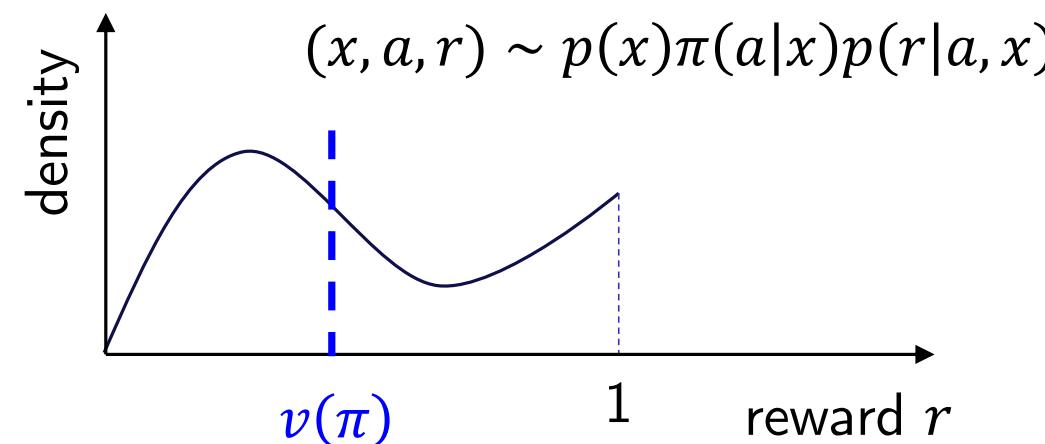


# Off-Policy Selection

- Straightforward if we estimate **the value function**  $\pi \mapsto v(\pi)$ , using  $\mathcal{D}$

$$\begin{aligned} v(\pi) &\stackrel{\text{def}}{=} \mathbb{E}_{p(x)\pi(a|x)p(r|a,x)}[r] \\ &= \mathbb{E}_{p(x)\pi_0(a|x)p(r|a,x)} \left[ \frac{\pi(a|x)}{\pi_0(a|x)} r \right] \end{aligned}$$

importance weight



# Off-Policy Selection

- Straightforward if we estimate **the value function**  $\pi \mapsto v(\pi)$ , using  $\mathcal{D}$

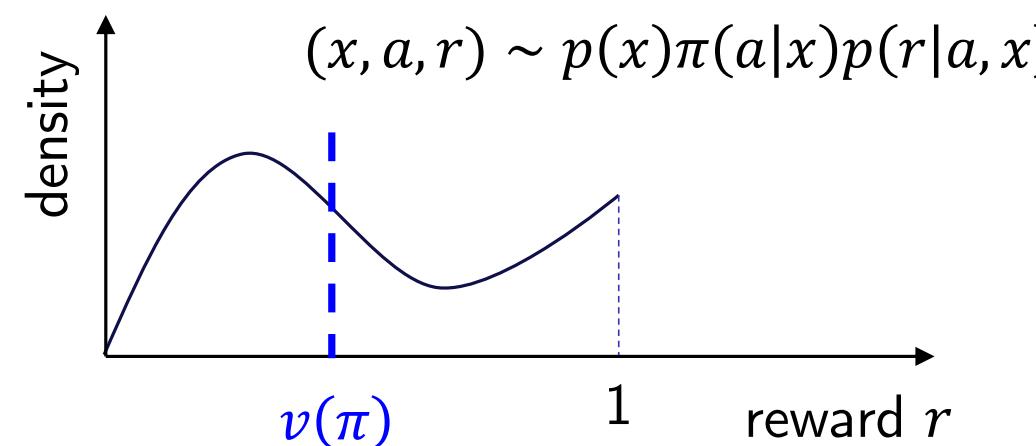
$$\begin{aligned} v(\pi) &\stackrel{\text{def}}{=} \mathbb{E}_{p(x)\pi(a|x)p(r|a,x)}[r] \\ &= \mathbb{E}_{p(x)\pi_0(a|x)p(r|a,x)} \left[ \frac{\pi(a|x)}{\pi_0(a|x)} r \right] \end{aligned}$$

*importance weight*

$$\approx \frac{1}{n} \sum_{t=1}^n \frac{\pi(a_t|x_t)}{\pi_0(a_t|x_t)} r_t \stackrel{\text{def}}{=} \text{IPS}(\pi)$$

**unbiased!**

inverse propensity score (IPS)  
inverse propensity weighting (IPW)  
importance weighted estimator (IW)



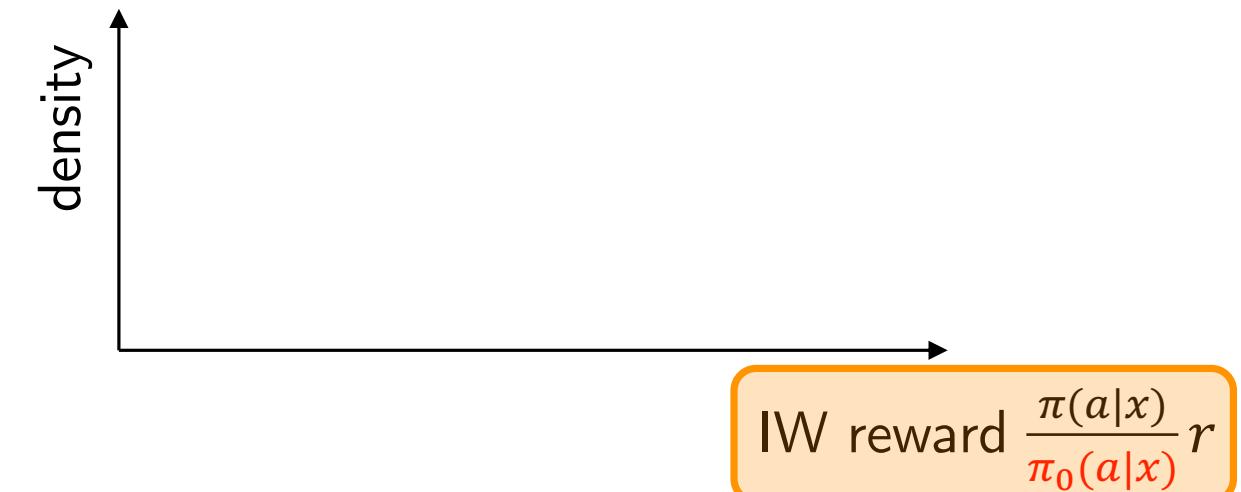
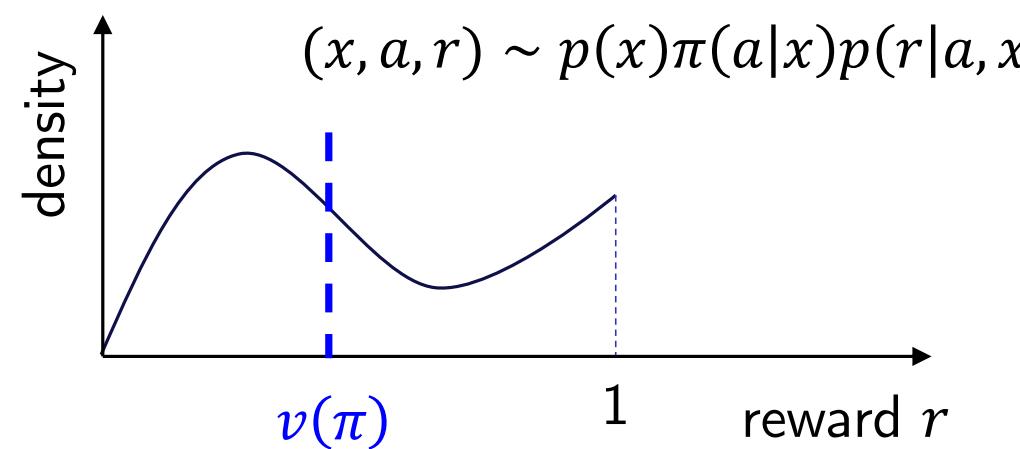
# Off-Policy Selection

- Straightforward if we estimate **the value function**  $\pi \mapsto v(\pi)$ , using  $\mathcal{D}$

$$\begin{aligned} v(\pi) &\stackrel{\text{def}}{=} \mathbb{E}_{p(x)\pi(a|x)p(r|a,x)}[r] \\ &= \mathbb{E}_{p(x)\pi_0(a|x)p(r|a,x)} \left[ \frac{\pi(a|x)}{\pi_0(a|x)} r \right] \quad \text{importance weight} \\ &\approx \frac{1}{n} \sum_{t=1}^n \frac{\pi(a_t|x_t)}{\pi_0(a_t|x_t)} r_t \stackrel{\text{def}}{=} \text{IPS}(\pi) \end{aligned}$$

unbiased!

inverse propensity score (IPS)  
inverse propensity weighting (IPW)  
importance weighted estimator (IW)



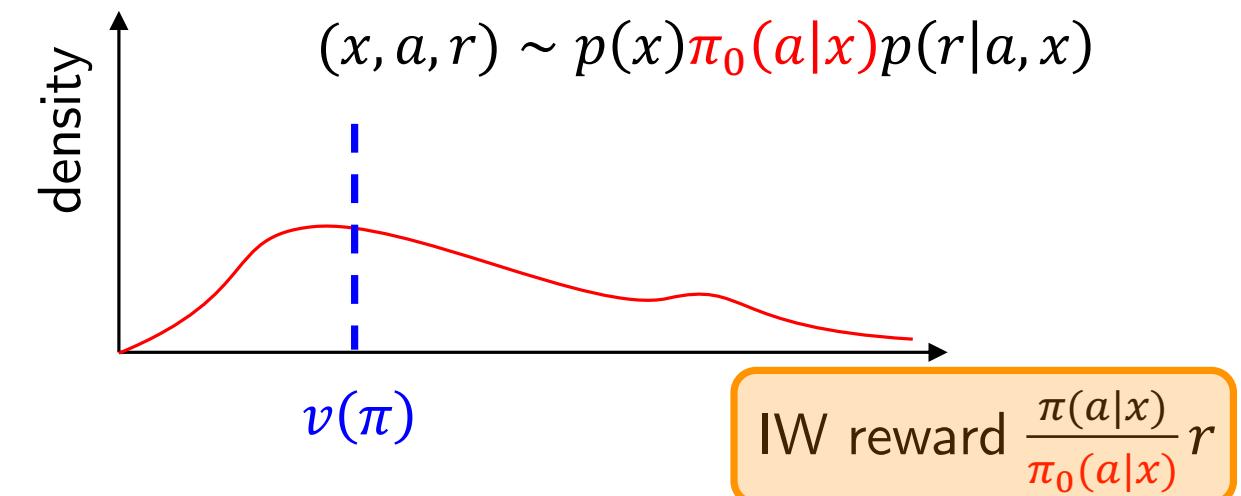
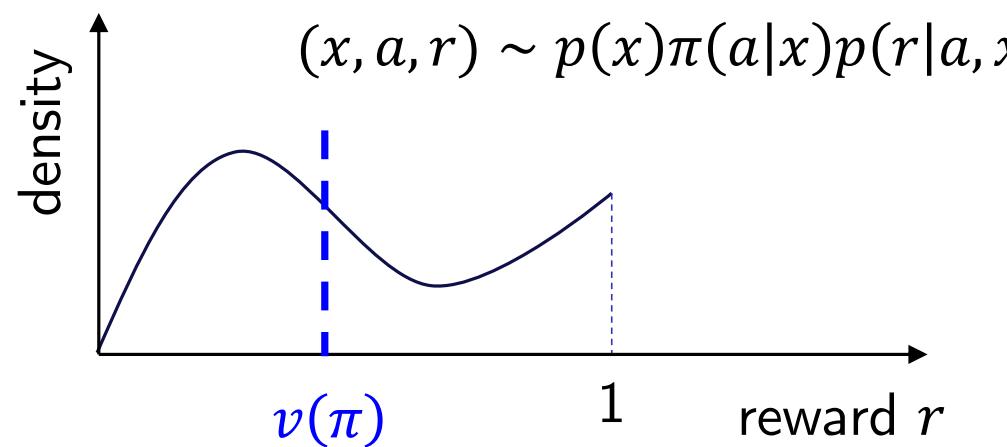
# Off-Policy Selection

- Straightforward if we estimate **the value function**  $\pi \mapsto v(\pi)$ , using  $\mathcal{D}$

$$\begin{aligned} v(\pi) &\stackrel{\text{def}}{=} \mathbb{E}_{p(x)\pi(a|x)p(r|a,x)}[r] \\ &= \mathbb{E}_{p(x)\pi_0(a|x)p(r|a,x)} \left[ \frac{\pi(a|x)}{\pi_0(a|x)} r \right] \xrightarrow{\text{importance weight}} \\ &\approx \frac{1}{n} \sum_{t=1}^n \frac{\pi(a_t|x_t)}{\pi_0(a_t|x_t)} r_t \stackrel{\text{def}}{=} \text{IPS}(\pi) \end{aligned}$$

unbiased!

inverse propensity score (IPS)  
inverse propensity weighting (IPW)  
importance weighted estimator (IW)



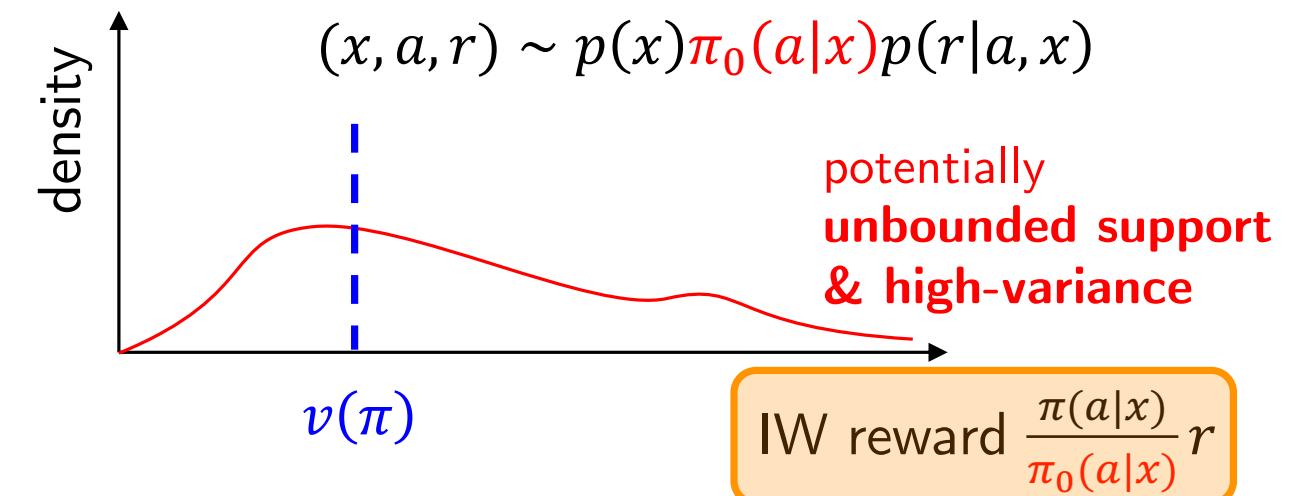
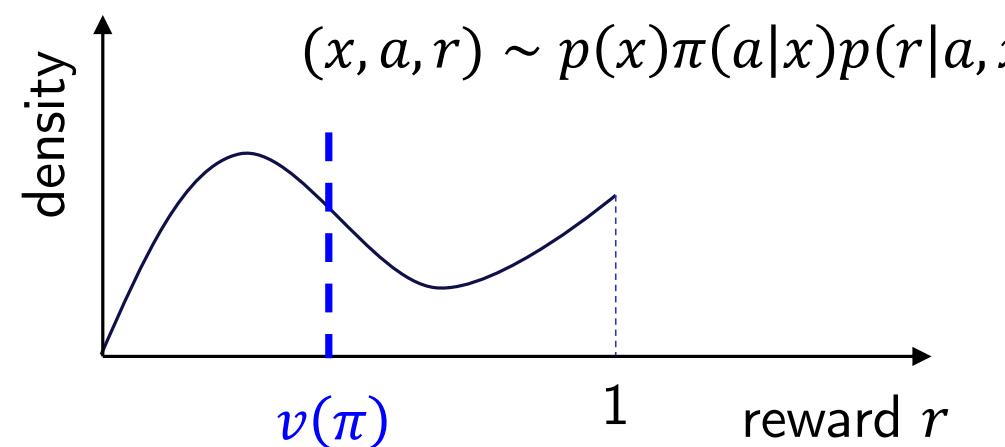
# Off-Policy Selection

- Straightforward if we estimate **the value function**  $\pi \mapsto v(\pi)$ , using  $\mathcal{D}$

$$\begin{aligned} v(\pi) &\stackrel{\text{def}}{=} \mathbb{E}_{p(x)\pi(a|x)p(r|a,x)}[r] \\ &= \mathbb{E}_{p(x)\pi_0(a|x)p(r|a,x)} \left[ \frac{\pi(a|x)}{\pi_0(a|x)} r \right] \xrightarrow{\text{importance weight}} \\ &\approx \frac{1}{n} \sum_{t=1}^n \frac{\pi(a_t|x_t)}{\pi_0(a_t|x_t)} r_t \stackrel{\text{def}}{=} \text{IPS}(\pi) \end{aligned}$$

unbiased!

inverse propensity score (IPS)  
inverse propensity weighting (IPW)  
importance weighted estimator (IW)

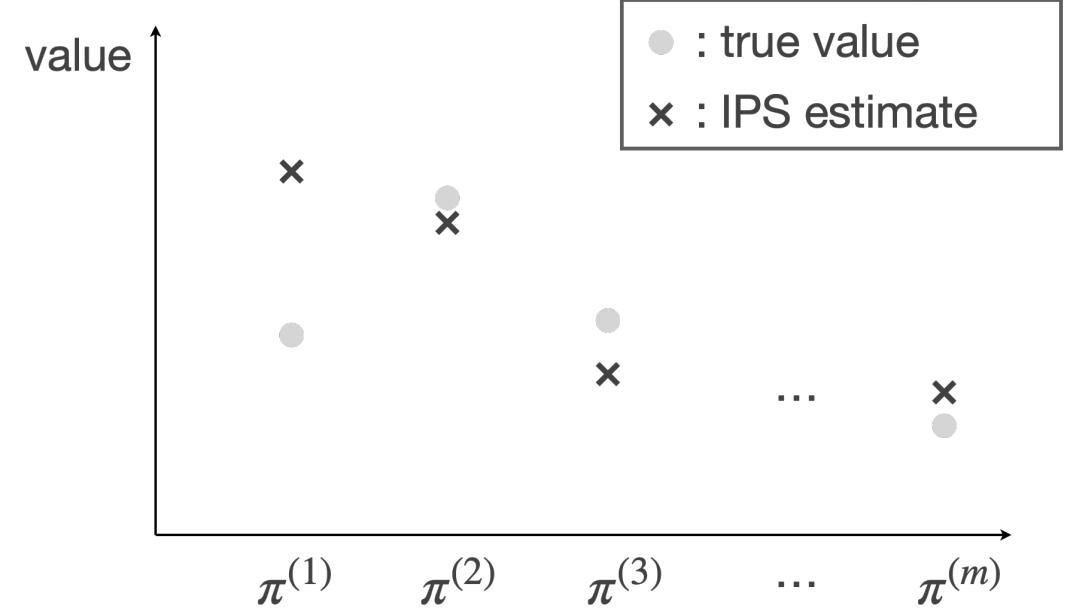


# Off-Policy Selection: IPS and Pessimism

- $\text{IPS}(\pi)$  is unbiased but might be high-variance!

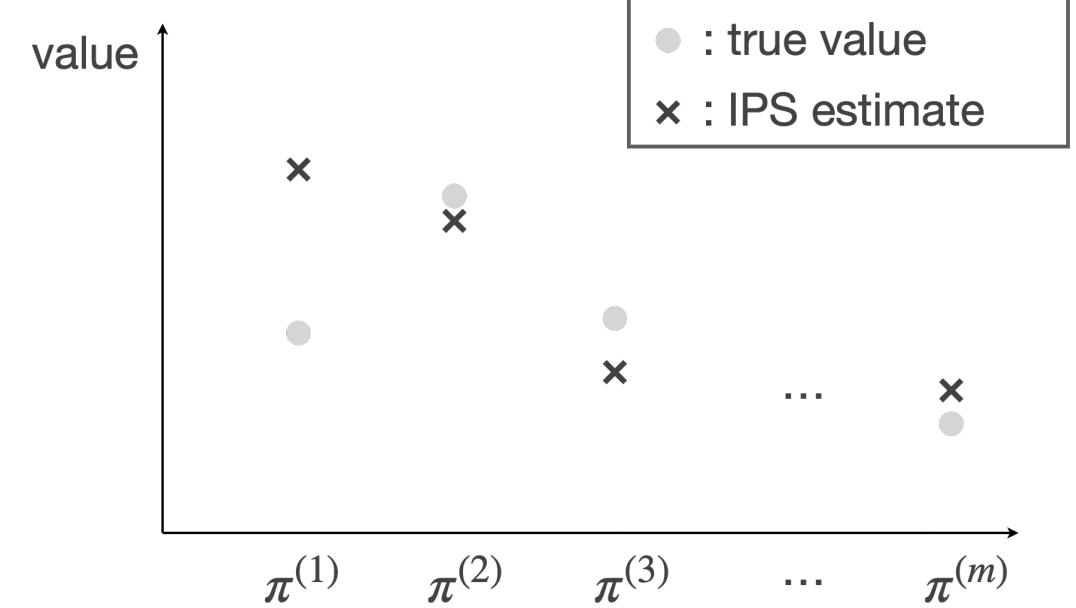
# Off-Policy Selection: IPS and Pessimism

- IPS( $\pi$ ) is unbiased but might be high-variance!



