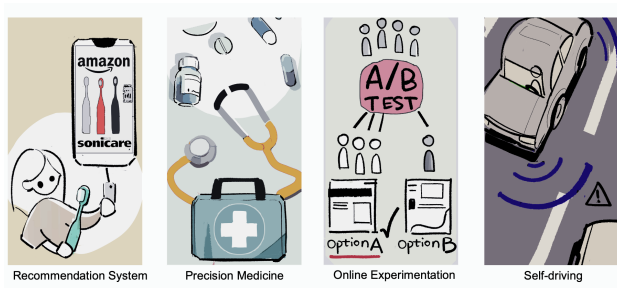


# ONLINE ESTIMATION AND INFERENCE FOR ROBUST POLICY EVALUATION IN REINFORCEMENT LEARNING

**Weidong Liu, Jiyuan Tu, Xi Chen, Yichen Zhang**

NeurIPS 2025

# SEQUENTIAL DECISION-MAKING



## ► Regret Minimization

- Bandit/RL algorithms are designed for regret minimization
- Regret: how much worse it performs compared to an offline oracle

## ► Statistical inference or uncertainty quantification

- Infer the effect of a policy; Gaining generalizable knowledge; Identifying new directions; Crucial for scientific discovery

# STATISTICAL INFERENCE IN SEQUENTIAL DECISION-MAKING

- ▶ Bandit /RL Algorithms Induce Dependence
  - Observations are not independent over time
    - ▶ Use past history to select the action in the new context
    - ▶ Data is adaptively collected
  - Consequences: Introduce a bias, or potentially non-normality
    - ▶ Villar, Bowden, Wason (2015); Deshpande et al (2021); Khamaru, Mackey, Wainwright (2021); Zhang, Janson, Murphy (2021, 2022); Chen, Lu, Song (2021ab)
    - ▶ Not fully-online algorithms: Heavy cost on computation and storage
- ▶ Design inference procedures that are simultaneously: **sample-efficient**, **computation-efficient**, valid for **adaptive collected** data?

## POLICY EVALUATION W. LINEAR FUNCTION APPROXIMATION

For an MDP  $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R})$ , we evaluate the cumulative reward for a policy  $\pi$

$$J^*(s) = \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k \mathcal{R}(s_k) \mid s_0 = s \right].$$

- ▶ Linear function approximation:  $\tilde{J}(s, \boldsymbol{\theta}) = \boldsymbol{\theta}^\top \boldsymbol{\phi}(s) = \sum_{l=1}^d \theta_l \phi_l(s)$ .
- ▶ The optimal parameter minimizes the MSPBE (Tsitsiklis and Van Roy, 1997):

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \left\| \tilde{J}(\boldsymbol{\theta}) - \Pi \mathcal{T}^\pi \tilde{J}(\boldsymbol{\theta}) \right\|_D^2,$$

where  $D$  denotes the stable distribution under  $\pi$ , and  $\mathcal{T}$  is the Bellman operator:

$$\mathcal{T}(Q) = \mathcal{R} + \gamma \mathcal{P}Q.$$

## LEAST-SQUARES TD(0)

- ▶ Consider an equivalent **fixed-point formulation** (Kolter and Ng, 2009):

$$\boldsymbol{\theta}^* = \operatorname{argmin}_{\boldsymbol{u} \in \mathbb{R}^d} \mathbb{E} |\boldsymbol{\phi}^\top(s) \boldsymbol{u} - (\mathcal{R}(s) + \gamma \boldsymbol{\phi}^\top(s') \boldsymbol{\theta}^*)|^2.$$

- ▶ Given a trajectory  $\{(s_n, \mathcal{R}(s_n), s_{n+1})\}_{n=1}^\infty$ , the linear TD update (Sutton, 1988):

$$\hat{\boldsymbol{\theta}}_{n+1} = \hat{\boldsymbol{\theta}}_n + \alpha_{n+1} \{ (\boldsymbol{\phi}(s_n) - \gamma \boldsymbol{\phi}^\top(s_{n+1}))^\top \hat{\boldsymbol{\theta}}_n - \mathcal{R}(s_n) \} \boldsymbol{\phi}(s_n).$$

- ▶ Problems:
  - Linear TD(0) algorithm is **vulnerable to corruption** in reward  $\mathcal{R}$ ;
  - How to conduct **online inference** on  $\boldsymbol{\theta}^*$ , or  $J^*(s)$ ?