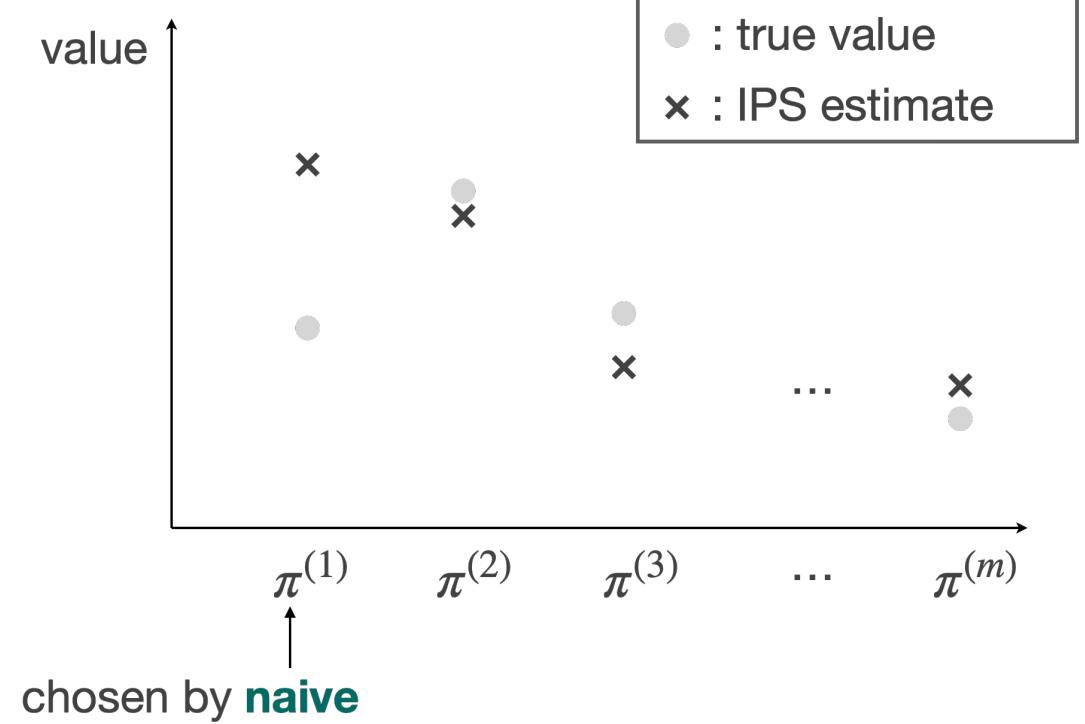
# Off-Policy Selection: IPS and Pessimism

- $\text{IPS}(\pi)$  is unbiased but might be high-variance!
- Naive selection with  $\text{IPS}(\pi)$  (i.e.,  $\underset{\pi \in \Pi}{\operatorname{argmax}} \text{IPS}(\pi)$ )



# Off-Policy Selection: IPS and Pessimism

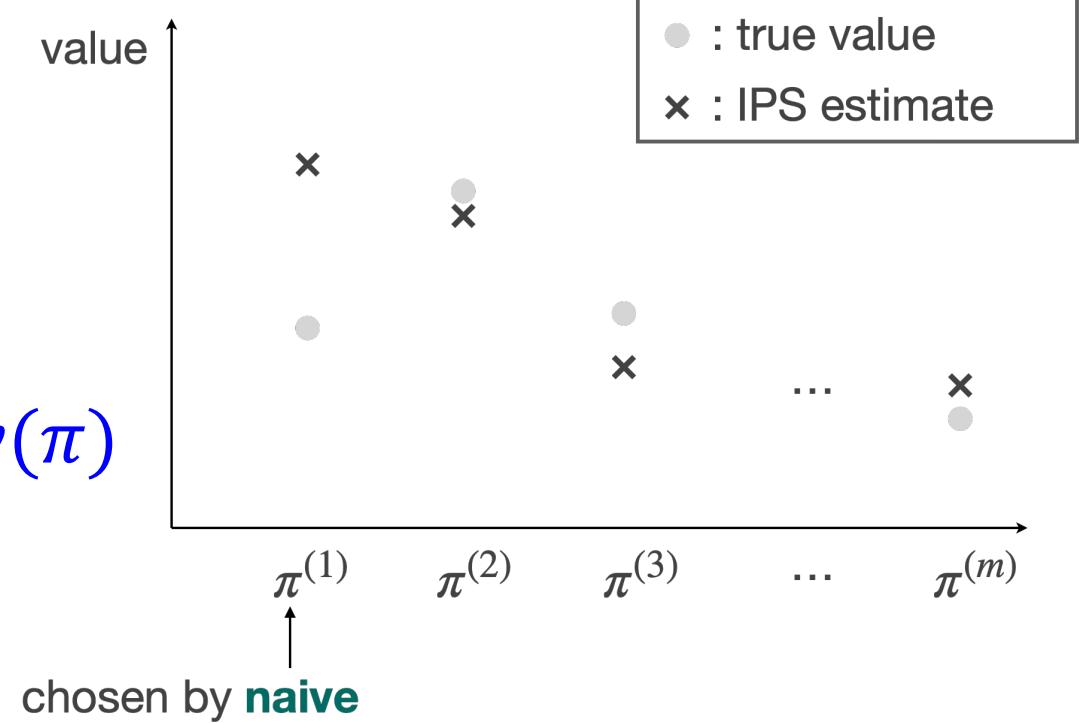
- $\text{IPS}(\pi)$  is unbiased but might be high-variance!
- Naive selection with  $\text{IPS}(\pi)$  (i.e.,  $\underset{\pi \in \Pi}{\text{argmax}} \text{IPS}(\pi)$ ) might pick up low-value, high-variance policy



# Off-Policy Selection: IPS and Pessimism

- $\text{IPS}(\pi)$  is unbiased but might be high-variance!
- Naive selection with  $\text{IPS}(\pi)$  (i.e.,  $\underset{\pi \in \Pi}{\text{argmax}} \text{IPS}(\pi)$ ) might pick up low-value, high-variance policy
- Pessimism: Use lower confidence bound (LCB) on  $v(\pi)$

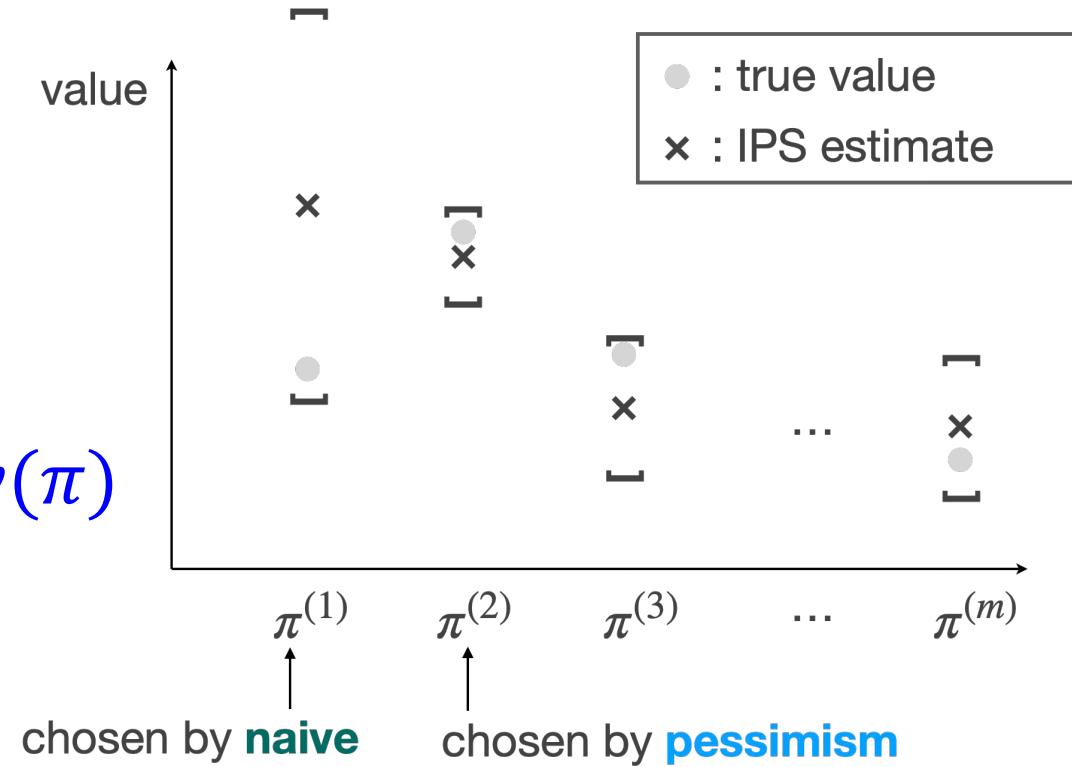
$$\underset{\pi \in \Pi}{\text{argmax}} \text{LCB}(\pi)$$



# Off-Policy Selection: IPS and Pessimism

- $\text{IPS}(\pi)$  is unbiased but might be high-variance!
- Naive selection with  $\text{IPS}(\pi)$  (i.e.,  $\underset{\pi \in \Pi}{\operatorname{argmax}} \text{IPS}(\pi)$ ) might pick up low-value, high-variance policy
- Pessimism: Use lower confidence bound (LCB) on  $v(\pi)$

$$\underset{\pi \in \Pi}{\operatorname{argmax}} \text{LCB}(\pi)$$



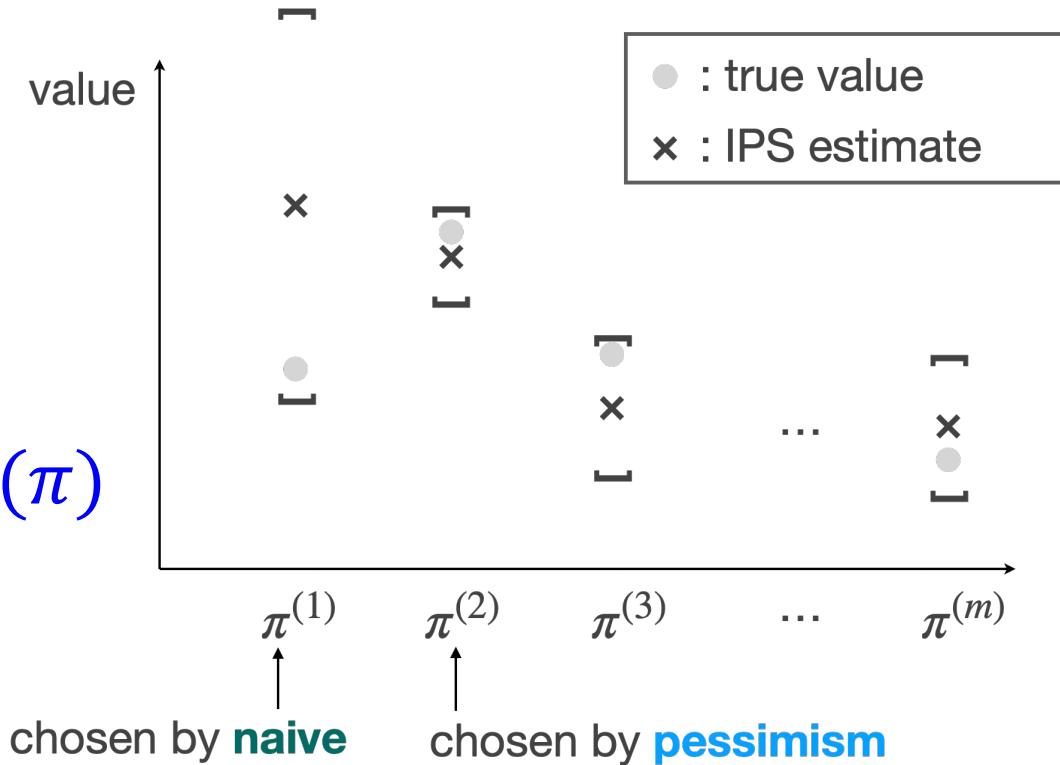
# Off-Policy Selection: IPS and Pessimism

- $\text{IPS}(\pi)$  is unbiased but might be high-variance! value
  - **Naive selection** with  $\text{IPS}(\pi)$  (i.e.,  $\underset{\pi \in \Pi}{\operatorname{argmax}} \text{IPS}(\pi)$ ) might pick up low-value, high-variance policy
  - **Pessimism**: Use lower confidence bound (LCB) on  $v(\pi)$

$$\underset{\pi \in \Pi}{\operatorname{argmax}} \text{LCB}(\pi)$$

# Can avoid over-optimization!

(caveat: high-value, high-variance policy might be still not visible)



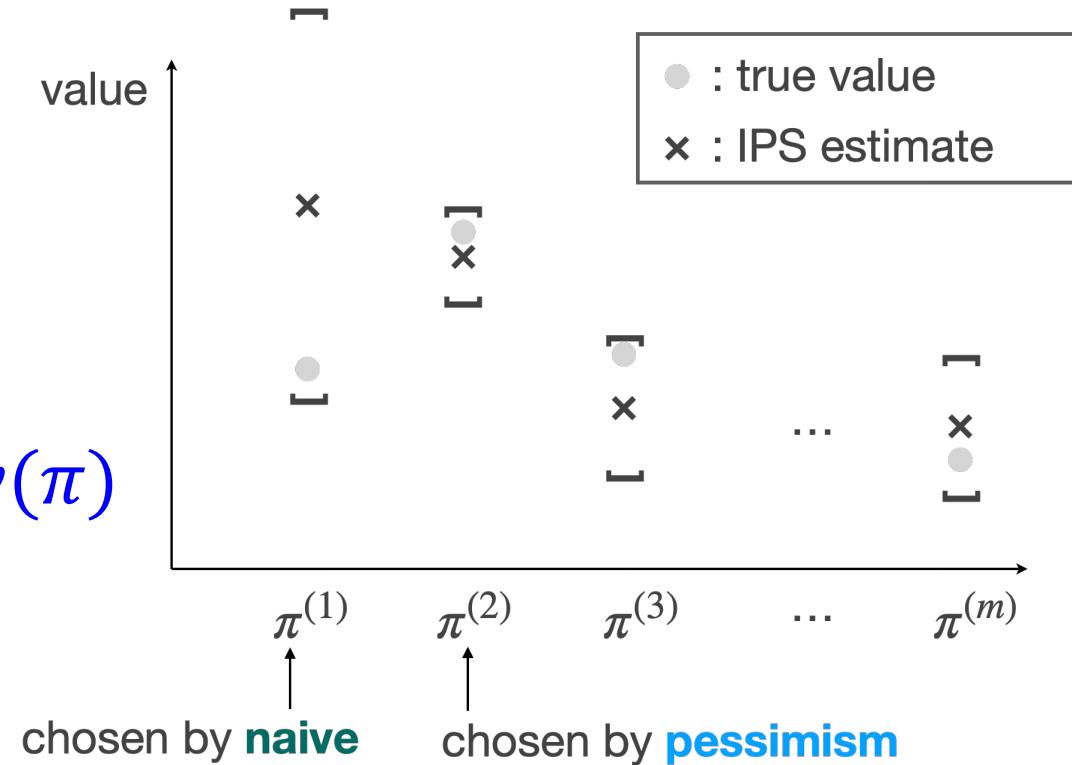
# Off-Policy Selection: IPS and Pessimism

- $\text{IPS}(\pi)$  is unbiased but might be high-variance!
- Naive selection with  $\text{IPS}(\pi)$  (i.e.,  $\underset{\pi \in \Pi}{\operatorname{argmax}} \text{IPS}(\pi)$ ) might pick up low-value, high-variance policy
- Pessimism: Use lower confidence bound (LCB) on  $v(\pi)$

$$\underset{\pi \in \Pi}{\operatorname{argmax}} \text{LCB}(\pi)$$

Can avoid over-optimization!

(caveat: high-value, high-variance policy might be still not visible)



## Our contribution

We propose a “tight” LCB that yields provably good selection via pessimism,  
(1) with variance-adaptive property, (2) without hyperparameter!

# Pessimism by semi-Unbounded coin Betting (PUB)

# Pessimism by semi-Unbounded coin Betting (PUB)

- New LCB based on betting + universal portfolio (will elaborate soon) 
$$\text{LCB}_n^{\text{UP}}(\pi)$$

# Pessimism by semi-Unbounded coin Betting (PUB)

- New LCB based on betting + universal portfolio (will elaborate soon) 
- **Pessimism by semi-Unbounded coin Betting (PUB)**  $\hat{\pi} = \arg \max_{\pi \in \Pi} \text{LCB}_n^{\text{UP}}(\pi)$

# Pessimism by semi-Unbounded coin Betting (PUB)

- New LCB based on betting + universal portfolio (will elaborate soon) 
- **Pessimism by semi-Unbounded coin Betting (PUB)**  $\hat{\pi} = \arg \max_{\pi \in \Pi} \text{LCB}_n^{\text{UP}}(\pi)$

$$\forall \pi^* \in \Pi, \quad v(\pi^*) - v(\hat{\pi}) \lesssim \sqrt{\text{Var}(\tilde{r}_1(\pi^*)) \frac{1}{n} \ln \frac{|\Pi|}{\delta}}$$

$\tilde{r}_t(\pi) \triangleq \frac{\pi(a_t|x_t)}{\pi_0(a_t|x_t)} r_t$

optimality gap

# Pessimism by semi-Unbounded coin Betting (PUB)

- New LCB based on betting + universal portfolio (will elaborate soon) 
- **Pessimism by semi-Unbounded coin Betting (PUB)**  $\hat{\pi} = \arg \max_{\pi \in \Pi} \text{LCB}_n^{\text{UP}}(\pi)$

$$\forall \pi^* \in \Pi, \quad v(\pi^*) - v(\hat{\pi}) \lesssim \sqrt{\underbrace{\text{Var}(\tilde{r}_1(\pi^*))}_{\text{variance-adaptive, w/o hyperparameter tuning!}}} \frac{1}{n} \ln \frac{|\Pi|}{\delta}$$

$\tilde{r}_t(\pi) \triangleq \frac{\pi(a_t|x_t)}{\pi_0(a_t|x_t)} r_t$

optimality gap

# Pessimism by semi-Unbounded coin Betting (PUB)

- New LCB based on betting + universal portfolio (will elaborate soon) 
- **Pessimism by semi-Unbounded coin Betting (PUB)**  $\hat{\pi} = \arg \max_{\pi \in \Pi} \text{LCB}_n^{\text{UP}}(\pi)$

$$\forall \pi^* \in \Pi, \quad v(\pi^*) - v(\hat{\pi}) \lesssim \underbrace{\sqrt{\text{Var}(\tilde{r}_1(\pi^*))}}_{\text{optimality gap}} \frac{1}{n} \ln \frac{|\Pi|}{\delta} \quad \tilde{r}_t(\pi) \triangleq \frac{\pi(a_t|x_t)}{\pi_0(a_t|x_t)} r_t$$

**variance-adaptive, w/o hyperparameter tuning!**

**Prior art:** Logarithmic smoothing (Sakhi+2024)

$$\lesssim (1 + \mathbb{E}[(\tilde{r}_1(\pi^*))^2]) \sqrt{\frac{1}{n} \ln \frac{|\Pi|}{\delta}}$$

# Pessimism by semi-Unbounded coin Betting (PUB)

- New LCB based on betting + universal portfolio (will elaborate soon) 
- **Pessimism by semi-Unbounded coin Betting (PUB)**  $\hat{\pi} = \arg \max_{\pi \in \Pi} \text{LCB}_n^{\text{UP}}(\pi)$

$$\forall \pi^* \in \Pi, \quad v(\pi^*) - v(\hat{\pi}) \lesssim \sqrt{\underbrace{\text{Var}(\tilde{r}_1(\pi^*))}_{\text{optimality gap}}} \frac{1}{n} \ln \frac{|\Pi|}{\delta} \quad \tilde{r}_t(\pi) \triangleq \frac{\pi(a_t|x_t)}{\pi_0(a_t|x_t)} r_t$$

**variance-adaptive, w/o hyperparameter tuning!**

**Prior art:** Logarithmic smoothing (Sakhi+2024)

$$\lesssim (1 + \mathbb{E}[(\tilde{r}_1(\pi^*))^2]) \sqrt{\frac{1}{n} \ln \frac{|\Pi|}{\delta}}$$

$$\sqrt{\text{Var}(\tilde{r}_1(\pi^*))} \leq 1 + \text{Var}(\tilde{r}_1(\pi^*)) \leq 1 + \mathbb{E}[(\tilde{r}_1(\pi^*))^2]$$

**two-step improvement!**

# Comparison to Prior Arts

	Requires boundedness / moment condition?	Rate guarantee (=confidence width)	Bound is adaptive to
Empirical Bernstein [MaurerP2009]	Y	Y	$\sqrt{\text{Var}(\tilde{r}(\pi^*))}$
Oracle-efficient pessimism [Wang+2024]	Y	Y	("pseudo-loss") $\times$ (worst-ratio factor)
[Waudby-Smith+2022]	N	unknown	unknown
Implicit exploration (IX) [Gabbianelli+2024]	N	Y	"smoothed coverage ratio"
Logarithmic smoothing (LS) [Sakhi+2024]	N	Y	$1 + \mathbb{E}[(\tilde{r}(\pi^*))^2]$
<b>PUB (ours)</b>	N	Y	$\sqrt{\text{Var}(\tilde{r}(\pi^*))}$

# Remarks on Coverage Measures

$$v(\pi) - v(\hat{\pi}) \lesssim \boxed{?} \sqrt{\frac{1}{n} \ln \frac{|\Pi|}{\delta}}$$

worst



average

# Remarks on Coverage Measures

$$v(\pi) - v(\hat{\pi}) \lesssim \boxed{?} \sqrt{\frac{1}{n} \ln \frac{|\Pi|}{\delta}}$$

- In RL, “worst-case” coverage

$$\max_{(a,x)} \frac{\pi(a|x)}{\pi_0(a|x)}$$

worst



average

# Remarks on Coverage Measures

$$v(\pi) - v(\hat{\pi}) \lesssim \boxed{?} \sqrt{\frac{1}{n} \ln \frac{|\Pi|}{\delta}}$$

- In RL, “worst-case” coverage
- Pseudo-loss [Wang+2024]

$$\max_{(a,x)} \frac{\pi(a|x)}{\pi_0(a|x)}$$

worst

$$\mathbb{E}_{p(x)} \left[ \sum_a \frac{\pi(a|x)}{\pi_0(a|x)} \right]$$



average

# Remarks on Coverage Measures

$$v(\pi) - v(\hat{\pi}) \lesssim ? \sqrt{\frac{1}{n} \ln \frac{|\Pi|}{\delta}}$$

- In RL, “worst-case” coverage
- Pseudo-loss [Wang+2024]
- Smoothed coverage ratio (IX) [Gabbianelli+2024]

$$\begin{aligned} & \max_{(a,x)} \frac{\pi(a|x)}{\pi_0(a|x)} \\ \xrightarrow{\text{E}_{p(x)} \left[ \sum_a \frac{\pi(a|x)}{\pi_0(a|x)} \right]} \quad & \mathbb{E}_{p(x)} \left[ \sum_a \frac{\pi(a|x)}{\pi_0(a|x) + \gamma \mathbb{E}_{p(r|x,a)}[r]} \right] \end{aligned}$$

worst



average

# Remarks on Coverage Measures

$$v(\pi) - v(\hat{\pi}) \lesssim \boxed{?} \sqrt{\frac{1}{n} \ln \frac{|\Pi|}{\delta}}$$

- In RL, “worst-case” coverage

$$\max_{(a,x)} \frac{\pi(a|x)}{\pi_0(a|x)}$$

worst

- Pseudo-loss [Wang+2024]

$$\mathbb{E}_{p(x)} \left[ \sum_a \frac{\pi(a|x)}{\pi_0(a|x)} \right]$$

- Smoothed coverage ratio (IX) [Gabbianelli+2024]

$$\mathbb{E}_{p(x)} \left[ \sum_a \frac{\pi(a|x)}{\pi_0(a|x) + \gamma} \mathbb{E}_{p(r|x,a)}[r] \right]$$

- Second moment (LS) [Sakhi+2024]

$$\mathbb{E}[(\tilde{r}_1(\pi^*))^2] = \mathbb{E}_{p(x)\pi_0(a|x)p(r|x,a)} \left[ \left( \frac{\pi(a|x)}{\pi_0(a|x)} r \right)^2 \right]$$



average

# Remarks on Coverage Measures

$$v(\pi) - v(\hat{\pi}) \lesssim \boxed{?} \sqrt{\frac{1}{n} \ln \frac{|\Pi|}{\delta}}$$

- In RL, “worst-case” coverage
- Pseudo-loss [Wang+2024]
- Smoothed coverage ratio (IX) [Gabbianelli+2024]
- Second moment (LS) [Sakhi+2024]
- Variance (EB, PUB (ours))

worst

$$\max_{(a,x)} \frac{\pi(a|x)}{\pi_0(a|x)}$$

$$\mathbb{E}_{p(x)} \left[ \sum_a \frac{\pi(a|x)}{\pi_0(a|x)} \right]$$

$$\mathbb{E}_{p(x)} \left[ \sum_a \frac{\pi(a|x)}{\pi_0(a|x) + \gamma} \mathbb{E}_{p(r|x,a)}[r] \right]$$

$$\mathbb{E}[(\tilde{r}_1(\pi^*))^2] = \mathbb{E}_{p(x)\pi_0(a|x)p(r|x,a)} \left[ \left( \frac{\pi(a|x)}{\pi_0(a|x)} r \right)^2 \right]$$

$$\text{Var}(\tilde{r}_1(\pi^*)) = \text{Var}_{p(x)\pi_0(a|x)p(r|x,a)} \left( \frac{\pi(a|x)}{\pi_0(a|x)} r \right)$$

average

# Remarks on Coverage Measures

$$v(\pi) - v(\hat{\pi}) \lesssim \boxed{?} \sqrt{\frac{1}{n} \ln \frac{|\Pi|}{\delta}}$$

- In RL, “worst-case” coverage

$$\max_{(a,x)} \frac{\pi(a|x)}{\pi_0(a|x)}$$

worst

- Pseudo-loss [Wang+2024]

$$\mathbb{E}_{p(x)} \left[ \sum_a \frac{\pi(a|x)}{\pi_0(a|x)} \right]$$

$$\mathbb{E}_{p(x)} \left[ \sum_a \frac{\pi(a|x)}{\pi_0(a|x) + \gamma} \mathbb{E}_{p(r|x,a)}[r] \right]$$

- Smoothed coverage ratio (IX) [Gabbianelli+2024]

- Second moment (LS) [Sakhi+2024]

$$\mathbb{E}[(\tilde{r}_1(\pi^*))^2] = \mathbb{E}_{p(x)\pi_0(a|x)p(r|x,a)} \left[ \left( \frac{\pi(a|x)}{\pi_0(a|x)} r \right)^2 \right]$$

- Variance (EB, PUB (ours))

$$\text{Var}(\tilde{r}_1(\pi^*)) = \text{Var}_{p(x)\pi_0(a|x)p(r|x,a)} \left( \frac{\pi(a|x)}{\pi_0(a|x)} r \right)$$



average

# IX vs. LS vs. PUB

- Smoothed coverage ratio  
(IX) [Gabbianelli+2024]

$$\mathbb{E}_{p(x)} \left[ \sum_a \frac{\pi(a|x)}{\pi_0(a|x) + \gamma} \mathbb{E}_{p(r|x,a)}[r] \right]$$

- Second moment  
(LS) [Sakhi+2024]

$$\mathbb{E}[(\tilde{r}_1(\pi^*))^2] = \mathbb{E}_{p(x)\pi_0(a|x)p(r|x,a)} \left[ \left( \frac{\pi(a|x)}{\pi_0(a|x)} r \right)^2 \right]$$

- Variance  
(EB, PUB (ours))

$$\text{Var}(\tilde{r}_1(\pi^*)) = \text{Var}_{p(x)\pi_0(a|x)p(r|x,a)} \left( \frac{\pi(a|x)}{\pi_0(a|x)} r \right)$$

# IX vs. LS vs. PUB

- Smoothed coverage ratio  
(IX) [Gabbianelli+2024]

$$\mathbb{E}_{p(x)\pi_0(a|x)p(r|x,a)} \left[ \frac{(\tilde{r}(\pi))^2}{(\pi_0(a|x) + \gamma)\tilde{r}(\pi)} \right]$$

- Smoothed** second moment  
(LS) [Sakhi+2024]

$$\mathbb{E}_{p(x)\pi_0(a|x)p(r|x,a)} \left[ \frac{(\tilde{r}(\pi))^2}{1 + \gamma\tilde{r}(\pi)} \right]$$

- Smoothed** variance  
(PUB (ours))

$$\mathbb{E}_{p(x)\pi_0(a|x)p(r|x,a)} \left[ \frac{(\tilde{r}(\pi) - v(\pi))^2}{1 + \gamma(\tilde{r}(\pi) - v(\pi))} \right]$$

# IX vs. LS vs. PUB

- Smoothed coverage ratio  
(IX) [Gabbianelli+2024]

$$\mathbb{E}_{p(x)\pi_0(a|x)p(r|x,a)} \left[ \frac{(\tilde{r}(\pi))^2}{(\pi_0(a|x) + \gamma)\tilde{r}(\pi)} \right]$$

VI

- Smoothed** second moment  
(LS) [Sakhi+2024]

$$\mathbb{E}_{p(x)\pi_0(a|x)p(r|x,a)} \left[ \frac{(\tilde{r}(\pi))^2}{1 + \gamma\tilde{r}(\pi)} \right]$$

VI

- Smoothed** variance  
(PUB (ours))

$$\mathbb{E}_{p(x)\pi_0(a|x)p(r|x,a)} \left[ \frac{(\tilde{r}(\pi) - v(\pi))^2}{1 + \gamma(\tilde{r}(\pi) - v(\pi))} \right]$$

# IX vs. LS vs. PUB

- Smoothed coverage ratio  
(IX) [Gabbianelli+2024]
- Smoothed** second moment  
(LS) [Sakhi+2024]
- Smoothed** variance  
(PUB (ours))

$$\mathbb{E}_{p(x)\pi_0(a|x)p(r|x,a)} \left[ \frac{(\tilde{r}(\pi))^2}{(\pi_0(a|x) + \gamma)\tilde{r}(\pi)} \right] \quad \text{VI}$$

$$\mathbb{E}_{p(x)\pi_0(a|x)p(r|x,a)} \left[ \frac{(\tilde{r}(\pi))^2}{1 + \gamma\tilde{r}(\pi)} \right] \quad \text{VI}$$

$$\mathbb{E}_{p(x)\pi_0(a|x)p(r|x,a)} \left[ \frac{(\tilde{r}(\pi) - v(\pi))^2}{1 + \gamma(\tilde{r}(\pi) - v(\pi))} \right]$$

$\gamma$  needs to be tuned;  
tuned with oracle,  
becomes  $\text{sqrt}(\dots)$

# IX vs. LS vs. PUB

- Smoothed coverage ratio  
(IX) [Gabbianelli+2024]
- Smoothed** second moment  
(LS) [Sakhi+2024]
- Smoothed** variance  
(PUB (ours))

$$\mathbb{E}_{p(x)\pi_0(a|x)p(r|x,a)} \left[ \frac{(\tilde{r}(\pi))^2}{(\pi_0(a|x) + \gamma)\tilde{r}(\pi)} \right] \quad \text{VI}$$

$$\mathbb{E}_{p(x)\pi_0(a|x)p(r|x,a)} \left[ \frac{(\tilde{r}(\pi))^2}{1 + \gamma\tilde{r}(\pi)} \right] \quad \text{VI}$$

$$\mathbb{E}_{p(x)\pi_0(a|x)p(r|x,a)} \left[ \frac{(\tilde{r}(\pi) - v(\pi))^2}{1 + \gamma(\tilde{r}(\pi) - v(\pi))} \right]$$

$\gamma$  needs to be tuned;  
tuned with oracle,  
becomes  $\text{sqrt}(\dots)$

↑  
we can achieve  $\text{sqrt}(\dots)$  without tuning!  
(universal portfolio has auto-tuning feature ☺)

# *Technical Interlude: On the LCB Construction via Betting*

---

# Lower Confidence Bound

- Observe  $\textcolor{blue}{Y}_1, \dots, Y_n \in [0, \infty)$  with  $\mathbb{E}[Y_t] = \mu$  (unknown mean) (e.g.,  $Y_t = \tilde{r}_t(\pi)$ )
- **Want:** construct a LCB  $\ell_n(\delta) = \ell_n(Y_{1:n}; \delta)$  such that  $\mathbb{P}(\mu \geq \ell_n(\delta)) \geq 1 - \delta$

# Lower Confidence Bound

- Observe  $Y_1, \dots, Y_n \in [0, \infty)$  with  $\mathbb{E}[Y_t] = \mu$  (unknown mean) (e.g.,  $Y_t = \tilde{r}_t(\pi)$ )
- Want: construct a LCB  $\ell_n(\delta) = \ell_n(Y_{1:n}; \delta)$  such that  $\mathbb{P}(\mu \geq \ell_n(\delta)) \geq 1 - \delta$

**Naive attempts:**

Empirical Bernstein?

$$\frac{1}{n} \sum_{t=1}^n Y_t - \sqrt{2\widehat{\text{Var}}_n \frac{1}{n} \ln \frac{|\Pi|}{\delta}}$$

# Lower Confidence Bound

- Observe  $Y_1, \dots, Y_n \in [0, \infty)$  with  $\mathbb{E}[Y_t] = \mu$  (unknown mean) (e.g.,  $Y_t = \tilde{r}_t(\pi)$ )
- Want: construct a LCB  $\ell_n(\delta) = \ell_n(Y_{1:n}; \delta)$  such that  $\mathbb{P}(\mu \geq \ell_n(\delta)) \geq 1 - \delta$

**Naive attempts:**

Empirical Bernstein?

$$\frac{1}{n} \sum_{t=1}^n Y_t - \sqrt{2\widehat{\text{Var}}_n \frac{1}{n} \ln \frac{|\Pi|}{\delta}}$$

**Issue 1.** theoretically, analysis relies on boundedness/tail condition on  $Y_t$   
**Issue 2.** practically, need finite fourth moment at least

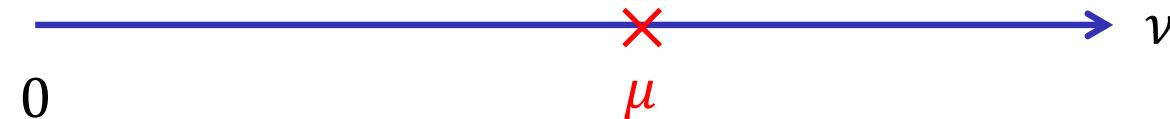
# Lower Confidence Bound via Testing

- Observe  $Y_1, \dots, Y_n \in [0, \infty)$  with  $\mathbb{E}[Y_t] = \mu$  (unknown mean) (e.g.,  $Y_t = \tilde{r}_t(\pi)$ )
- **Want:** construct a LCB  $\ell_n(\delta) = \ell_n(Y_{1:n}; \delta)$  such that  $\mathbb{P}(\mu \geq \ell_n(\delta)) \geq 1 - \delta$

# Lower Confidence Bound via Testing

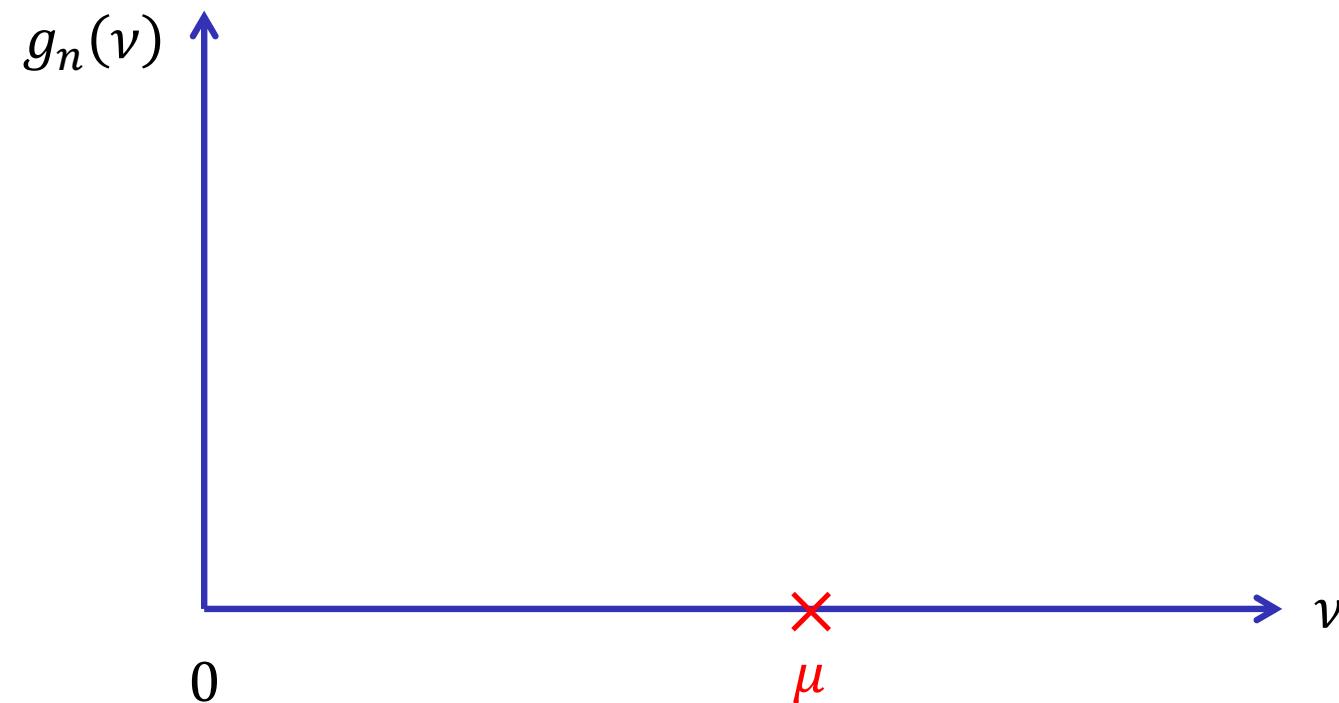
- Observe  $Y_1, \dots, Y_n \in [0, \infty)$  with  $\mathbb{E}[Y_t] = \mu$  (unknown mean) (e.g.,  $Y_t = \tilde{r}_t(\pi)$ )
- Want: construct a LCB  $\ell_n(\delta) = \ell_n(Y_{1:n}; \delta)$  such that  $\mathbb{P}(\mu \geq \ell_n(\delta)) \geq 1 - \delta$

For each  $\nu \in (0, \infty)$ , test: is  $\nu = \mu$ ?



# Lower Confidence Bound via Testing

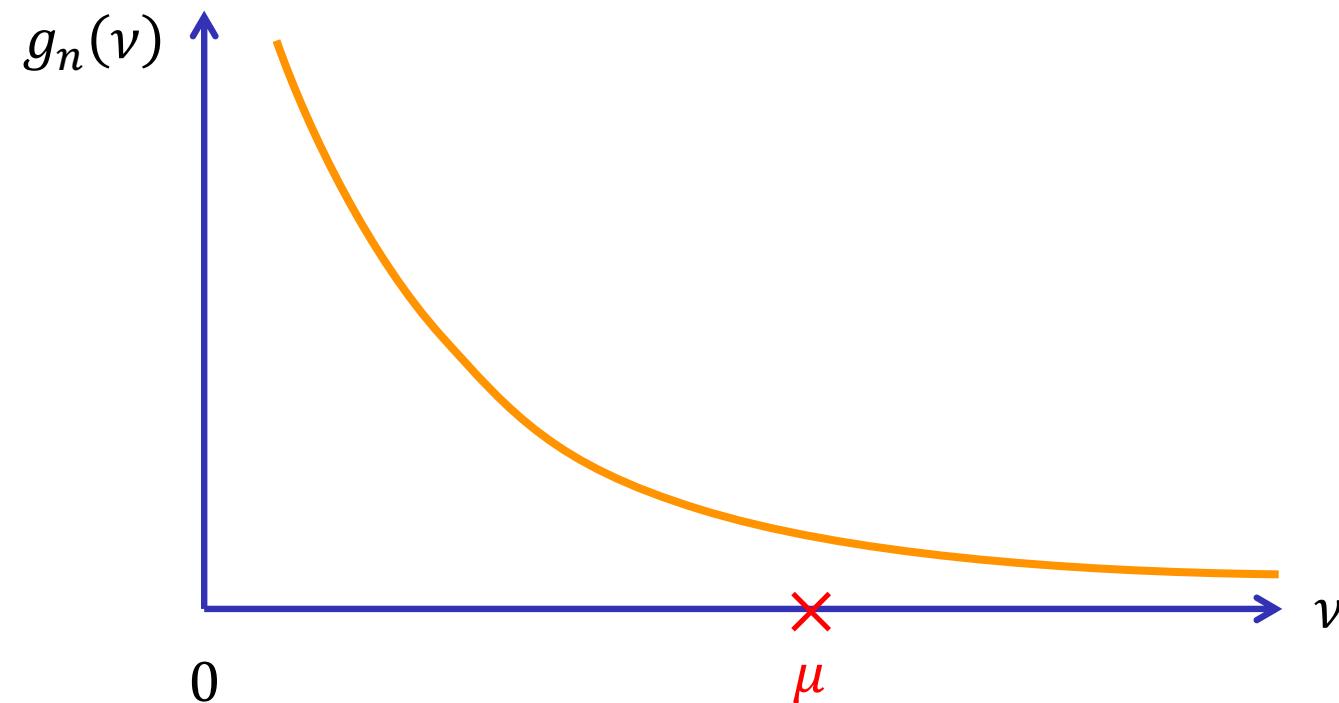
- Observe  $Y_1, \dots, Y_n \in [0, \infty)$  with  $\mathbb{E}[Y_t] = \mu$  (unknown mean) (e.g.,  $Y_t = \tilde{r}_t(\pi)$ )
- Want: construct a LCB  $\ell_n(\delta) = \ell_n(Y_{1:n}; \delta)$  such that  $\mathbb{P}(\mu \geq \ell_n(\delta)) \geq 1 - \delta$



For each  $v \in (0, \infty)$ , test: is  $v = \mu$ ?  
Design a "potential function"  $v \mapsto g_n(v)$ :

# Lower Confidence Bound via Testing

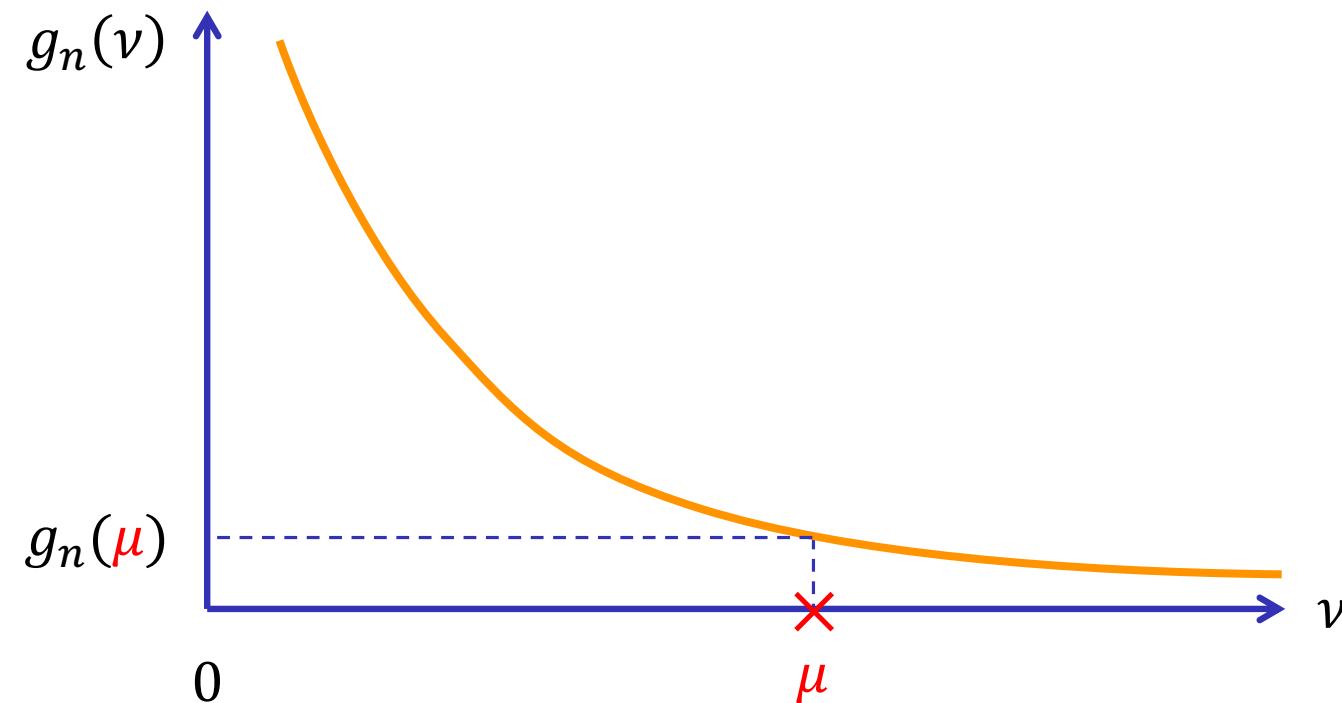
- Observe  $Y_1, \dots, Y_n \in [0, \infty)$  with  $\mathbb{E}[Y_t] = \mu$  (unknown mean) (e.g.,  $Y_t = \tilde{r}_t(\pi)$ )
- Want: construct a LCB  $\ell_n(\delta) = \ell_n(Y_{1:n}; \delta)$  such that  $\mathbb{P}(\mu \geq \ell_n(\delta)) \geq 1 - \delta$



For each  $v \in (0, \infty)$ , test: is  $v = \mu$ ?  
Design a "potential function"  $v \mapsto g_n(v)$ :  
1) monotonically decreasing;

# Lower Confidence Bound via Testing

- Observe  $Y_1, \dots, Y_n \in [0, \infty)$  with  $\mathbb{E}[Y_t] = \mu$  (unknown mean) (e.g.,  $Y_t = \tilde{r}_t(\pi)$ )
- Want: construct a LCB  $\ell_n(\delta) = \ell_n(Y_{1:n}; \delta)$  such that  $\mathbb{P}(\mu \geq \ell_n(\delta)) \geq 1 - \delta$

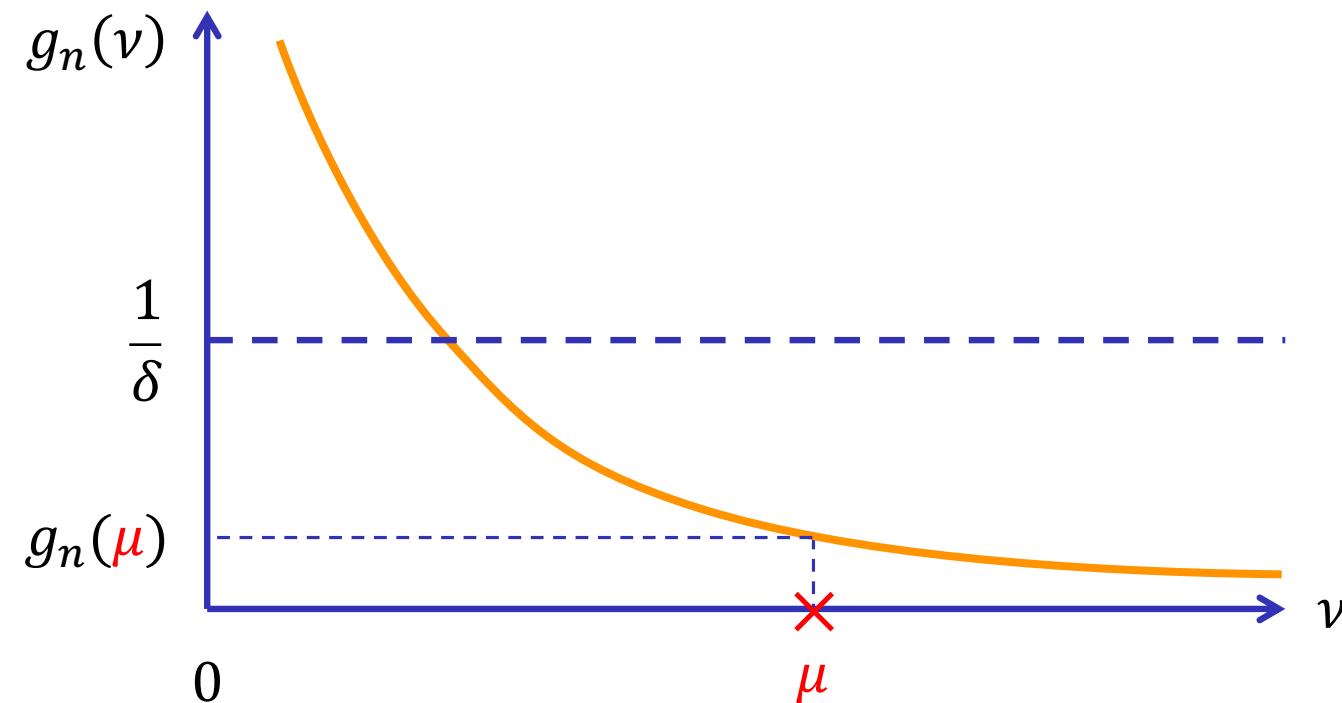


For each  $v \in (0, \infty)$ , test: is  $v = \mu$ ?  
Design a "potential function"  $v \mapsto g_n(v)$ :

- 1) monotonically decreasing;
- 2)  $g_n(\mu)$  must be "small" with high prob.