# ROBUST TEMPORAL DIFFERENCE

- ▶ Reward  $\mathcal{R}$  comes from distribution  $(1 - \alpha_n)P + \alpha_n Q$ :
  - $Q$  denotes arbitrary outlier distribution;
  - $P$  has  $(1 + \delta)$ -moment;
- ▶ Pseudo Huber loss  $f_\tau(x) = \tau^2(\sqrt{1 + (x/\tau)^2} - 1)$ ;
  - Smoothed version of Huber: near quadratic locally, Lipschitz globally.
- ▶ Solve the robust fixed-point equation

$$\theta_\tau^* = \operatorname{argmin}_{u \in \mathbb{R}^d} \mathbb{E} f_\tau(\phi^\top(s)u - (\mathcal{R}(s) + \gamma \phi^\top(s')\theta_\tau^*));$$

- ▶ No longer quadratic  $\Rightarrow$  Gradient methods performs nonlinear updates
- ▶ Equivalent to  $\mathbb{E}[\phi(s)g_\tau((\phi^\top(s) - \gamma \phi^\top(s'))\theta_\tau^* - \mathcal{R}(s))] = 0$ , with  $g_\tau(x) = f'_\tau(x)$ .

## ONLINE NEWTON METHOD

- ▶ Denote  $\mathbf{X}_i = \phi(s_i)$ ,  $\mathbf{Z}_i = \phi(s_i) - \gamma\phi(s_{i+1})$ , and  $b_i = \mathcal{R}(s_i)$ ;
- ▶ Target problem can be written as

$$\mathbb{E}[\mathbf{X}g_\tau(\mathbf{Z}^\top \boldsymbol{\theta}_\tau^* - b)] = 0;$$

- ▶ Propose the following online Newton update:

$$\hat{\boldsymbol{\theta}}_{n+1} = \frac{1}{n+1} \sum_{i=0}^n \hat{\boldsymbol{\theta}}_i - \hat{\mathbf{H}}_{n+1}^{-1} \frac{1}{n+1} \sum_{i=0}^n \mathbf{X}_{i+1} g_{\tau_{i+1}}(\mathbf{Z}_{i+1}^\top \hat{\boldsymbol{\theta}}_i - b_{i+1}),$$

$$\text{with } \hat{\mathbf{H}}_{n+1} = \frac{1}{n+1} \sum_{i=0}^n \mathbf{X}_{i+1} \mathbf{Z}_{i+1}^\top g'_{\tau_{i+1}}(\mathbf{Z}_{i+1}^\top \hat{\boldsymbol{\theta}}_i - b_{i+1}).$$

# ONLINE NEWTON METHOD

- ▶ Reformulate the update as

$$\hat{\boldsymbol{\theta}}_{n+1} = \bar{\boldsymbol{\theta}}_{n+1} - \hat{\mathbf{H}}_{n+1}^{-1} \mathbf{G}_{n+1};$$

- ▶ Online gradient update

$$\mathbf{G}_{n+1} = \frac{n}{n+1} \mathbf{G}_n + \frac{1}{n+1} \mathbf{X}_{n+1} g_{\tau_{n+1}}(\mathbf{Z}_{n+1}^\top \hat{\boldsymbol{\theta}}_n - b_{n+1});$$

- ▶ Online Hessian update (by Sherman-Morrison formula) – only scalar inverse:

$$\begin{aligned} \hat{\mathbf{H}}_{n+1}^{-1} &= \frac{n+1}{n} \hat{\mathbf{H}}_n^{-1} - \frac{n+1}{n^2} \hat{\mathbf{H}}_n^{-1} \mathbf{X}_{n+1} \\ &\quad \times \left[ \frac{1}{n} \mathbf{Z}_{n+1}^\top \hat{\mathbf{H}}_n^{-1} \mathbf{X}_{n+1} + \{g'_{\tau_{n+1}}(\mathbf{Z}_{n+1}^\top \hat{\boldsymbol{\theta}}_n - b_{n+1})\}^{-1} \right]^{-1} \mathbf{Z}_{n+1}^\top \hat{\mathbf{H}}_n^{-1}. \end{aligned}$$