# Lower Confidence Bound via Testing

- Observe  $Y_1, \dots, Y_n \in [0, \infty)$  with  $\mathbb{E}[Y_t] = \mu$  (unknown mean) (e.g.,  $Y_t = \tilde{r}_t(\pi)$ )
- Want: construct a LCB  $\ell_n(\delta) = \ell_n(Y_{1:n}; \delta)$  such that  $\mathbb{P}(\mu \geq \ell_n(\delta)) \geq 1 - \delta$

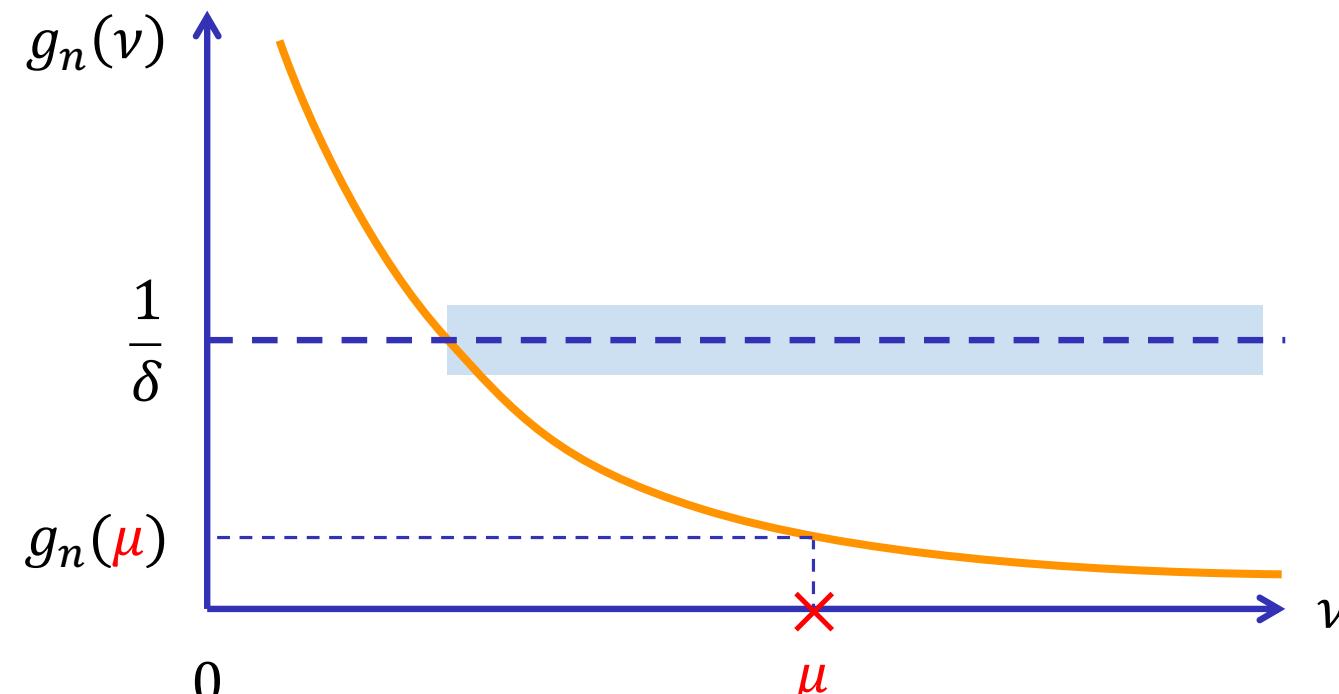


For each  $v \in (0, \infty)$ , test: is  $v = \mu$ ?  
Design a "potential function"  $v \mapsto g_n(v)$ :

- 1) monotonically decreasing;
- 2)  $g_n(\mu) \leq \frac{1}{\delta}$  with prob.  $1 - \delta$

# Lower Confidence Bound via Testing

- Observe  $Y_1, \dots, Y_n \in [0, \infty)$  with  $\mathbb{E}[Y_t] = \mu$  (unknown mean) (e.g.,  $Y_t = \tilde{r}_t(\pi)$ )
- Want: construct a LCB  $\ell_n(\delta) = \ell_n(Y_{1:n}; \delta)$  such that  $\mathbb{P}(\mu \geq \ell_n(\delta)) \geq 1 - \delta$



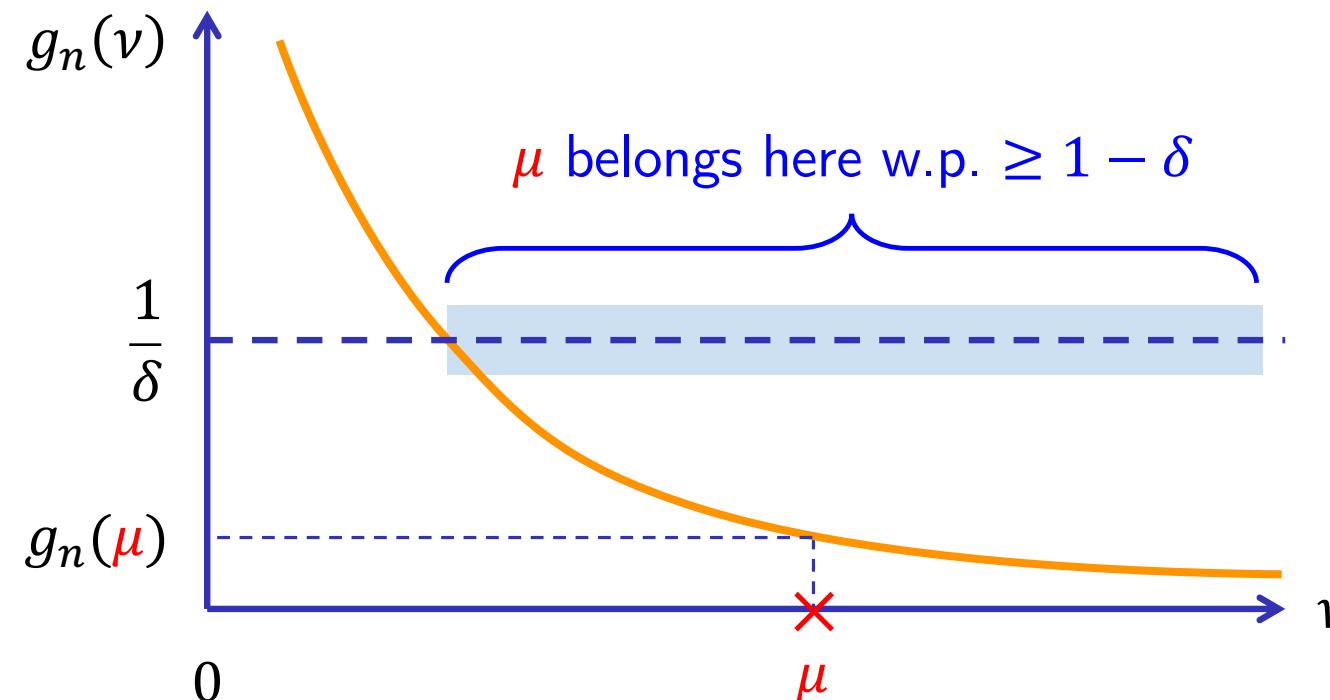
$$\left\{ v \in (0, \infty) \mid g_n(v) \leq \frac{1}{\delta} \right\}$$

For each  $v \in (0, \infty)$ , test: is  $v = \mu$ ?  
Design a "potential function"  $v \mapsto g_n(v)$ :

- 1) monotonically decreasing;
- 2)  $g_n(\mu) \leq \frac{1}{\delta}$  with prob.  $1 - \delta$

# Lower Confidence Bound via Testing

- Observe  $Y_1, \dots, Y_n \in [0, \infty)$  with  $\mathbb{E}[Y_t] = \mu$  (unknown mean) (e.g.,  $Y_t = \tilde{r}_t(\pi)$ )
- Want: construct a LCB  $\ell_n(\delta) = \ell_n(Y_{1:n}; \delta)$  such that  $\mathbb{P}(\mu \geq \ell_n(\delta)) \geq 1 - \delta$



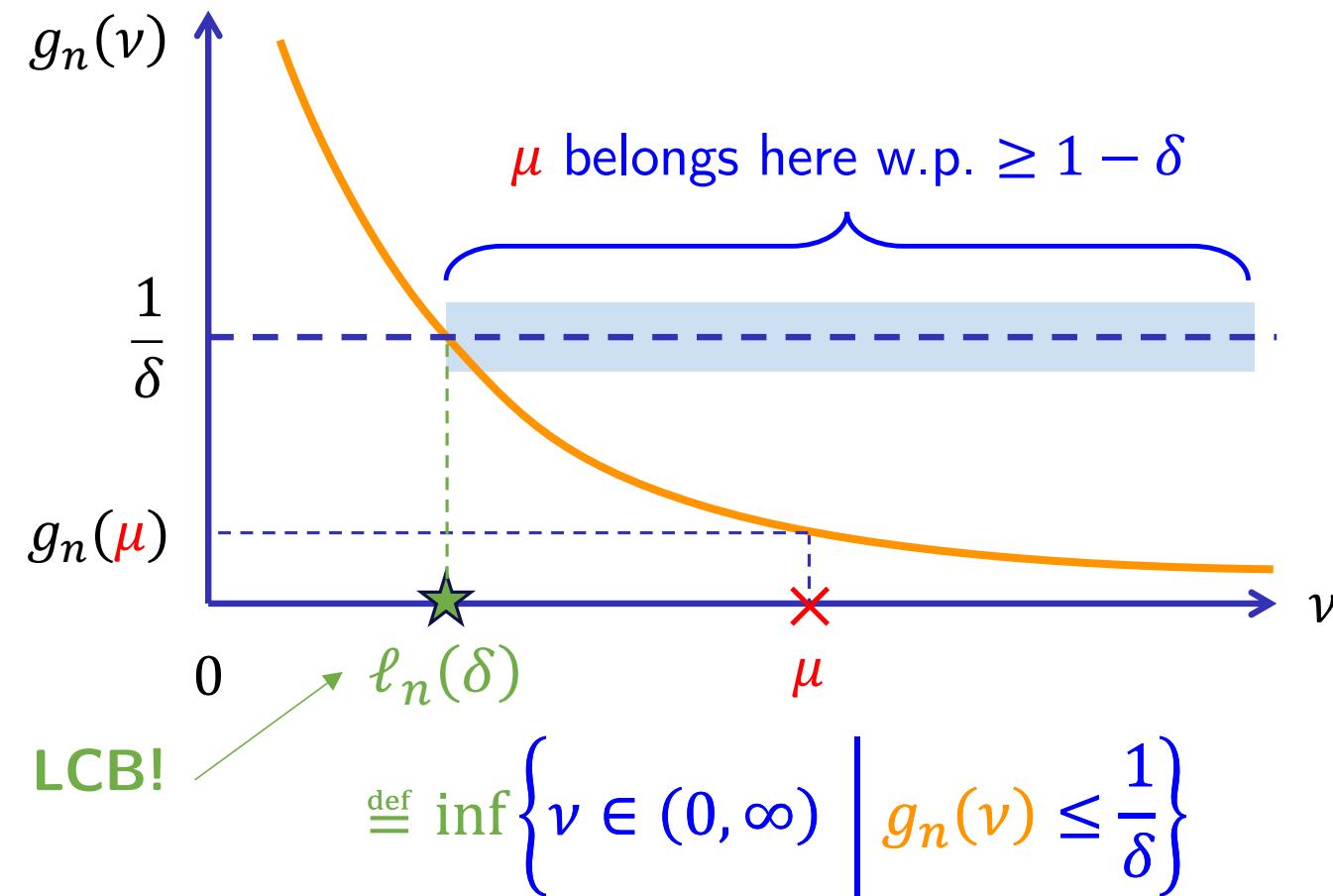
$$\left\{ v \in (0, \infty) \mid g_n(v) \leq \frac{1}{\delta} \right\}$$

For each  $v \in (0, \infty)$ , test: is  $v = \mu$ ?  
Design a "potential function"  $v \mapsto g_n(v)$ :

- 1) monotonically decreasing;
- 2)  $g_n(\mu) \leq \frac{1}{\delta}$  with prob.  $1 - \delta$

# Lower Confidence Bound via Testing

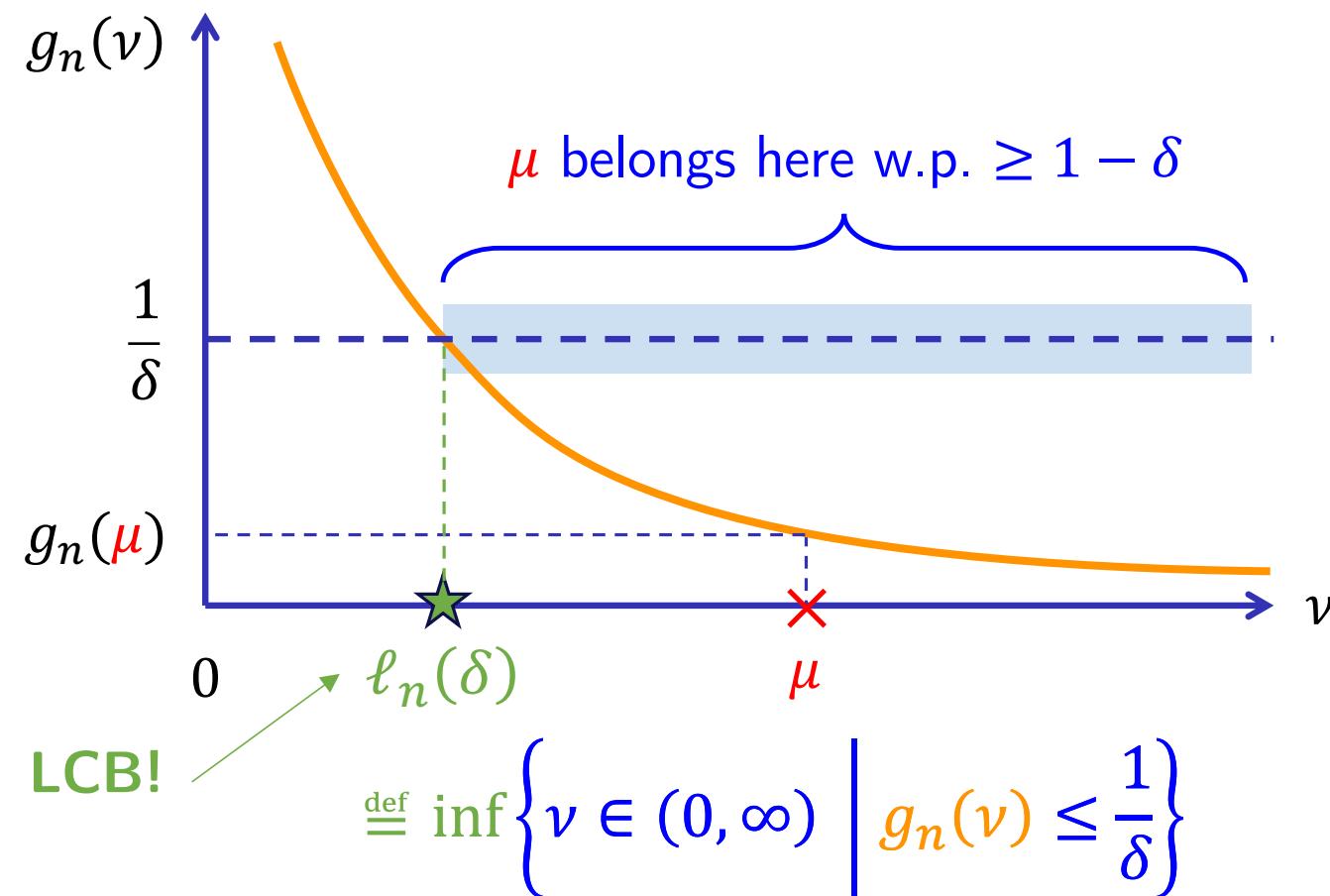
- Observe  $Y_1, \dots, Y_n \in [0, \infty)$  with  $\mathbb{E}[Y_t] = \mu$  (unknown mean) (e.g.,  $Y_t = \tilde{r}_t(\pi)$ )
  - **Want:** construct a LCB  $\ell_n(\delta) = \ell_n(Y_{1:n}; \delta)$  such that  $\mathbb{P}(\mu \geq \ell_n(\delta)) \geq 1 - \delta$



For each  $\nu \in (0, \infty)$ , test: is  $\nu = \mu$ ?  
Design a "potential function"  $\nu \mapsto g_n(\nu)$ :  
1) monotonically decreasing;  
2)  $g_n(\mu) \leq \frac{1}{\delta}$  with prob.  $1 - \delta$

# Lower Confidence Bound via Testing

- Observe  $Y_1, \dots, Y_n \in [0, \infty)$  with  $\mathbb{E}[Y_t] = \mu$  (unknown mean) (e.g.,  $Y_t = \tilde{r}_t(\pi)$ )
- Want: construct a LCB  $\ell_n(\delta) = \ell_n(Y_{1:n}; \delta)$  such that  $\mathbb{P}(\mu \geq \ell_n(\delta)) \geq 1 - \delta$



For each  $v \in (0, \infty)$ , test: is  $v = \mu$ ?

Design a "potential function"  $v \mapsto g_n(v)$ :

- 1) monotonically decreasing;
- 2)  $g_n(\mu) \leq \frac{1}{\delta}$  with prob.  $1 - \delta$

We can design such a potential function as a wealth function from "betting" against a synthetic stock market ☺

# Potential Function Design by Betting

# Potential Function Design by Betting

- $g_n(v) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b_t + (1 - b_t) \frac{Y_t}{v} \right)$  satisfies **the desiderata**:

# Potential Function Design by Betting

- $g_n(v) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b_t + (1 - b_t) \frac{Y_t}{v} \right)$  satisfies **the desiderata**:
  - 1) monotonically decreasing;

# Potential Function Design by Betting

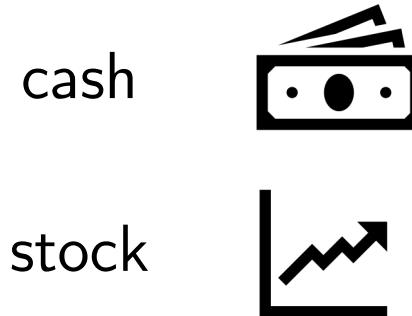
- $g_n(v) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b_t + (1 - b_t) \frac{Y_t}{v} \right)$  satisfies **the desiderata**:
  - 1) monotonically decreasing;
  - 2)  $(g_t(\mu))_{t=1}^n$  is a **martingale**  $\Rightarrow \sup_{t \geq 1} g_t(\mu) \leq \frac{1}{\delta}$  with prob.  $1 - \delta$  (by Ville's inequality)

# Potential Function Design by Betting

- $g_n(v) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b_t + (1 - b_t) \frac{Y_t}{v} \right)$  satisfies **the desiderata**:
  - 1) monotonically decreasing;
  - 2)  $(g_t(\mu))_{t=1}^n$  is a **martingale**  $\Rightarrow \sup_{t \geq 1} g_t(\mu) \leq \frac{1}{\delta}$  with prob.  $1 - \delta$  (by Ville's inequality)
- This is a “cumulative wealth” from a synthetic two-stock market setup

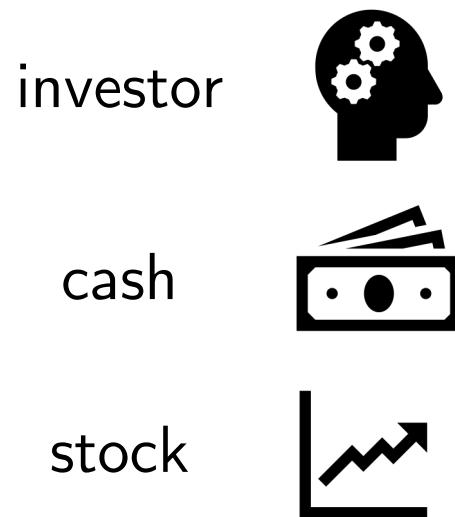
# Potential Function Design by Betting

- $g_n(v) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b_t + (1 - b_t) \frac{Y_t}{v} \right)$  satisfies **the desiderata**:
  - 1) monotonically decreasing;
  - 2)  $(g_t(\mu))_{t=1}^n$  is a **martingale**  $\Rightarrow \sup_{t \geq 1} g_t(\mu) \leq \frac{1}{\delta}$  with prob.  $1 - \delta$  (by Ville's inequality)
- This is a “cumulative wealth” from a synthetic two-stock market setup



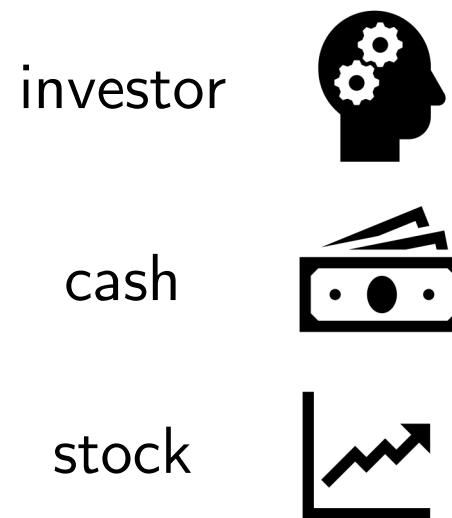
# Potential Function Design by Betting

- $g_n(v) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b_t + (1 - b_t) \frac{Y_t}{v} \right)$  satisfies **the desiderata**:
  - 1) monotonically decreasing;
  - 2)  $(g_t(\mu))_{t=1}^n$  is a **martingale**  $\Rightarrow \sup_{t \geq 1} g_t(\mu) \leq \frac{1}{\delta}$  with prob.  $1 - \delta$  (by Ville's inequality)
- This is a “cumulative wealth” from a synthetic two-stock market setup



# Potential Function Design by Betting

- $g_n(v) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b_t + (1 - b_t) \frac{Y_t}{v} \right)$  satisfies **the desiderata**:
  - 1) monotonically decreasing;
  - 2)  $(g_t(\mu))_{t=1}^n$  is a **martingale**  $\Rightarrow \sup_{t \geq 1} g_t(\mu) \leq \frac{1}{\delta}$  with prob.  $1 - \delta$  (by Ville's inequality)
- This is a “cumulative wealth” from a synthetic two-stock market setup



# Potential Function Design by Betting

- $g_n(v) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b_t + (1 - b_t) \frac{Y_t}{v} \right)$  satisfies **the desiderata**:
  - 1) monotonically decreasing;
  - 2)  $(g_t(\mu))_{t=1}^n$  is a **martingale**  $\Rightarrow \sup_{t \geq 1} g_t(\mu) \leq \frac{1}{\delta}$  with prob.  $1 - \delta$  (by Ville's inequality)
- This is a “cumulative wealth” from a synthetic two-stock market setup

investor		$b_t \in [0,1]$
cash		$b_t \times \text{wealth}_{t-1}$
stock		$(1 - b_t) \times \text{wealth}_{t-1}$

# Potential Function Design by Betting

- $g_n(v) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b_t + (1 - b_t) \frac{Y_t}{v} \right)$  satisfies **the desiderata**:
  - 1) monotonically decreasing;
  - 2)  $(g_t(\mu))_{t=1}^n$  is a **martingale**  $\Rightarrow \sup_{t \geq 1} g_t(\mu) \leq \frac{1}{\delta}$  with prob.  $1 - \delta$  (by Ville's inequality)
- This is a “cumulative wealth” from a synthetic two-stock market setup

investor		$b_t \in [0,1]$	$\frac{\text{price}_t}{\text{price}_{t-1}}$
cash		$b_t \times \text{wealth}_{t-1}$	
stock		$(1 - b_t) \times \text{wealth}_{t-1}$	

# Potential Function Design by Betting

- $g_n(v) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b_t + (1 - b_t) \frac{Y_t}{v} \right)$  satisfies **the desiderata**:
  - 1) monotonically decreasing;
  - 2)  $(g_t(\mu))_{t=1}^n$  is a **martingale**  $\Rightarrow \sup_{t \geq 1} g_t(\mu) \leq \frac{1}{\delta}$  with prob.  $1 - \delta$  (by Ville's inequality)
- This is a “cumulative wealth” from a synthetic two-stock market setup

investor		$b_t \in [0,1]$	$\frac{\text{price}_t}{\text{price}_{t-1}}$
cash		$b_t \times \text{wealth}_{t-1}$	1
stock		$(1 - b_t) \times \text{wealth}_{t-1}$	$\frac{Y_t}{v}$

# Potential Function Design by Betting

- $g_n(\nu) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b_t + (1 - b_t) \frac{Y_t}{\nu} \right)$  satisfies **the desiderata**:
  - 1) monotonically decreasing;
  - 2)  $(g_t(\mu))_{t=1}^n$  is a **martingale**  $\Rightarrow \sup_{t \geq 1} g_t(\mu) \leq \frac{1}{\delta}$  with prob.  $1 - \delta$  (by Ville's inequality)
- This is a “cumulative wealth” from a synthetic two-stock market setup

investor		$b_t \in [0,1]$	$\frac{\text{price}_t}{\text{price}_{t-1}}$	
cash		$b_t \times \text{wealth}_{t-1}$	$1$	
stock		$(1 - b_t) \times \text{wealth}_{t-1}$	$\frac{Y_t}{\nu}$	$\left\{ \frac{\text{wealth}_t(\nu)}{\text{wealth}_{t-1}(\nu)} = b_t + (1 - b_t) \frac{Y_t}{\nu} \right.$

# Potential Function Design by Betting

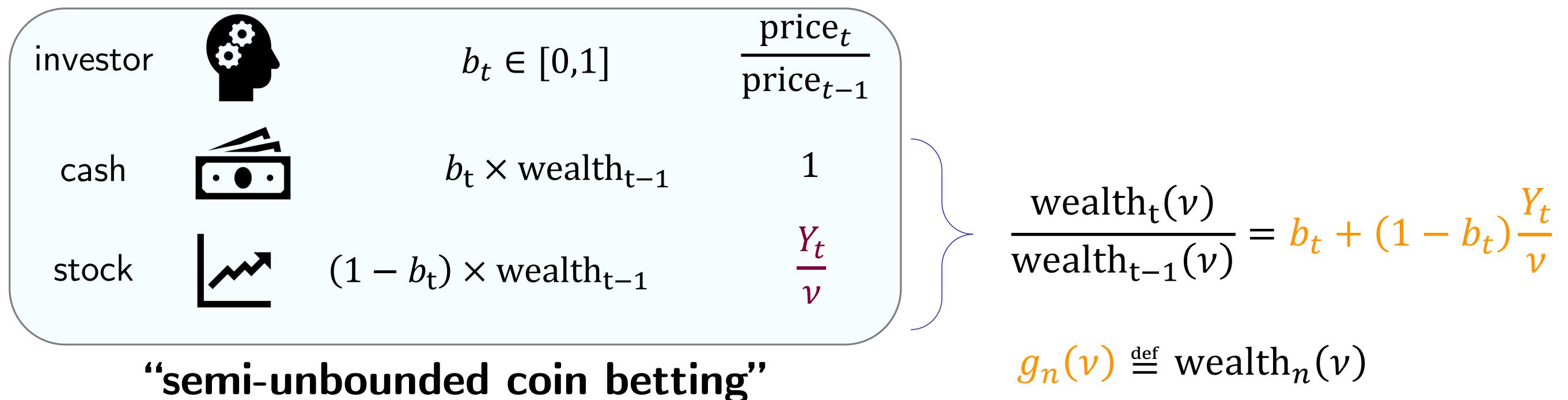
- $g_n(v) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b_t + (1 - b_t) \frac{Y_t}{v} \right)$  satisfies **the desiderata**:
  - 1) monotonically decreasing;
  - 2)  $(g_t(\mu))_{t=1}^n$  is a **martingale**  $\Rightarrow \sup_{t \geq 1} g_t(\mu) \leq \frac{1}{\delta}$  with prob.  $1 - \delta$  (by Ville's inequality)
- This is a “cumulative wealth” from a synthetic two-stock market setup

investor		$b_t \in [0,1]$	$\frac{\text{price}_t}{\text{price}_{t-1}}$	
cash		$b_t \times \text{wealth}_{t-1}$	$1$	$\frac{Y_t}{v}$
stock		$(1 - b_t) \times \text{wealth}_{t-1}$		$\left. \frac{\text{wealth}_t(v)}{\text{wealth}_{t-1}(v)} = b_t + (1 - b_t) \frac{Y_t}{v} \right\}$

$g_n(v) \stackrel{\text{def}}{=} \text{wealth}_n(v)$

# Potential Function Design by Betting

- $g_n(\nu) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b_t + (1 - b_t) \frac{Y_t}{\nu} \right)$  satisfies **the desiderata**:
  - 1) monotonically decreasing;
  - 2)  $(g_t(\mu))_{t=1}^n$  is a **martingale**  $\Rightarrow \sup_{t \geq 1} g_t(\mu) \leq \frac{1}{\delta}$  with prob.  $1 - \delta$  (by Ville's inequality)
- This is a “cumulative wealth” from a synthetic two-stock market setup



# Potential Function Design by Betting

- $g_n(v) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b_t + (1 - b_t) \frac{Y_t}{v} \right)$  satisfies **the desiderata**:
  - 1) monotonically decreasing;
  - 2)  $(g_t(\mu))_{t=1}^n$  is a **martingale**  $\Rightarrow \sup_{t \geq 1} g_t(\mu) \leq \frac{1}{\delta}$  with prob.  $1 - \delta$  (by Ville's inequality)
- This is a “cumulative wealth” from a synthetic two-stock market setup

## Historical Bits

- [W-S&R24], [OJ24] studied betting framework for **bounded variables**.
- [W-S+22] proposed this setup for nonnegative variables **in disguise**.
- [RyuB24] proposed **two-stock market formulation**.
- [RyuW24] proposed  **$k$ -stock market variants for bounded vectors**.

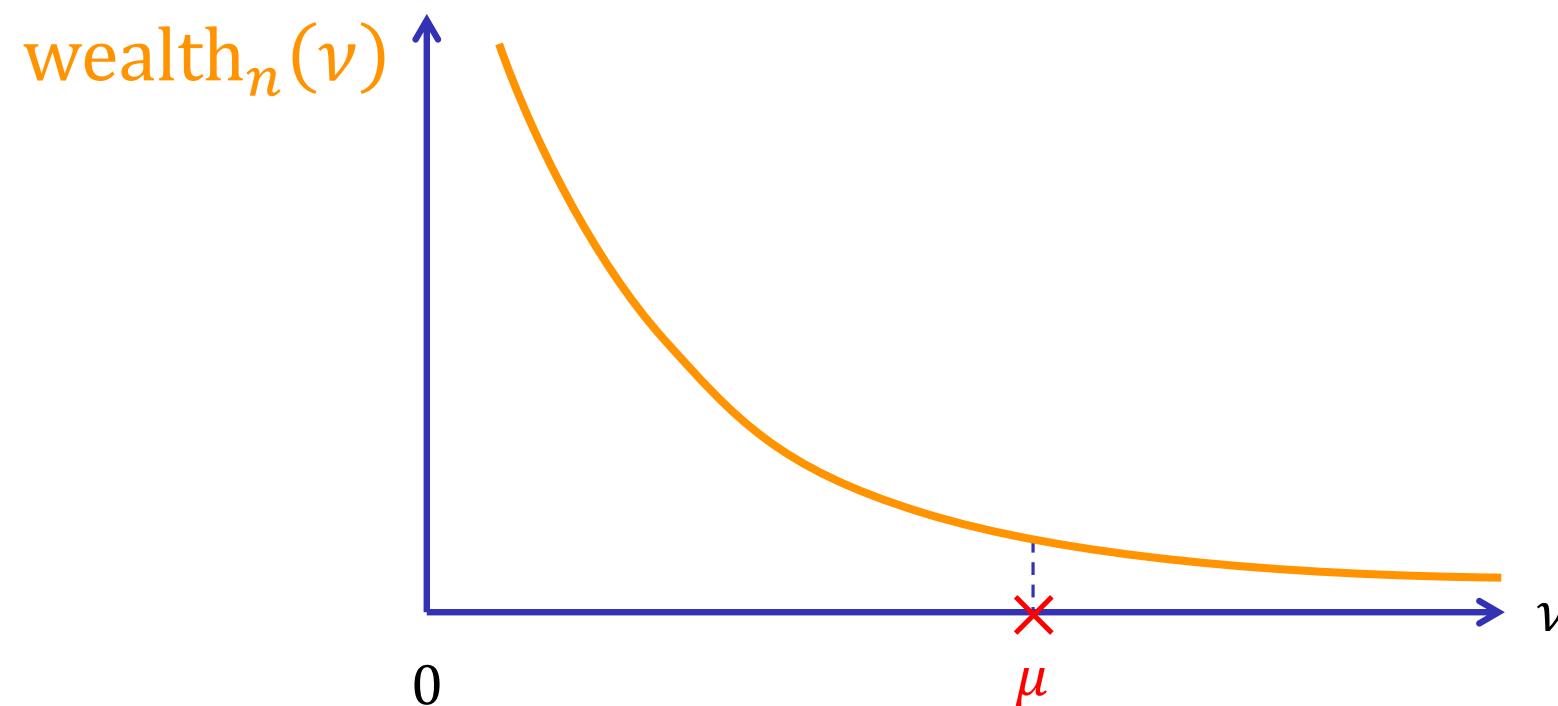
# Template: LCB via Betting

# Template: LCB via Betting

- Choose a (causal) betting strategy   $(b_t(Y_{1:t-1}))_{t \geq 1}$

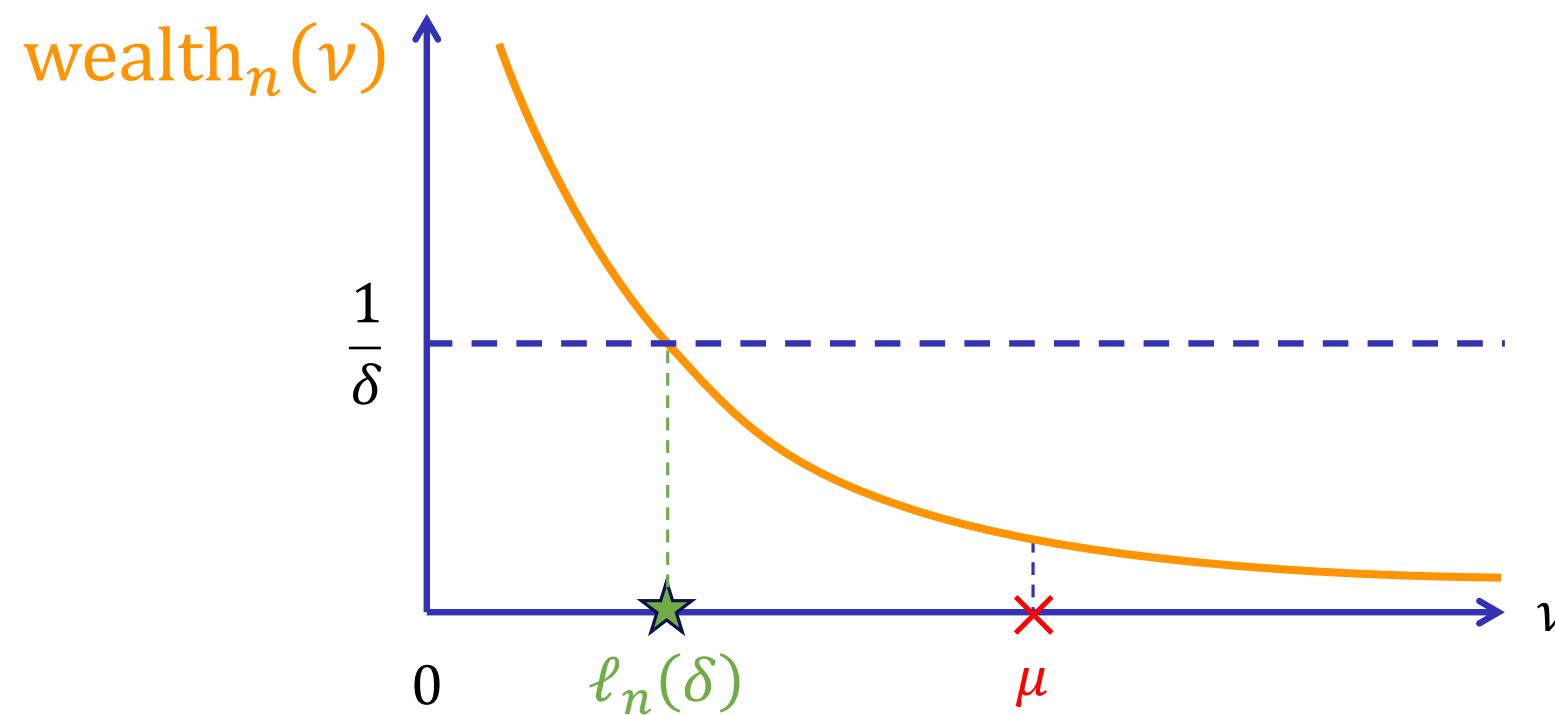
# Template: LCB via Betting

- Choose a **(causal) betting strategy**   $(b_t(Y_{1:t-1}))_{t \geq 1}$
- Compute the induced wealth function  $v \mapsto \text{wealth}_n(v)$  from the two-stock market



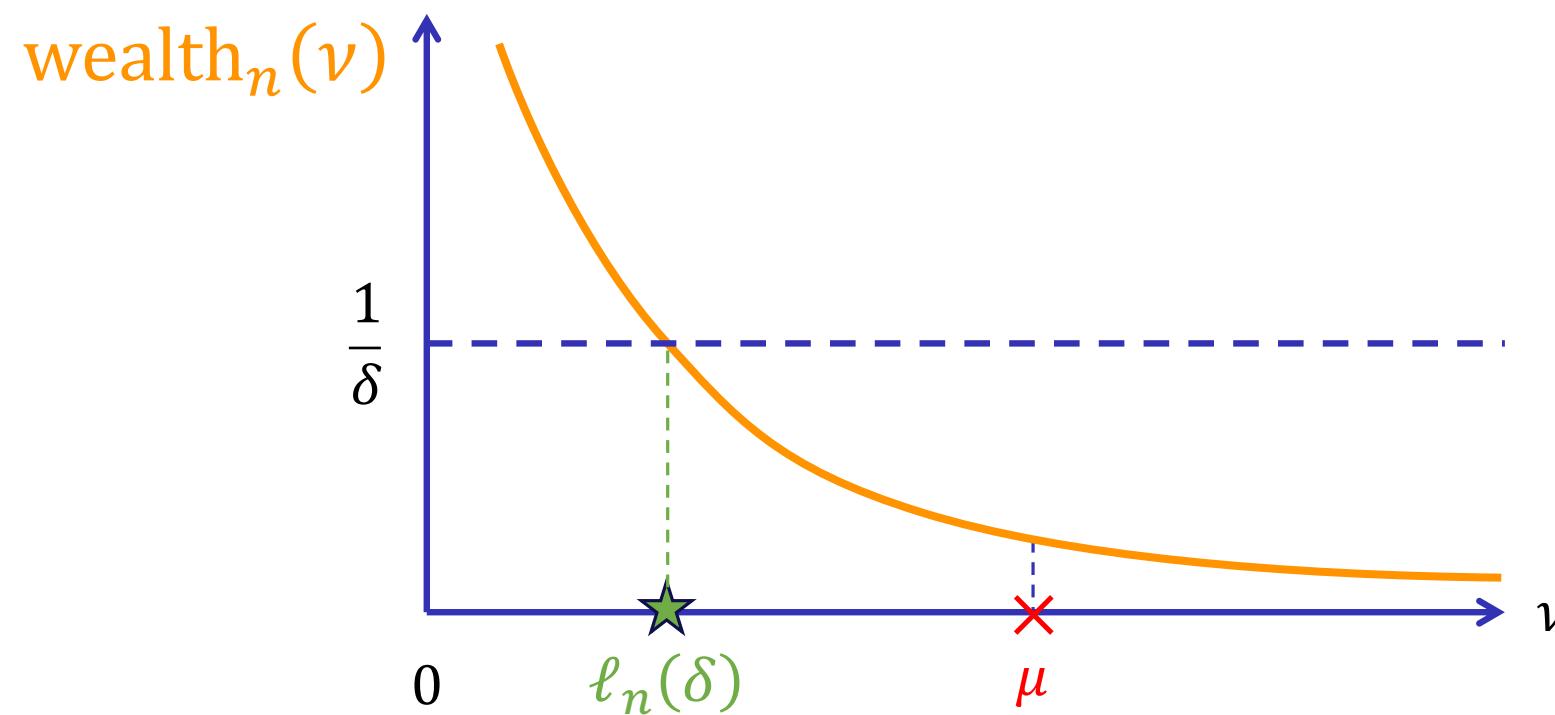
# Template: LCB via Betting

- Choose a (causal) betting strategy   $(b_t(Y_{1:t-1}))_{t \geq 1}$
- Compute the induced wealth function  $v \mapsto \text{wealth}_n(v)$  from the two-stock market
- Construct LCB via testing with  $v \mapsto \text{wealth}_n(v)$



# Template: LCB via Betting

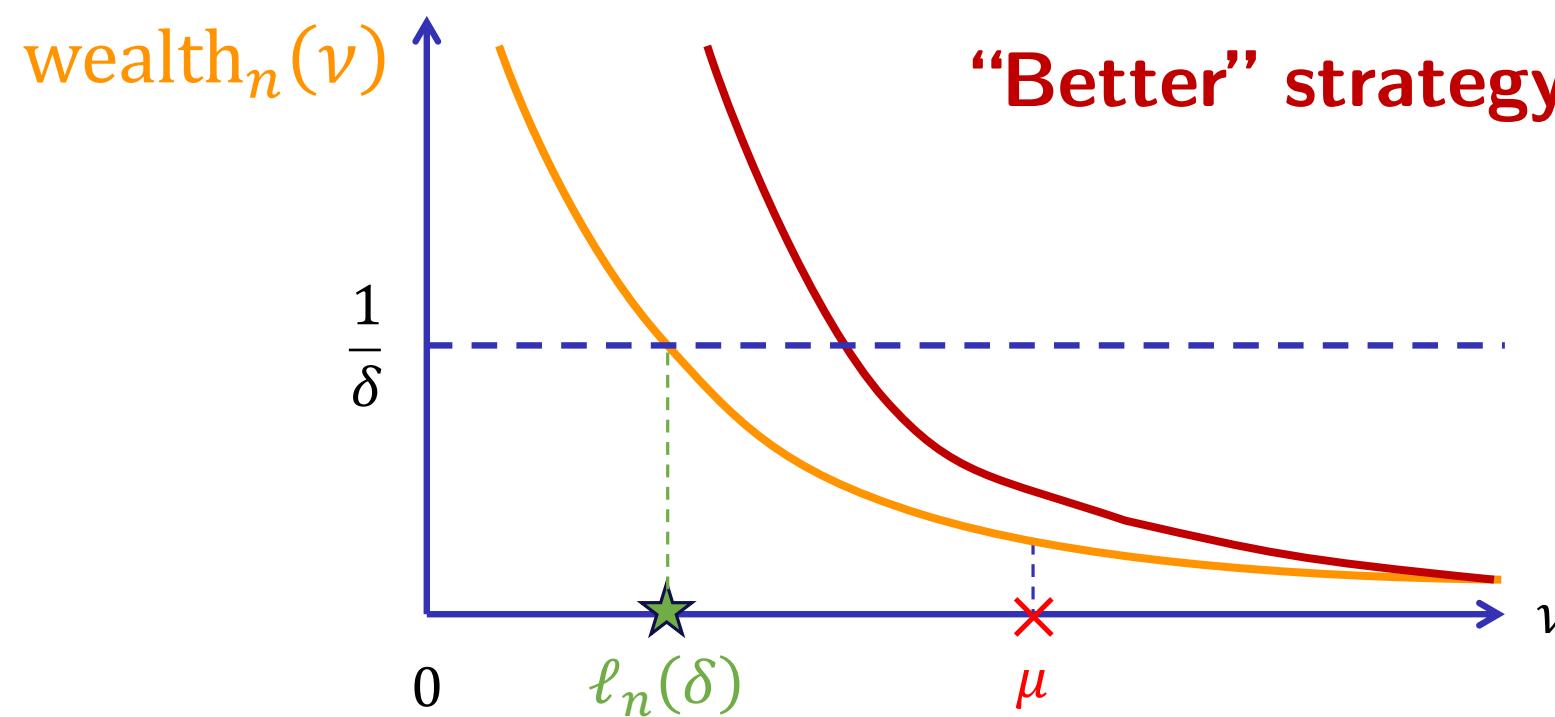
- Choose a (causal) betting strategy   $(b_t(Y_{1:t-1}))_{t \geq 1}$
- Compute the induced wealth function  $v \mapsto \text{wealth}_n(v)$  from the two-stock market
- Construct LCB via testing with  $v \mapsto \text{wealth}_n(v)$



**Root finding** can be done in  $O\left(\log \frac{1}{\epsilon}\right)$   
(modulo function evaluation)

# Template: LCB via Betting

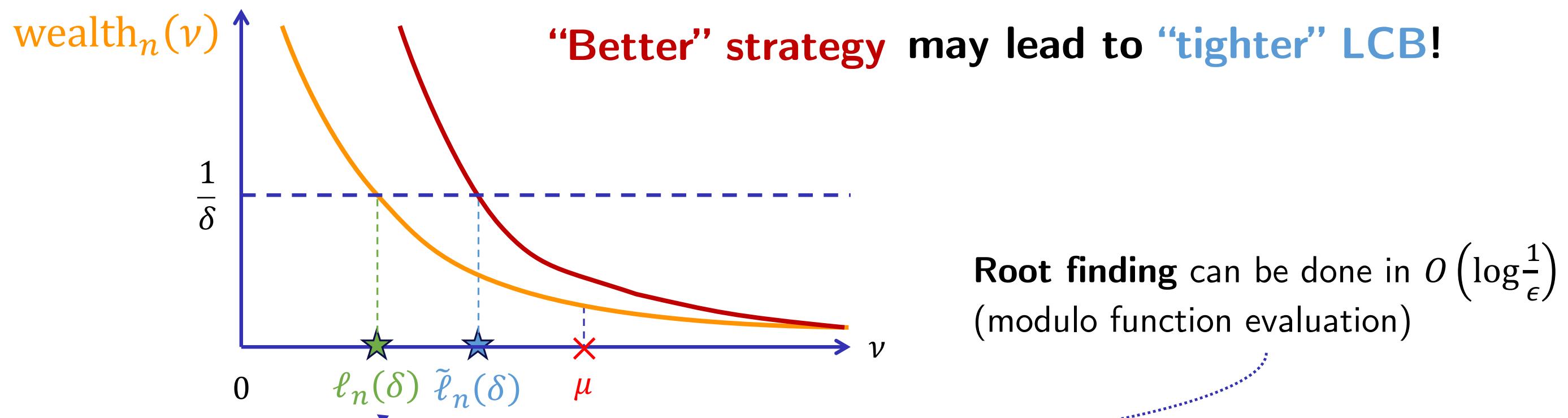
- Choose a (causal) betting strategy   $(b_t(Y_{1:t-1}))_{t \geq 1}$
- Compute the induced wealth function  $v \mapsto \text{wealth}_n(v)$  from the two-stock market
- Construct LCB via testing with  $v \mapsto \text{wealth}_n(v)$



**Root finding** can be done in  $O\left(\log \frac{1}{\epsilon}\right)$   
(modulo function evaluation)

# Template: LCB via Betting

- Choose a (causal) betting strategy   $(b_t(Y_{1:t-1}))_{t \geq 1}$
- Compute the induced wealth function  $v \mapsto \text{wealth}_n(v)$  from the two-stock market
- Construct LCB via testing with  $v \mapsto \text{wealth}_n(v)$



# What Strategy To Use?

# What Strategy To Use?

- **Simple strategy:** Constant Rebalanced Portfolio (CRP) ( $b, 1-b$ ) for  $b \in (0,1)$

# What Strategy To Use?

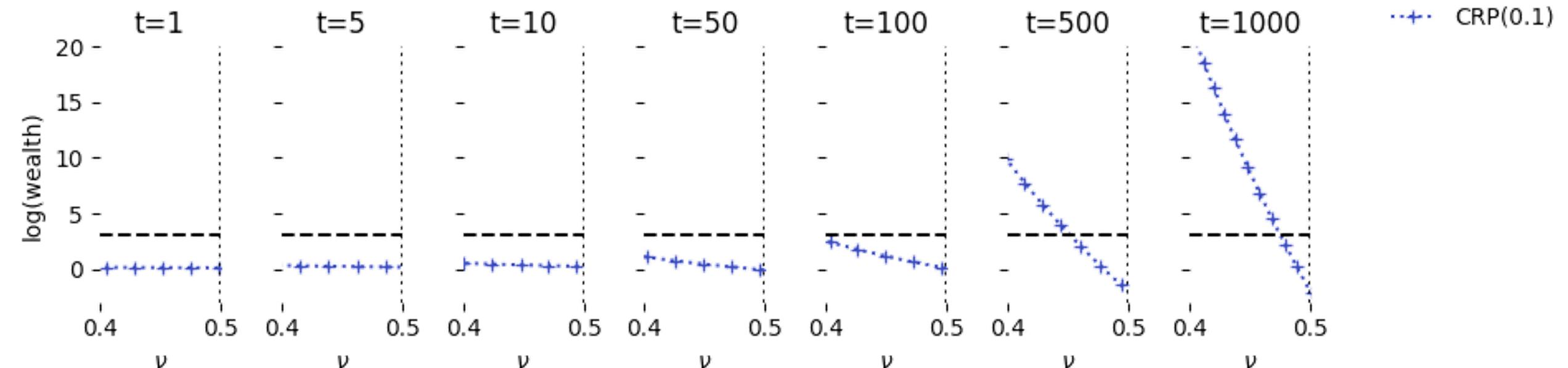
$$\text{wealth}_t^b(\nu) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b + (1-b) \frac{Y_t}{\nu} \right)$$

- **Simple strategy:** Constant Rebalanced Portfolio (CRP) ( $b, 1-b$ ) for  $b \in (0,1)$

# What Strategy To Use?

$$\text{wealth}_t^b(\nu) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b + (1-b) \frac{Y_t}{\nu} \right)$$