# CONVERGENCE RATE

## Theorem 1

Take the thresholding parameter  $\tau_i = C_\tau i^\beta$ . Assume  $n_0$  is sufficiently large and the initial value  $|\hat{\boldsymbol{\theta}}_0 - \boldsymbol{\theta}^*|_2 \leq c_0$  for some  $c_0 < 1$ . Then there is

$$\mathbb{P}\left(\cap_{i=n_0}^n \{|\hat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}^*|_2 \geq C\sqrt{d}e_i\}\right) \geq 1 - cn_0^{-\nu},$$

where 
$$e_n = \alpha_n \tau_n + \tau_n^{-\min(\delta, 2)} + \sqrt{\frac{\tau_n^{(1-\delta)+} \log n}{n}} + \frac{\tau_n \log^2 n}{n} + \frac{1}{\sqrt{d}}(c_0)^{2^{n-n_0}}.$$

Assume  $m_n = \alpha_n n$  outliers among  $n$  samples. We further have

$$|\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^*|_2 = O_{\mathbb{P}}\left(\sqrt{d}\left(\sqrt{\frac{\log n}{n}} + \frac{\log^2 n}{n^{1-\beta}} + \frac{m_n}{n^{1-\beta}} + \frac{1}{n^{2\beta}}\right)\right).$$

# CONVERGENCE RATE (NO CONTAMINATION)

## Corollary 2

When the contamination rate  $\alpha_n = 0$ ,

- ▶ When  $\delta \in (0, 1]$ , we specify  $\tau_i = C_\tau(i/\log i)^{1/(1+\delta)}$ . Then

$$|\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^*|_2 = O_{\mathbb{P}} \left( \sqrt{d} \left( \frac{\log n}{n} \right)^{\frac{\delta}{1+\delta}} \right).$$

- ▶ When  $\delta > 1$ , we specify  $\tau_i = C_\tau(i/\log i)^{1/2}$ . Then

$$|\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^*|_2 = O_{\mathbb{P}} \left( \sqrt{d} \left( \frac{\log n}{n} \right)^{1/2} \right).$$

# ASYMPTOTIC NORMALITY

## Theorem 3

If the contamination rate  $\alpha_n = o(1/(\sqrt{n}\tau_n))$  and  $n^{1/4} = o(\tau_n)$ ,

$$\frac{\sqrt{n}}{\sigma_v} \mathbf{v}^\top (\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^*) \xrightarrow{d} \mathcal{N}(0, 1), \text{ where } \sigma_v^2 = \mathbf{v}^\top \mathbf{H}^{-1} \boldsymbol{\Sigma} (\mathbf{H}^\top)^{-1} \mathbf{v},$$

as  $n \rightarrow \infty$ . Here  $\boldsymbol{\Sigma} = \sum_{k=-\infty}^{\infty} \mathbb{E}[\mathbf{X}_0 \mathbf{X}_k^\top (\mathbf{Z}_0^\top \boldsymbol{\theta}^* - b_0)(\mathbf{Z}_k^\top \boldsymbol{\theta}^* - b_k)]$ .

When  $\alpha_n = 0$  (no contamination) and  $\tau_i = C_\tau i^\beta$  for  $\beta \geq 3/4$ ,

$$\sqrt{n} \mathbf{v}^\top (\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^*) = \mathbf{v}^\top \mathbf{H}^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{X}_i (\mathbf{Z}_i^\top \boldsymbol{\theta}^* - b_i) + O_{\mathbb{P}}\left(\frac{d \log n}{\sqrt{n}}\right).$$

- ▶ Asymptotic Normality is valid when  $d^2/n = o(1)$ , ignoring logarithm terms.
- ▶ For first-order methods, best remainder rate is  $O_{\mathbb{P}}(dn^{-1/3})$ , even under i.i.d.

# ONLINE STATISTICAL INFERENCE

- Asymptotic covariance matrix  $\mathbf{H}^{-1}\boldsymbol{\Sigma}(\mathbf{H}^\top)^{-1}$  with

$$\boldsymbol{\Sigma} = \sum_{k=-\infty}^{\infty} \mathbb{E}[\mathbf{X}_0 \mathbf{X}_k^\top (\mathbf{Z}_0^\top \boldsymbol{\theta}^* - b_0)(\mathbf{Z}_k^\top \boldsymbol{\theta}^* - b_k)].$$

- Construct the estimator

$$\hat{\boldsymbol{\Sigma}}_n = \frac{1}{n} \sum_{i=1}^n \sum_{k=-\lceil \lambda \log i \rceil \wedge (i-1)}^{\lceil \lambda \log i \rceil \wedge (i-1)} \mathbf{X}_i \mathbf{X}_{i-k}^\top g_{\tau_i}(\mathbf{Z}_i^\top \hat{\boldsymbol{\theta}}_{i-1} - b_i) g_{\tau_{i-k}}(\mathbf{Z}_{i-k}^\top \hat{\boldsymbol{\theta}}_{i-k-1} - b_{i-k}).$$

- Requires only  $O(\lceil \log n \rceil)$  memory due to  $\phi$ -mixing property.