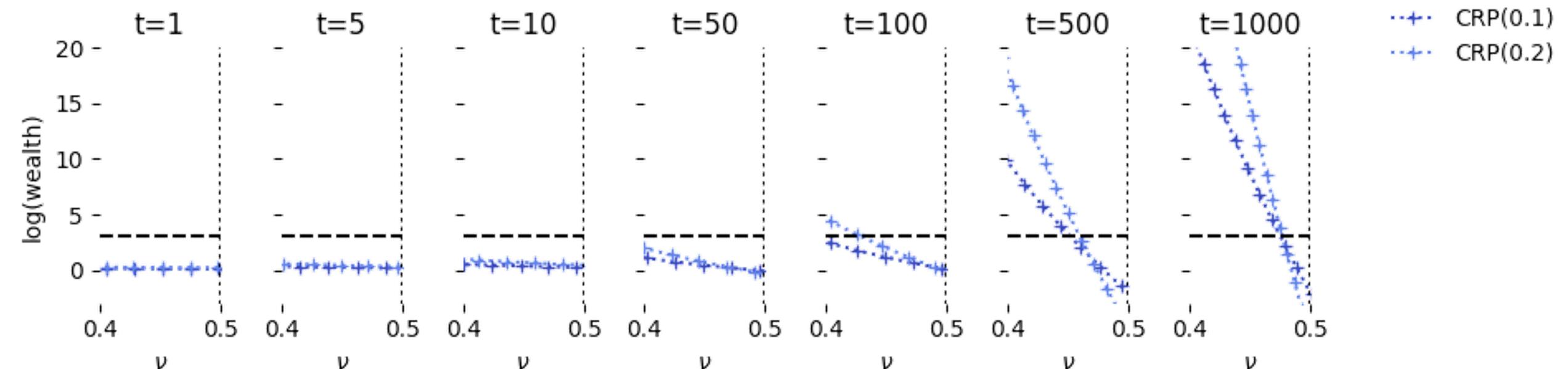
- **Simple strategy:** Constant Rebalanced Portfolio (CRP) ( $b, 1-b$ ) for  $b \in (0,1)$



# What Strategy To Use?

$$\text{wealth}_t^b(\nu) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b + (1-b) \frac{Y_t}{\nu} \right)$$

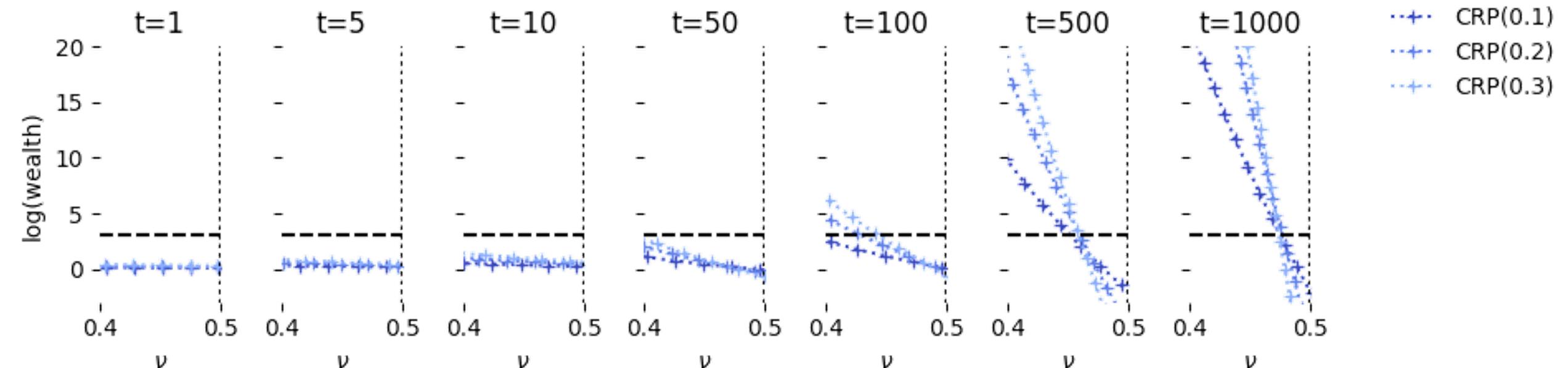
- **Simple strategy:** Constant Rebalanced Portfolio (CRP) ( $b, 1-b$ ) for  $b \in (0,1)$



# What Strategy To Use?

$$\text{wealth}_t^b(\nu) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b + (1-b) \frac{Y_t}{\nu} \right)$$

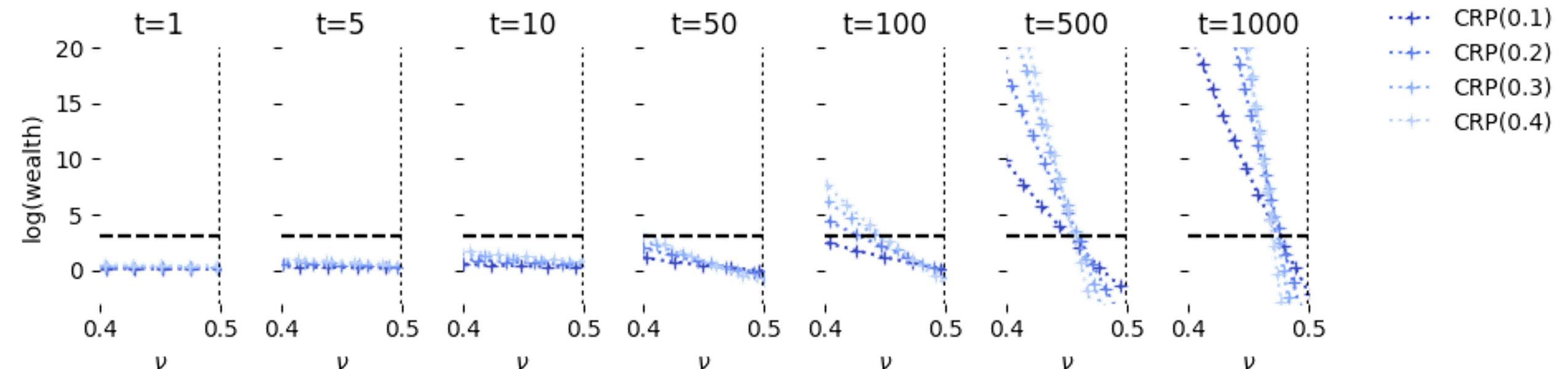
- **Simple strategy:** Constant Rebalanced Portfolio (CRP) ( $b$ ,  $1-b$ ) for  $b \in (0,1)$



# What Strategy To Use?

$$\text{wealth}_t^b(\nu) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b + (1-b) \frac{Y_t}{\nu} \right)$$

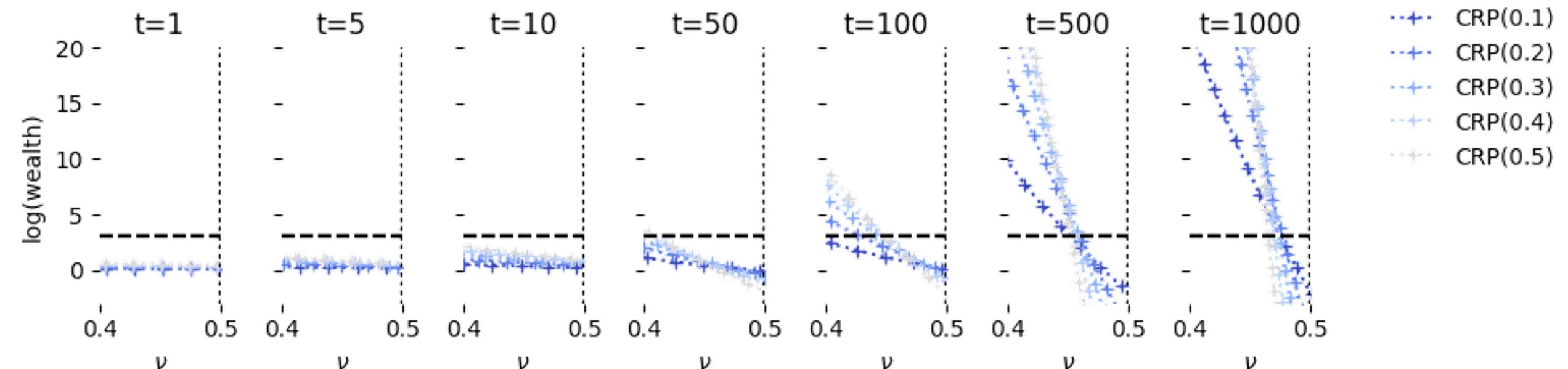
- **Simple strategy:** Constant Rebalanced Portfolio (CRP) ( $b$ ,  $1-b$ ) for  $b \in (0,1)$



# What Strategy To Use?

$$\text{wealth}_t^b(\nu) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b + (1-b) \frac{Y_t}{\nu} \right)$$

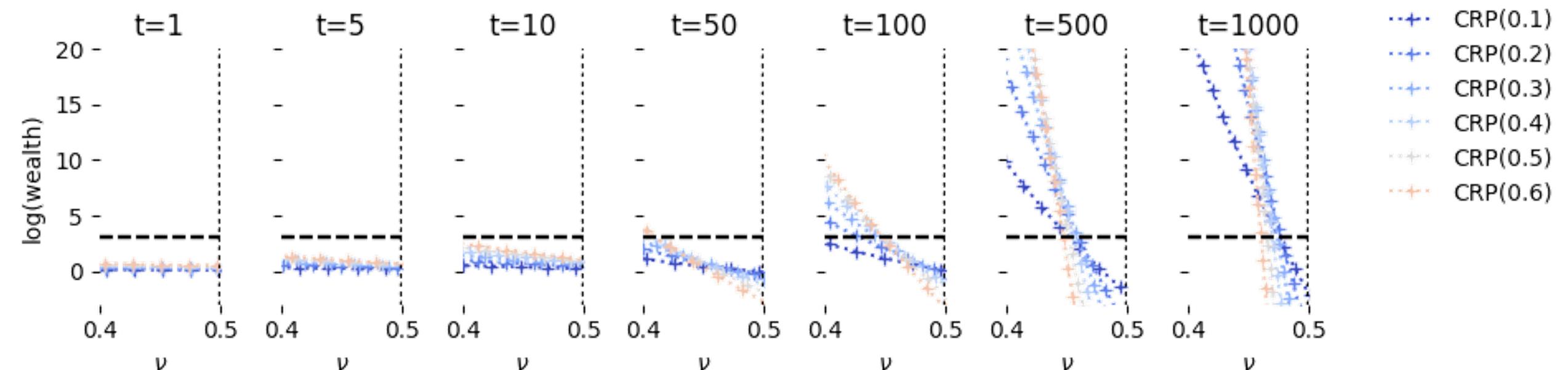
- **Simple strategy:** Constant Rebalanced Portfolio (CRP) ( $b, 1-b$ ) for  $b \in (0,1)$



# What Strategy To Use?

$$\text{wealth}_t^b(\nu) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b + (1-b) \frac{Y_t}{\nu} \right)$$

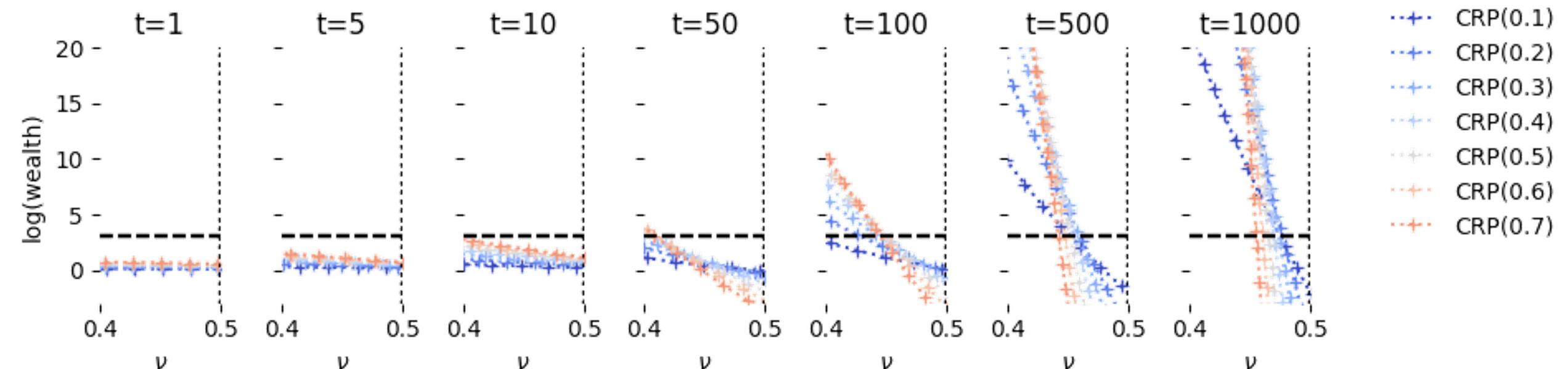
- **Simple strategy:** Constant Rebalanced Portfolio (CRP) ( $b, 1-b$ ) for  $b \in (0,1)$



# What Strategy To Use?

$$\text{wealth}_t^b(\nu) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b + (1-b) \frac{Y_t}{\nu} \right)$$

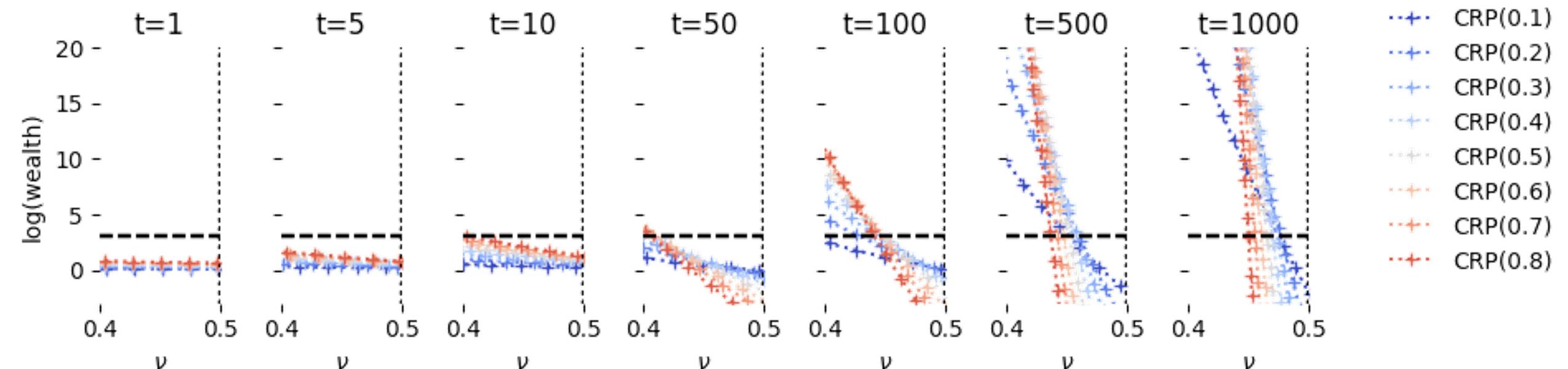
- **Simple strategy:** Constant Rebalanced Portfolio (CRP) ( $b, 1-b$ ) for  $b \in (0,1)$



# What Strategy To Use?

$$\text{wealth}_t^b(\nu) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b + (1-b) \frac{Y_t}{\nu} \right)$$

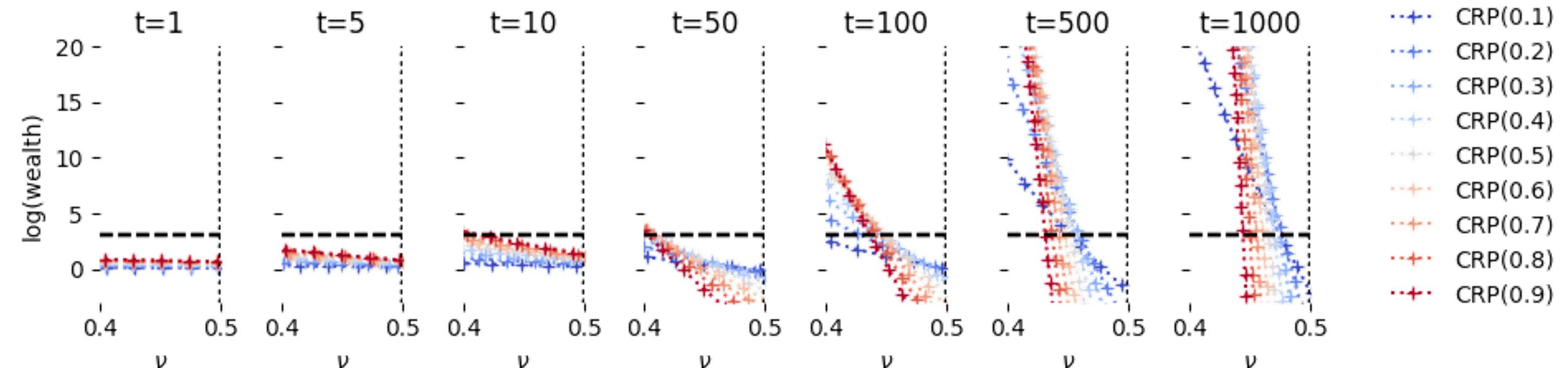
- **Simple strategy:** Constant Rebalanced Portfolio (CRP) ( $b, 1-b$ ) for  $b \in (0,1)$



# What Strategy To Use?

$$\text{wealth}_t^b(\nu) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b + (1-b) \frac{Y_t}{\nu} \right)$$

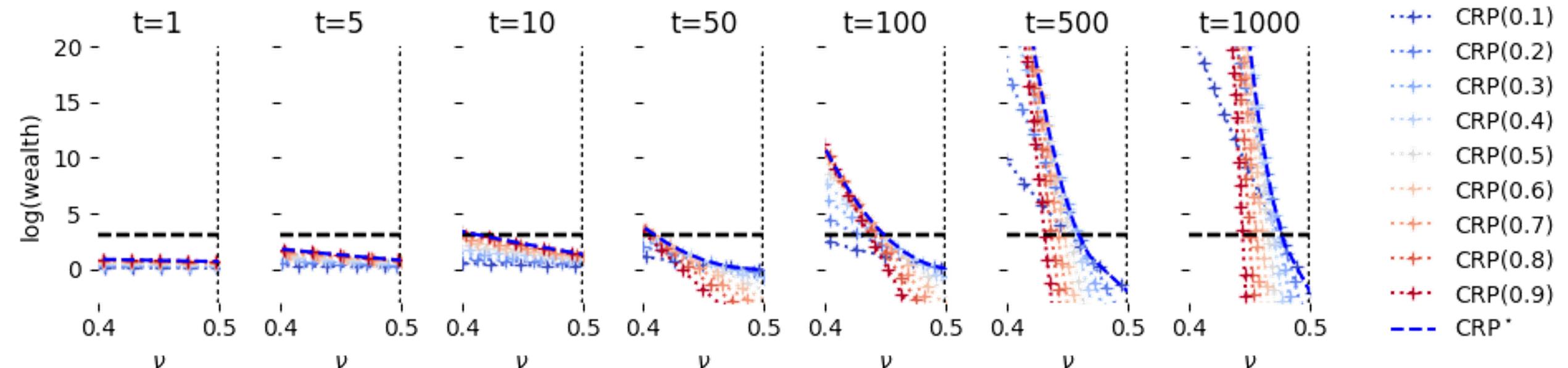
- **Simple strategy:** Constant Rebalanced Portfolio (CRP) ( $b$ ,  $1-b$ ) for  $b \in (0,1)$



# What Strategy To Use?

$$\text{wealth}_t^b(\nu) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b + (1-b) \frac{Y_t}{\nu} \right)$$

- **Simple strategy:** Constant Rebalanced Portfolio (CRP) ( $b, 1-b$ ) for  $b \in (0,1)$

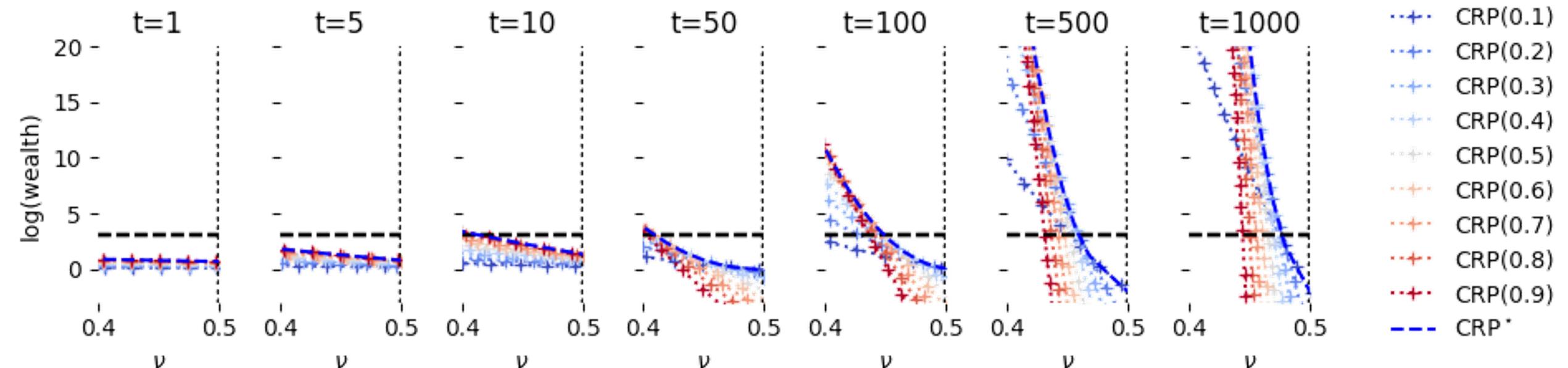


Q. Can we attain the maximum wealth of CRPs **in hindsight**?  $\nu \mapsto \max_{b \in (0,1)} \text{wealth}_t^b(\nu)$

# What Strategy To Use?

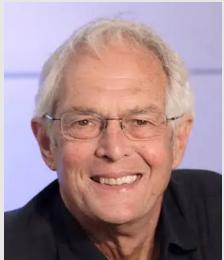
$$\text{wealth}_t^b(\nu) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b + (1-b) \frac{Y_t}{\nu} \right)$$

- **Simple strategy:** Constant Rebalanced Portfolio (CRP) ( $b, 1-b$ ) for  $b \in (0,1)$



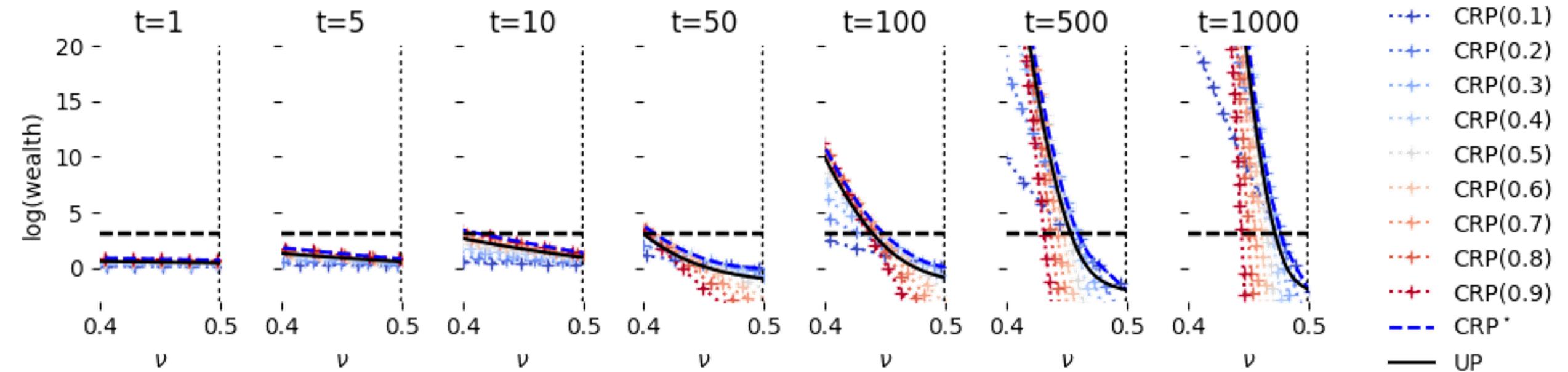
Q. Can we attain the maximum wealth of CRPs **in hindsight**?  $\nu \mapsto \max_{b \in (0,1)} \text{wealth}_t^b(\nu)$

# What Strategy To Use?



$$! \quad \text{wealth}_t^b(\nu) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b + (1-b) \frac{Y_t}{\nu} \right)$$