# ONLINE ESTIMATION OF LONG-RUN COVARIANCE MATRIX

## Theorem 4

The covariance estimator  $\hat{\Sigma}_n$  satisfies

$$\|\hat{\Sigma}_n - \Sigma\| = O_{\mathbb{P}} \left( \sqrt{d} \tau_n^2 \alpha_n + \sqrt{d} \tau_n^{-1} + \tau_n \sqrt{\frac{d \log n}{n}} + \frac{d \tau_n^2 \log^2 n}{n} \right).$$

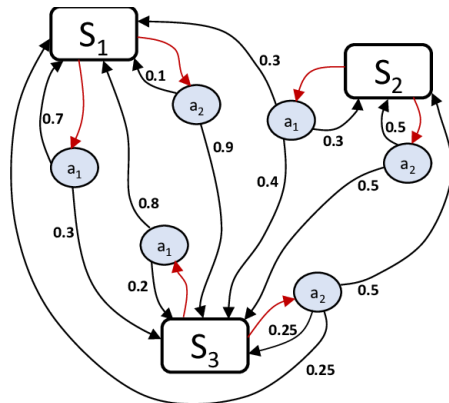
Given a unit vector  $v$ , a confidence interval with nominal level  $(1 - \xi)$  is

$$\left[ v^{\top} \hat{\theta}_n - q_{1-\xi/2} \hat{\sigma}_v, v^{\top} \hat{\theta}_n + q_{1-\xi/2} \hat{\sigma}_v \right],$$

with  $\hat{\sigma}_v^2 = v^{\top} \hat{H}_n^{-1} \hat{\Sigma}_n (\hat{H}_n^{\top})^{-1} v$  and  $q_{1-\xi/2} = \Phi^{-1}(1 - \xi/2)$ .

# SIMULATION STUDIES

Infinite-Horizon MDP

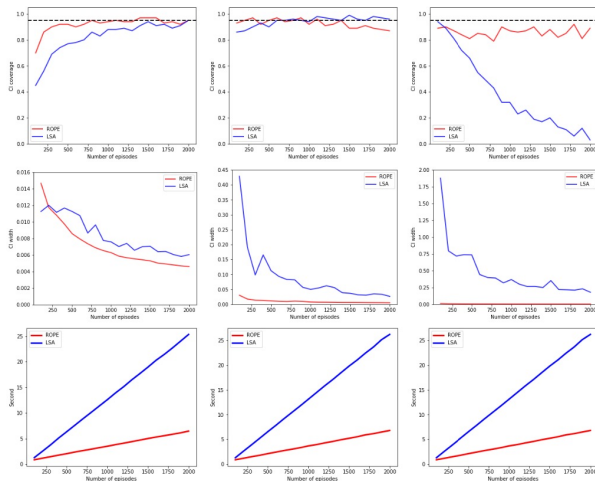


FrozenLake Environment



# FROZENLAKE ENVIRONMENT

## COMPARATIVE ANALYSIS



- ▶ Comparison between ROPE and LSA;
- ▶ Coverage probability (the first row)  
the width of confidence interval (the second row)  
computing time (the third row) ;
- ▶ Contamination rate  $\alpha_n$  in  $\{0, n^{-1}, 0.05n^{-1/2}\}$ .

## CONCLUSIONS

- ▶ An online policy evaluation **robust to heavy-tailedness and outliers**;
- ▶ Technically difficult due to **non-linear** stochastic gradient update;
- ▶ An **online statistical inference** procedure for policy evaluation;
- ▶ 2nd-order approach for improved efficiency without computation loss;
- ▶ **Improved** higher-order convergence to the asymptotic normality.
- ▶ The *Annals of Statistics* (to appear), arXiv: 2310.02581.

Thanks!