- **Simple strategy:** Constant Rebalanced Portfolio (CRP) ( $b, 1-b$ ) for  $b \in (0,1)$



Q. Can we attain the maximum wealth of CRPs **in hindsight**?  $\nu \mapsto \max_{b \in (0,1)} \text{wealth}_t^b(\nu)$

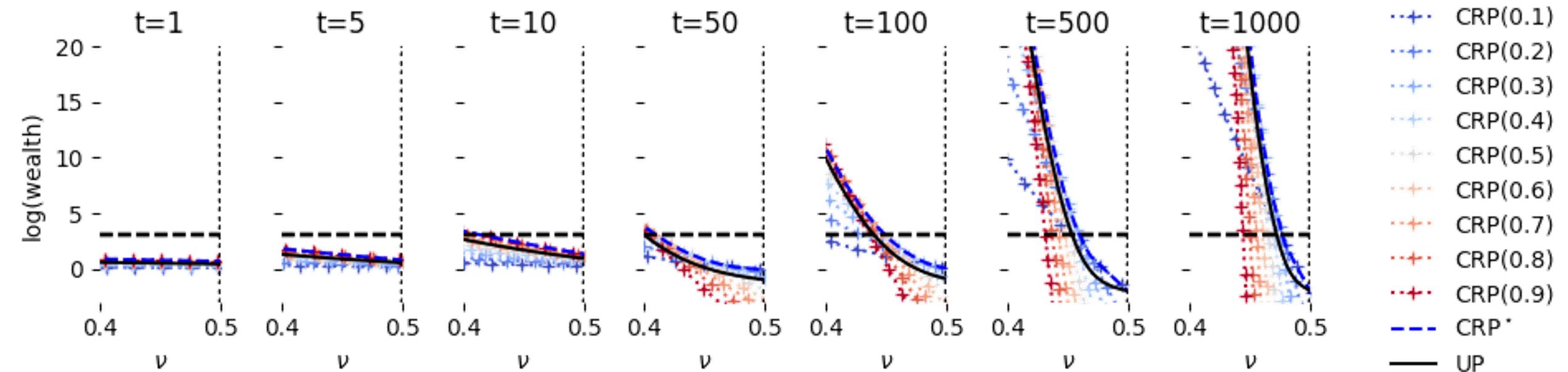
A. Yes, by [Cover91]'s **universal portfolio (UP)**, defined as **mixture wealth**

# What Strategy To Use?



$$! \quad \text{wealth}_t^b(\nu) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b + (1-b) \frac{Y_t}{\nu} \right)$$

- **Simple strategy:** Constant Rebalanced Portfolio (CRP) ( $b, 1-b$ ) for  $b \in (0,1)$



Q. Can we attain the maximum wealth of CRPs **in hindsight**?  $\nu \mapsto \max_{b \in (0,1)} \text{wealth}_t^b(\nu)$

A. Yes, by [Cover91]'s **universal portfolio (UP)**, defined as **mixture wealth**

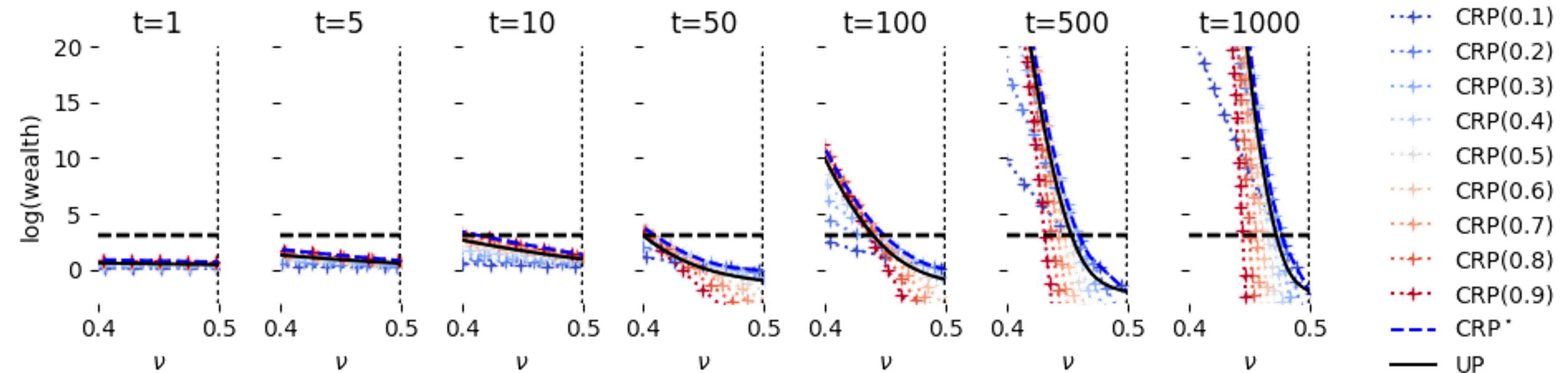
$$\text{wealth}_t^{\text{UP}}(\nu) \stackrel{\text{def}}{=} \int_0^1 \text{wealth}_t^b(\nu) \pi(b) db$$

# What Strategy To Use?



$$! \quad \text{wealth}_t^b(\nu) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b + (1-b) \frac{Y_t}{\nu} \right)$$

- **Simple strategy:** Constant Rebalanced Portfolio (CRP) ( $b, 1-b$ ) for  $b \in (0,1)$



**Q.** Can we attain the maximum wealth of CRPs in hindsight?  $\nu \mapsto \max_{b \in (0,1)} \text{wealth}_t^b(\nu)$

**A.** Yes, by [Cover91]'s universal portfolio (UP), defined as mixture wealth

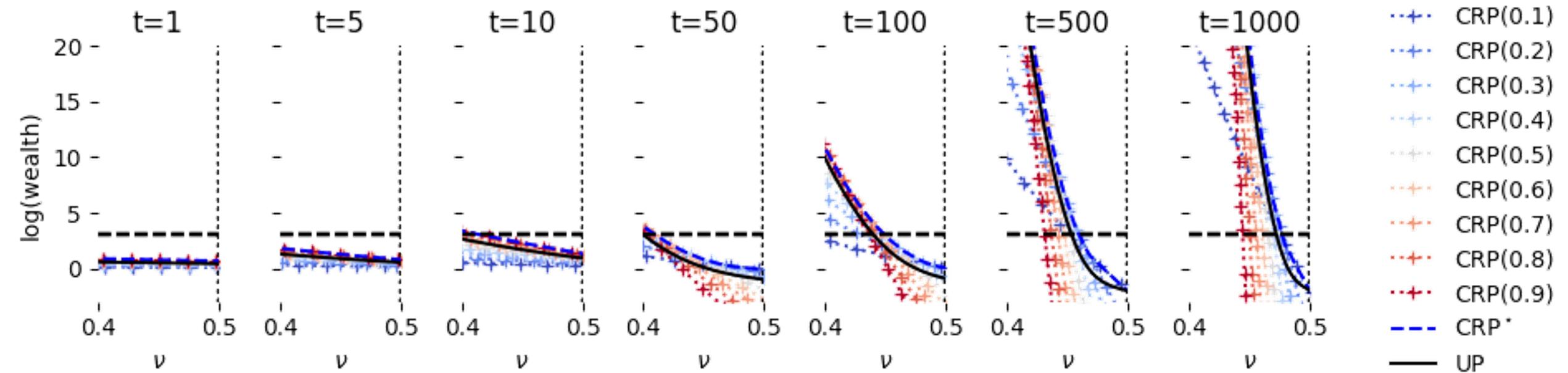
$$\text{wealth}_t^{\text{UP}}(\nu) \stackrel{\text{def}}{=} \int_0^1 \text{wealth}_t^b(\nu) \pi(b) db \stackrel{[\text{Cover96}]}{\geq} \frac{1}{\sqrt{\pi(t+1)}} \max_{b \in (0,1)} \text{wealth}_t^b(\nu)$$

# What Strategy To Use?



$$! \quad \text{wealth}_t^b(\nu) \stackrel{\text{def}}{=} \prod_{t=1}^n \left( b + (1-b) \frac{Y_t}{\nu} \right)$$

- **Simple strategy:** Constant Rebalanced Portfolio (CRP) ( $b, 1-b$ ) for  $b \in (0,1)$



**Q.** Can we attain the maximum wealth of CRPs in hindsight?  $\nu \mapsto \max_{b \in (0,1)} \text{wealth}_t^b(\nu)$

**A.** Yes, by [Cover91]'s universal portfolio (UP), defined as mixture wealth

$$\text{wealth}_t^{\text{UP}}(\nu) \stackrel{\text{def}}{=} \int_0^1 \text{wealth}_t^b(\nu) \pi(b) db \stackrel{[\text{Cover96}]}{\geq} \frac{1}{\sqrt{\pi(t+1)}} \max_{b \in (0,1)} \text{wealth}_t^b(\nu) \quad [\text{OJ24}] \quad [\text{RyuB24}]$$

[RyuW24]

# UP-LCB

- **UP-LCB** = semi-unbounded coin betting + Cover's universal portfolio

# UP-LCB

- **UP-LCB** = semi-unbounded coin betting + Cover's universal portfolio

$$\begin{aligned} \text{wealth}_n^{\text{UP}}(\nu) &\stackrel{\text{def}}{=} \int_0^1 \text{wealth}_n^b(\nu) \pi(b) db \\ &\geq \text{wealth}_n^{\text{pCRP}^*}(\nu) \stackrel{\text{def}}{=} \frac{1}{\sqrt{\pi(n+1)}} \max_{b \in (0,1)} \text{wealth}_n^b(\nu) \end{aligned}$$

# UP-LCB

- **UP-LCB** = semi-unbounded coin betting + Cover's universal portfolio

$$\begin{aligned} \text{wealth}_n^{\text{UP}}(\nu) &\stackrel{\text{def}}{=} \int_0^1 \text{wealth}_n^b(\nu) \pi(b) db && \longrightarrow O(n^2) \\ &\geq \text{wealth}_n^{\text{pCRP}^*}(\nu) \stackrel{\text{def}}{=} \frac{1}{\sqrt{\pi(n+1)}} \max_{b \in (0,1)} \text{wealth}_n^b(\nu) && \longrightarrow O(n) \end{aligned}$$

# UP-LCB

- **UP-LCB** = semi-unbounded coin betting + Cover's universal portfolio

$$\begin{aligned} \text{wealth}_n^{\text{UP}}(\nu) &\stackrel{\text{def}}{=} \int_0^1 \text{wealth}_n^b(\nu) \pi(b) db && \longrightarrow O(n^2) \\ &\geq \text{wealth}_n^{\text{pCRP}^*}(\nu) \stackrel{\text{def}}{=} \frac{1}{\sqrt{\pi(n+1)}} \max_{b \in (0,1)} \text{wealth}_n^b(\nu) && \longrightarrow O(n) \end{aligned}$$

**Theorem** (informal) With probability  $\geq 1 - \delta$ ,

$$0 \leq \mu - \ell_n^{\text{UP}}(\delta) \leq \mu - \ell_n^{\text{pCRP}^*}(\delta) \lesssim \sqrt{\text{Var}(Y_1) \frac{1}{n} \log \frac{1}{\delta}}$$

# UP-LCB

- **UP-LCB** = semi-unbounded coin betting + Cover's universal portfolio

$$\begin{aligned} \text{wealth}_n^{\text{UP}}(\nu) &\stackrel{\text{def}}{=} \int_0^1 \text{wealth}_n^b(\nu) \pi(b) db && \longrightarrow O(n^2) \\ &\geq \text{wealth}_n^{\text{pCRP}^*}(\nu) \stackrel{\text{def}}{=} \frac{1}{\sqrt{\pi(n+1)}} \max_{b \in (0,1)} \text{wealth}_n^b(\nu) && \longrightarrow O(n) \end{aligned}$$

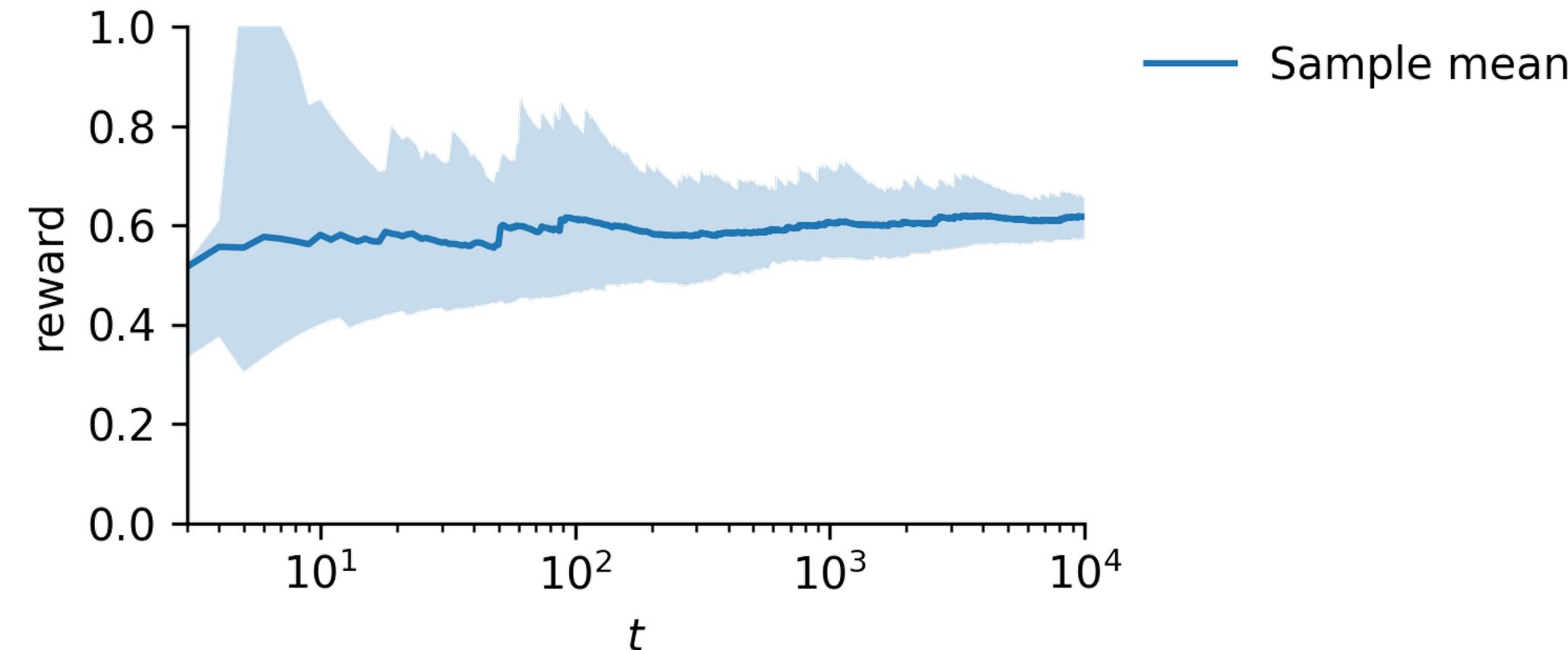
**Theorem** (informal) With probability  $\geq 1 - \delta$ ,

$$0 \leq \mu - \ell_n^{\text{UP}}(\delta) \leq \mu - \ell_n^{\text{pCRP}^*}(\delta) \lesssim \sqrt{\text{Var}(Y_1) \frac{1}{n} \log \frac{1}{\delta}}$$

The variance-adaptive selection guarantee immediately follows as a corollary!

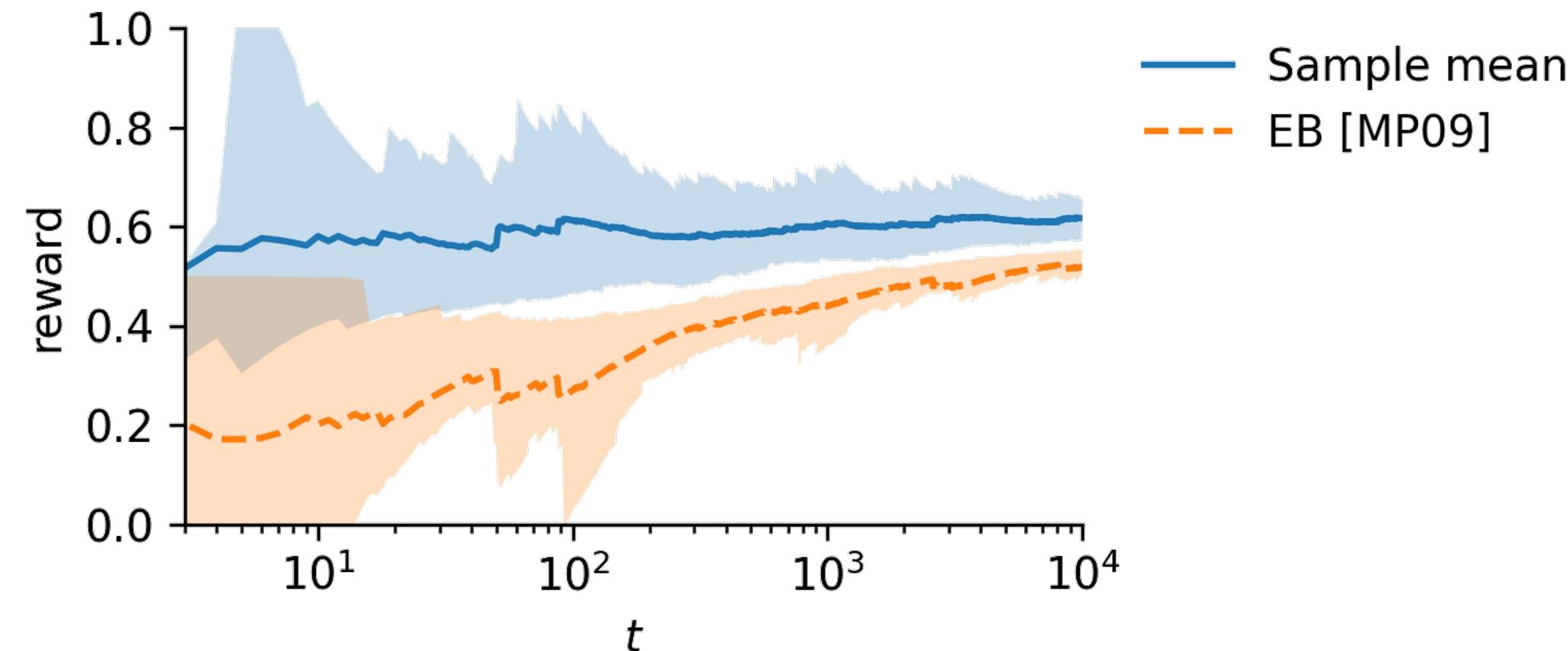
# Experiment: UP-LCB

- Synthetic reward processes with **infinite** 4th moment
- Empirical Bernstein (EB) = Plug-in empirical variance (heuristic)



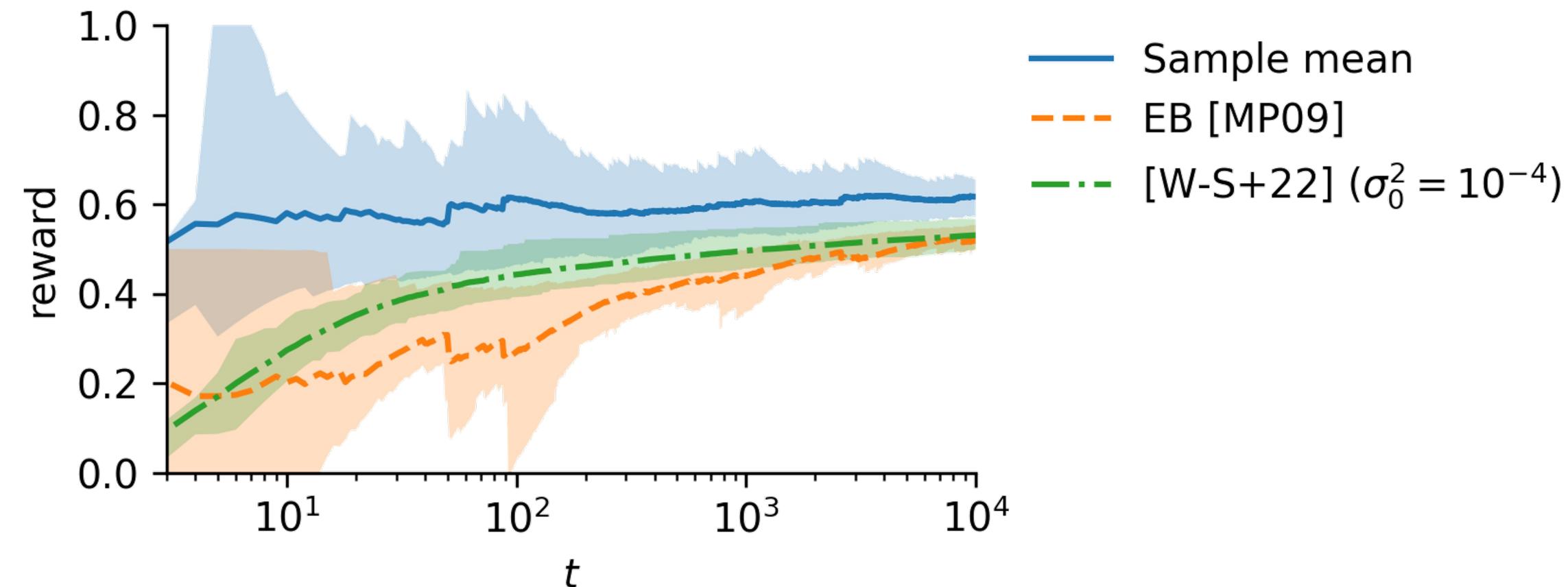
# Experiment: UP-LCB

- Synthetic reward processes with **infinite** 4th moment
- Empirical Bernstein (EB) = Plug-in empirical variance (heuristic)



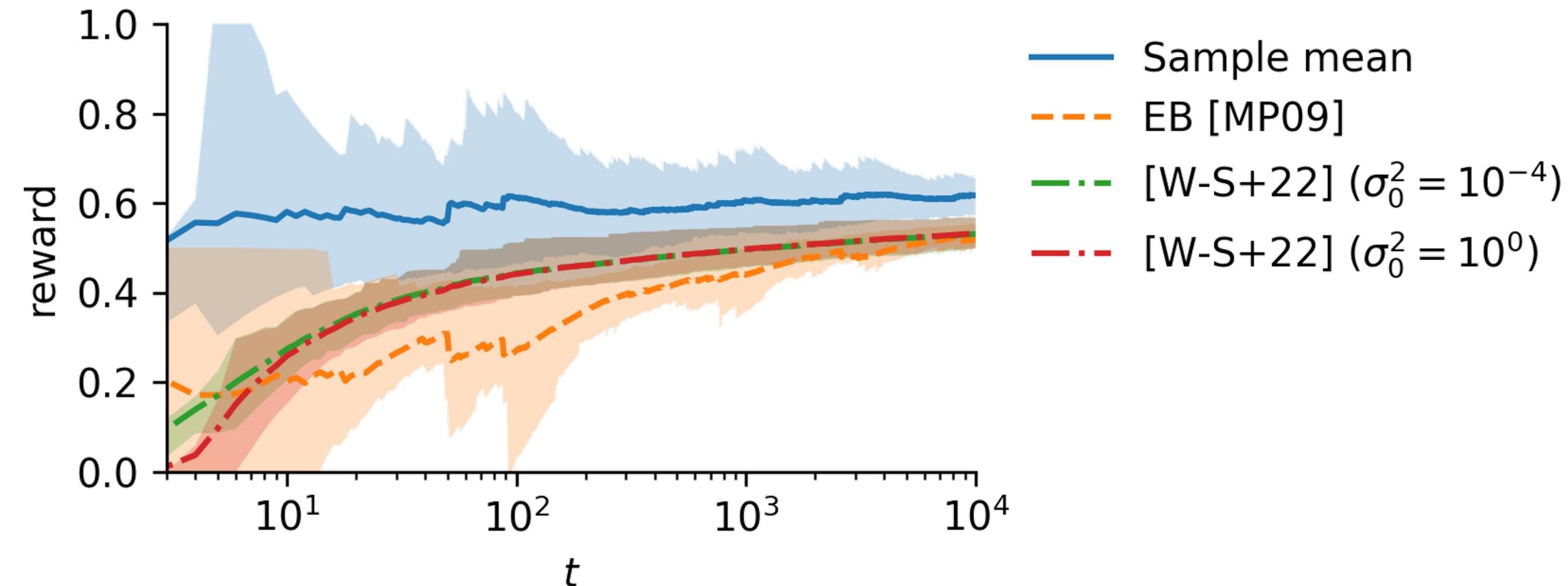
# Experiment: UP-LCB

- Synthetic reward processes with **infinite** 4th moment
- Empirical Bernstein (EB) = Plug-in empirical variance (heuristic)



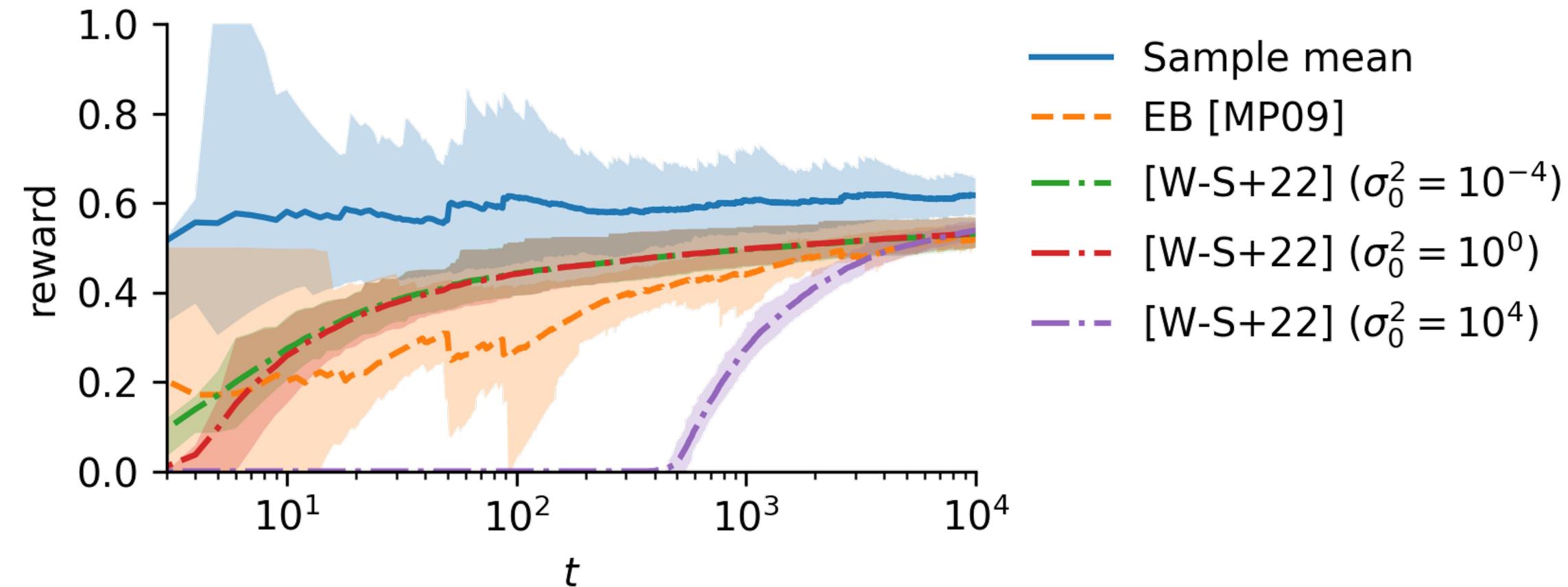
# Experiment: UP-LCB

- Synthetic reward processes with **infinite** 4th moment
- Empirical Bernstein (EB) = Plug-in empirical variance (heuristic)



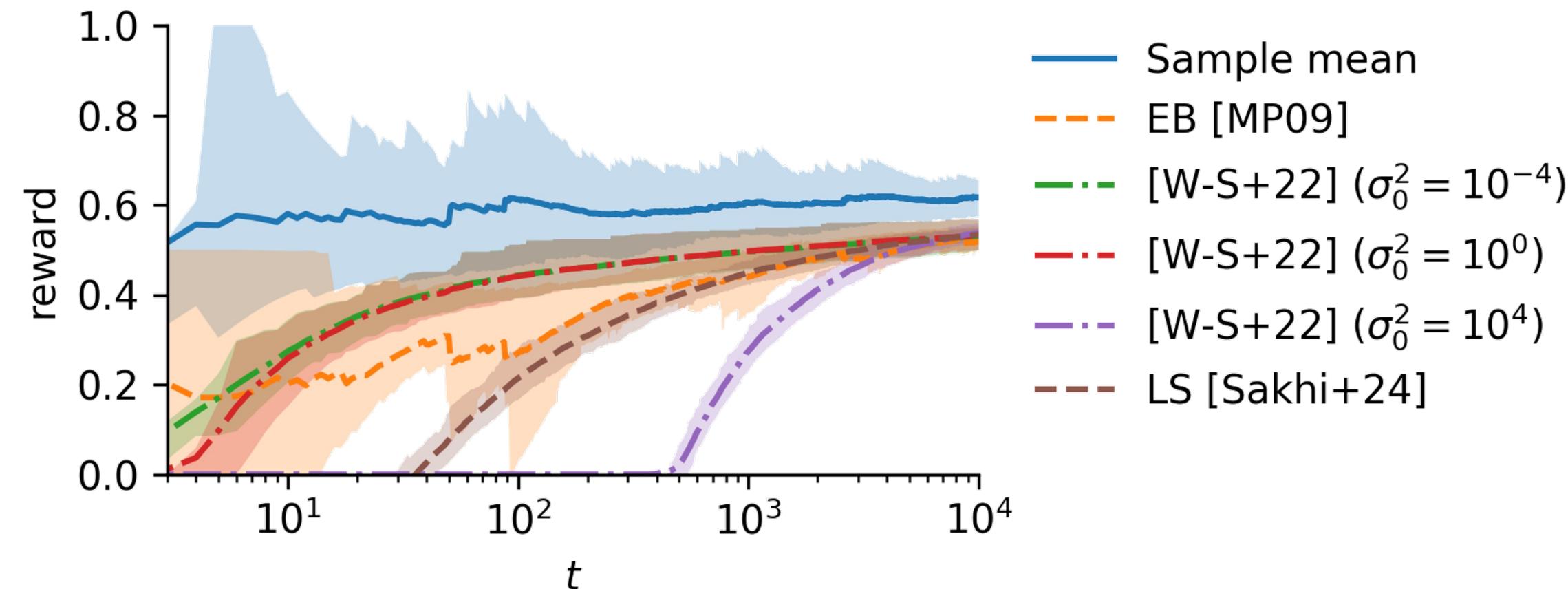
# Experiment: UP-LCB

- Synthetic reward processes with **infinite** 4th moment
- Empirical Bernstein (EB) = Plug-in empirical variance (heuristic)



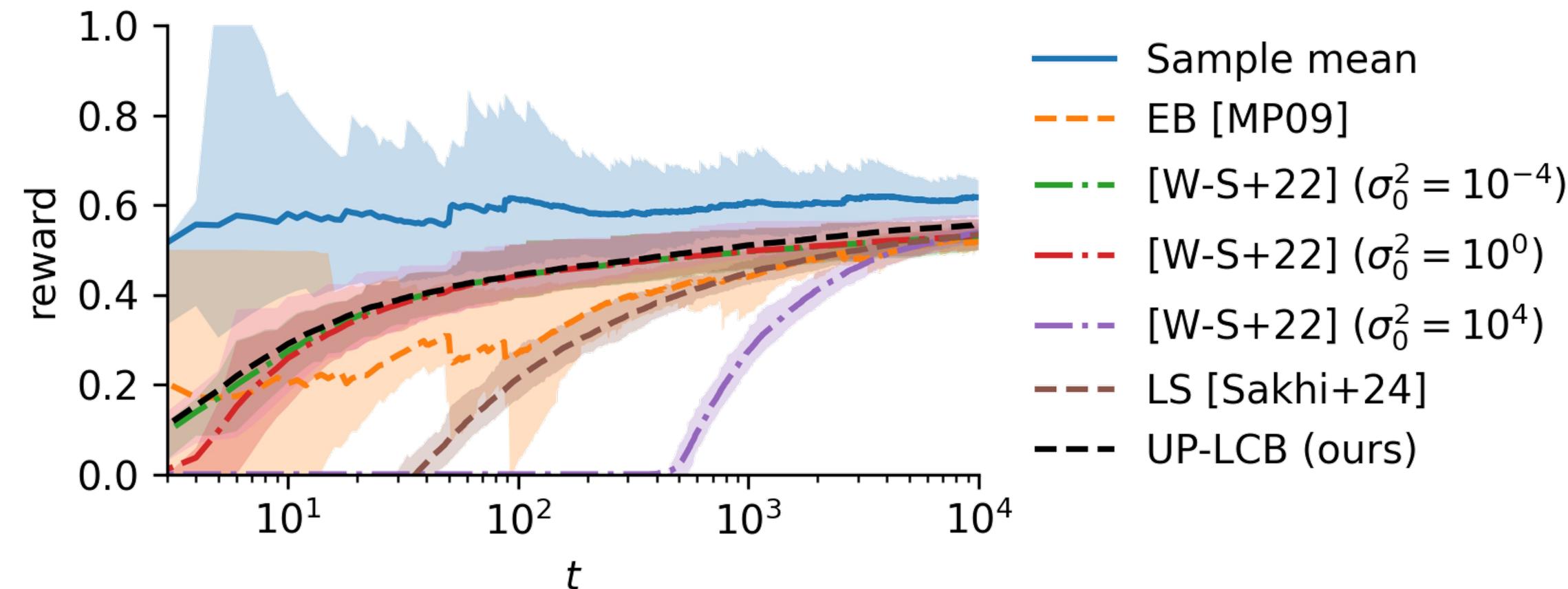
# Experiment: UP-LCB

- Synthetic reward processes with **infinite** 4th moment
- Empirical Bernstein (EB) = Plug-in empirical variance (heuristic)



# Experiment: UP-LCB

- Synthetic reward processes with **infinite** 4th moment
- Empirical Bernstein (EB) = Plug-in empirical variance (heuristic)



# Experiment: UP-LCB + Pessimism

- Experiment setup from [WangKS2024]

Relative improvement from ‘Naive’ (higher = better)

Dataset	PenDigits			SatImage			JPVowel			
	0.01	0.1	1	0.01	0.1	1	0.01	0.1	1	
Size	34.93	<b>23.60</b>	3.02	30.88	26.26	<b>17.84</b>	40.60	<b>27.30</b>	<b>0.87</b>	
[Sakhi+24]	→ LS	34.93	<b>23.60</b>	3.02	30.88	26.26	<b>17.84</b>	40.60	<b>27.30</b>	<b>0.87</b>
	→ EB	41.34	<b>22.67</b>	<b>3.77</b>	33.76	26.79	<b>17.68</b>	<b>44.58</b>	<b>27.84</b>	<b>1.59</b>
	PUB (ours)	<b>44.02</b>	<b>22.74</b>	<b>4.01</b>	<b>36.14</b>	<b>28.16</b>	<b>16.35</b>	<b>44.51</b>	<b>28.59</b>	<b>1.40</b>

**Empirical Bernstein:** A heuristic based on the central limit theorem.

# Experiment: UP-LCB + Pessimism

- Experiment setup from [WangKS2024]

Relative improvement from ‘Naive’ (higher = better)

Dataset	PenDigits			SatImage			JPVowel			
	0.01	0.1	1	0.01	0.1	1	0.01	0.1	1	
Size	34.93	<b>23.60</b>	3.02	30.88	26.26	<b>17.84</b>	40.60	<b>27.30</b>	<b>0.87</b>	
[Sakhi+24]	→ LS	34.93	<b>23.60</b>	3.02	30.88	26.26	<b>17.84</b>	40.60	<b>27.30</b>	<b>0.87</b>
	→ EB	41.34	<b>22.67</b>	<b>3.77</b>	33.76	26.79	<b>17.68</b>	<b>44.58</b>	<b>27.84</b>	<b>1.59</b>
	PUB (ours)	<b>44.02</b>	<b>22.74</b>	<b>4.01</b>	<b>36.14</b>	<b>28.16</b>	<b>16.35</b>	<b>44.51</b>	<b>28.59</b>	<b>1.40</b>

**Empirical Bernstein:** A heuristic based on the central limit theorem.

Always either **the best or statistically indistinguishable from the best!**

# Part 2. Off-Policy Learning

---

# Off-Policy Learning (=Optimization)

- Policy class is infinitely large:  $|\Pi| = \infty$  (like neural networks)

# Off-Policy Learning (=Optimization)

- Policy class is infinitely large:  $|\Pi| = \infty$  (like neural networks)
- Using LCB-based pessimism? Root-finding via optimization is possible, but...

# Off-Policy Learning (=Optimization)

- Policy class is infinitely large:  $|\Pi| = \infty$  (like neural networks)
- Using LCB-based pessimism? Root-finding via optimization is possible, but...

$$\max_{\pi \in \Pi} \max_{\alpha \geq 0} \min_{\nu} \left\{ \nu + \alpha \left( \frac{1}{\sqrt{\pi(n+1)}} \max_{b \in (0,1)} \log \text{wealth}_t^b(\nu; \pi) - \log \frac{|\Pi|}{\delta} \right) \right\}$$

# Off-Policy Learning (=Optimization)

- Policy class is infinitely large:  $|\Pi| = \infty$  (like neural networks)
- Using LCB-based pessimism? Root-finding via optimization is possible, but...

$$\max_{\pi \in \Pi} \max_{\alpha \geq 0} \min_{\nu} \left\{ \nu + \alpha \left( \frac{1}{\sqrt{\pi(n+1)}} \max_{b \in (0,1)} \log \text{wealth}_t^b(\nu; \pi) - \log \frac{|\Pi|}{\delta} \right) \right\}$$

DIFFICULT!

# Off-Policy Learning (=Optimization)

- Policy class is infinitely large:  $|\Pi| = \infty$  (like neural networks)
- Using LCB-based pessimism? Root-finding via optimization is possible, but...

$$\max_{\pi \in \Pi} \max_{\alpha \geq 0} \min_{\nu} \left\{ \nu + \alpha \left( \frac{1}{\sqrt{\pi(n+1)}} \max_{b \in (0,1)} \log \text{wealth}_t^b(\nu; \pi) - \log \frac{|\Pi|}{\delta} \right) \right\}$$

DIFFICULT!

- Oracle-efficiency: objective is **optimizable** efficiently with an optimization oracle

# Off-Policy Learning (=Optimization)

- **Our proposal:** A family of objectives achieving the SOTA guarantee

$$\max_{\pi \in \Pi} \sum_{t=1}^n \phi(\beta \tilde{r}_t(\pi))$$

↑  
hyperparameter

# Off-Policy Learning (=Optimization)

- **Our proposal:** A family of objectives achieving the SOTA guarantee

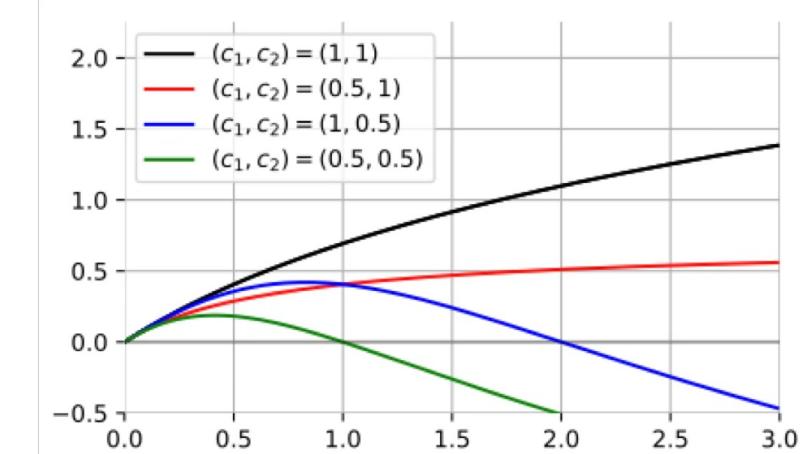
$$\max_{\pi \in \Pi} \sum_{t=1}^n \phi(\beta \tilde{r}_t(\pi))$$

↑  
hyperparameter

**Observation:** any  $\phi(\cdot)$  satisfying

$$-\log\left(1 - x + \frac{x^2}{c_1 + c_2 x}\right) \leq \phi(x) \leq \log(1 + x)$$

implements pessimism ( $\approx$  underestimation)



# Off-Policy Learning (=Optimization)

- **Our proposal:** A family of objectives achieving the SOTA guarantee

$$\max_{\pi \in \Pi} \sum_{t=1}^n \phi(\beta \tilde{r}_t(\pi))$$

↑  
hyperparameter

Logarithmic smoothing:  
[Sakhi+24]

Freezing:  
(Ours)

**Observation:** any  $\phi(\cdot)$  satisfying

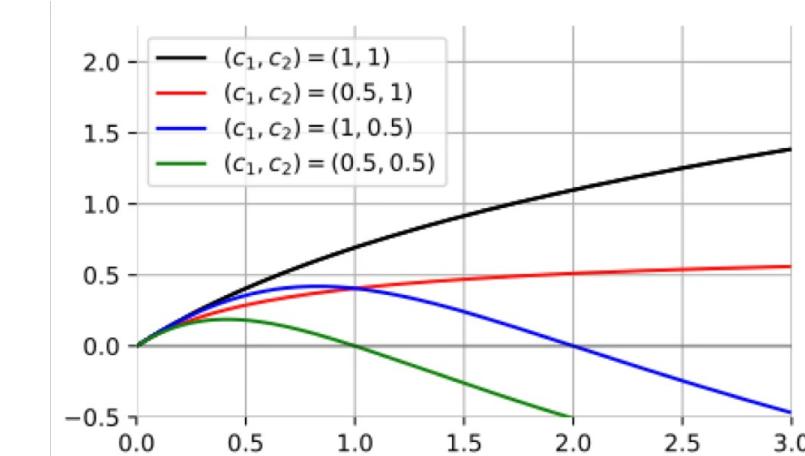
$$-\log\left(1 - x + \frac{x^2}{c_1 + c_2 x}\right) \leq \phi(x) \leq \log(1 + x)$$

implements pessimism ( $\approx$  underestimation)

$$\phi(x) = \log(1 + x)$$

$$\phi(x) = \log(1 + \mathbb{I}\{x \leq 1\}x)$$

more aggressively reduces variance;  
advantageous with small sample size



# Off-Policy Learning (=Optimization)

- **Our proposal:** A family of objectives achieving the SOTA guarantee

$$\max_{\pi \in \Pi} \sum_{t=1}^n \phi(\beta \tilde{r}_t(\pi))$$

↑  
hyperparameter

Logarithmic smoothing:  
[Sakhi+24]

Freezing:  
(Ours)

Theorem

$$\mu(\pi^*) - \mu(\hat{\pi}_n) \leq \beta \mathbb{E} \left[ \frac{(\tilde{r}^{\pi^*})^2}{c_1 + c_2 \tilde{r}^{\pi^*}} \right] + \frac{2}{\beta n} \ln \frac{|\Pi|}{\delta} - F_\beta(\phi)$$

**Observation:** any  $\phi(\cdot)$  satisfying

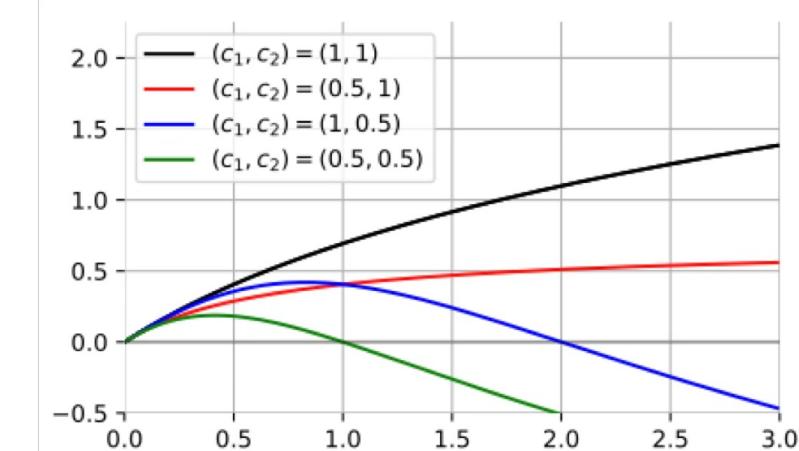
$$-\log \left( 1 - x + \frac{x^2}{c_1 + c_2 x} \right) \leq \phi(x) \leq \log(1 + x)$$

implements pessimism ( $\approx$  underestimation)

$$\phi(x) = \log(1 + x)$$

$$\phi(x) = \log(1 + \mathbb{I}\{x \leq 1\}x)$$

more aggressively reduces variance;  
advantageous with small sample size

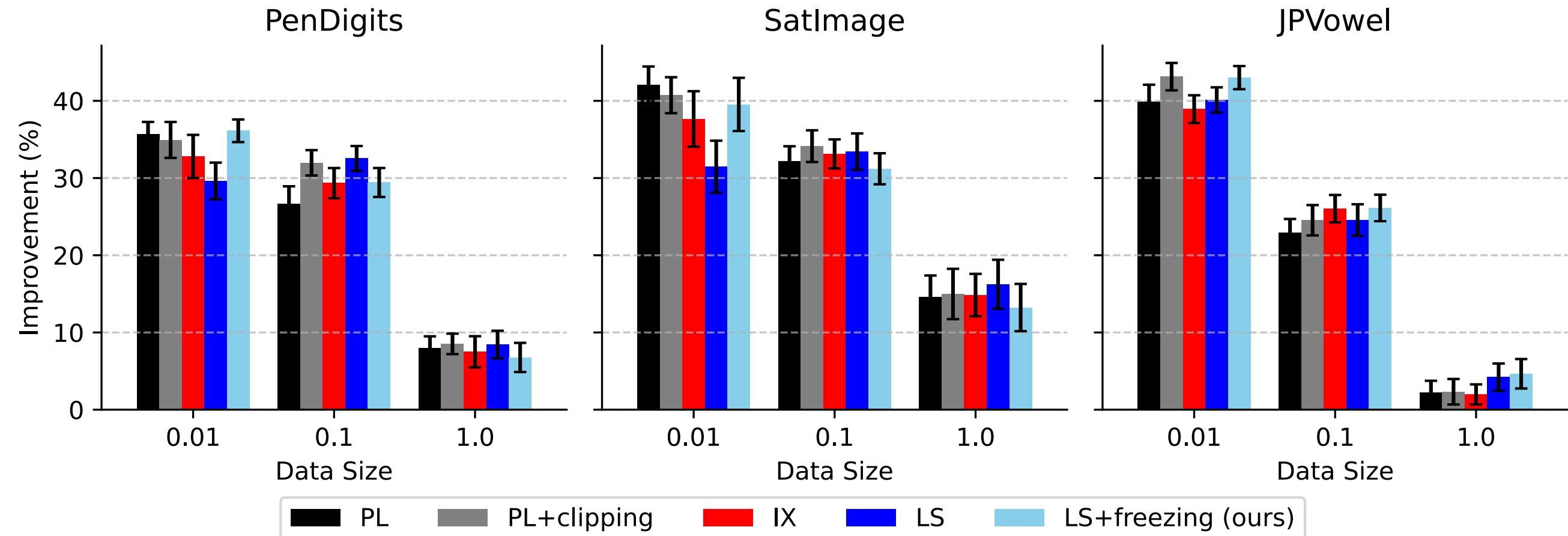


“negative influence”

$$F_\beta(\phi^{\text{freezing}}) \geq F_\beta(\phi^{\text{clipping}}) \geq F_\beta(\phi^{\text{LS}})$$

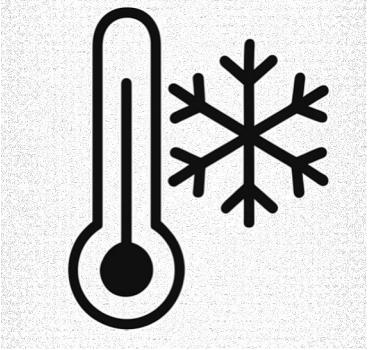
# Experiment: Off-Policy Learning

- Experiment setup from [WangKS2024] (value regression setup)
- Again, best or almost the best!



# Concluding Remarks

- **Betting** for off-policy selection (and evaluation too ☺)
  - **Freezing** for off-policy learning
- ⇒ Improved second-order bounds for optimality gap!
- 
- **Efficient offline learning** will become increasingly important
  - **Future research directions**
    - Betting-type idea for  $|\Pi| = \infty$ ?
    - Similarly, more sophisticated techniques for RL?



**Thank You!**