

# A Statistical Analysis of Polyak-Ruppert Averaged Q-learning

Xiang Li, Wenhao Yang, Jiadong Liang, Zhihua Zhang, Michael I. Jordan

Peking University, University of California Berkeley

## Introduction

The classic Q-Learning:

$$Q_t = (1 - \eta_t)Q_{t-1} + \eta_t \widehat{T}_t(Q_{t-1}) \text{ where}$$

$$\widehat{T}_t(Q)(s, a) = r_t(s, a) + \gamma \max_{a' \in \mathcal{A}} Q(s_t, a'),$$

is the empirical estimate of the Bellman operator

$$\mathcal{T}(Q)(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \max_{a' \in \mathcal{A}} Q(s', a').$$

- Many analysis for the last-iterate  $Q_T$ , while the averaged iterate  $\bar{Q}_T$  is less well understood,

$$\bar{Q}_T = \frac{1}{T} \sum_{t=1}^T Q_t.$$

- Classic results imply under mild regularity conditions,

$$\sqrt{T}(\bar{Q}_T - Q^*) \xrightarrow{d} \mathcal{Z} \sim \mathcal{N}(\mathbf{0}, \text{Var}_Q),$$

where the asymptotic variance  $\text{Var}_Q$  is given by

$$\text{Var}_Q := (\mathbf{I} - \gamma \mathbf{P}^{\pi^*})^{-1} \text{Var}(\mathbf{Z})(\mathbf{I} - \gamma \mathbf{P}^{\pi^*})^{-\top} \quad (1)$$

with  $\mathbf{Z} \in \mathbb{R}^D$ ,  $\mathbf{Z}(s, a) = [\widehat{T}_t(Q) - \mathcal{T}(Q)](s, a)$ .

## Questions Studied

An observation from  $\sqrt{T}(\bar{Q}_T - Q^*) \xrightarrow{d} \mathcal{Z}$ :

$$\sqrt{T} \mathbb{E} \|\bar{Q}_T - Q^*\|_{\infty} \rightarrow \mathbb{E} \|\mathcal{Z}\|_{\infty} \approx \sqrt{\ln D} \sqrt{\|\text{diag}(\text{Var}_Q)\|_{\infty}}.$$

- How to obtain a valid non-asymptotic bound?
- Is this asymptotic variance  $\text{Var}_Q$  optimal?
- Can we do statistical inference without estimating  $\text{Var}_Q$ ?

## Discounted infinite-horizon MDPs

- $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \gamma, P, R, r)$  with  $\gamma \in [0, 1)$ .
- $P: \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ : probability transition kernel.
- $R: \mathcal{S} \times \mathcal{A} \rightarrow [0, \infty)$ : random reward,  $r = \mathbb{E}R$ .
- A policy  $\pi: \mathcal{S} \rightarrow \mathcal{A}$  with its (Q-)value defined by

$$Q^{\pi}(s, a) = \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a \right].$$

- Assume unique optimal policy  $\pi^*$ .

## Non-asymptotic Convergence

Assume  $0 \leq R(s, a) \leq 1$  for all  $(s, a)$ .

- If  $\eta_t = t^{-\alpha}$  with  $\alpha \in (0.5, 1)$ ,  $\mathbb{E} \|\bar{Q}_T - Q^*\|_{\infty} =$

$$\mathcal{O} \left( \sqrt{\|\text{diag}(\text{Var}_Q)\|_{\infty}} \sqrt{\frac{\ln D}{T}} + \frac{\sqrt{\ln D}}{(1-\gamma)^3 T^{1-\frac{\alpha}{2}}} \right) + \tilde{\mathcal{O}} \left( \frac{1}{(1-\gamma)^{3+\frac{2}{1-\alpha}} T} + \frac{\gamma}{(1-\gamma)^{4+\frac{1}{1-\alpha}} T^{\alpha}} \right).$$

- The dominant term (in red) is of the same magnitude as the variance-reduced Q-learning.

- If  $\eta_t = \frac{1}{1+(1-\gamma)t}$ ,  $\mathbb{E} \|\bar{Q}_T - Q^*\|_{\infty} =$

$$\mathcal{O} \left( \sqrt{\frac{\|\text{Var}(\mathbf{Z})\|_{\infty}}{(1-\gamma)^2}} \sqrt{\frac{\ln D}{T}} \right) + \tilde{\mathcal{O}} \left( \frac{1}{(1-\gamma)^6 T} \right).$$

- Because  $\|\text{diag}(\text{Var}_Q)\|_{\infty} \leq \frac{1}{(1-\gamma)^2} \|\text{Var}(\mathbf{Z})\|_{\infty}$ , this rate is slightly loose. (How to tighten it? Unclear).

## Conclusion

- Averaged Q-learning achieves both the worst-case and instance-dependent optimality asymptotically.
- The asymptotic variance of averaged Q-learning is optimal among all RAL estimators.
- We established a FCLT that facilitates online statistical inference.

## Semiparametric Statistics

- Our MDP model  $\mathcal{M}$  has parameter  $\theta = (P, R)$ .
- The transition  $P$  is parametric due to discrete action-state space, while the random reward is totally non-parametric.
- (Unformal) An estimator is **regular** if its limiting distribution is unaffected by local changes in the data generating process.

- (Unformal) An estimator  $\widehat{Q}_T$  is **asymptotically linear** with a measurable random function  $\phi(\mathbf{r}_t, \mathbf{P}_t) \in \mathbb{R}^D$  such that

$$\sqrt{T}(\widehat{Q}_T - Q^*) = \frac{1}{\sqrt{T}} \sum_{t=1}^T \phi(\mathbf{r}_t, \mathbf{P}_t) + o_{\mathbb{P}}(1), \quad (2)$$

where  $\mathcal{D} = \{(\mathbf{r}_t, \mathbf{P}_t)\}_{t \in [T]}$  is the collected i.i.d. data and  $\phi(\cdot, \cdot)$  is referred to as an influence function satisfying that  $\mathbb{E}\phi(\mathbf{r}_t, \mathbf{P}_t) = \mathbf{0}$  and  $\mathbb{E}\phi(\mathbf{r}_t, \mathbf{P}_t)\phi(\mathbf{r}_t, \mathbf{P}_t)^{\top}$  is finite and nonsingular.

## Numeral Validation

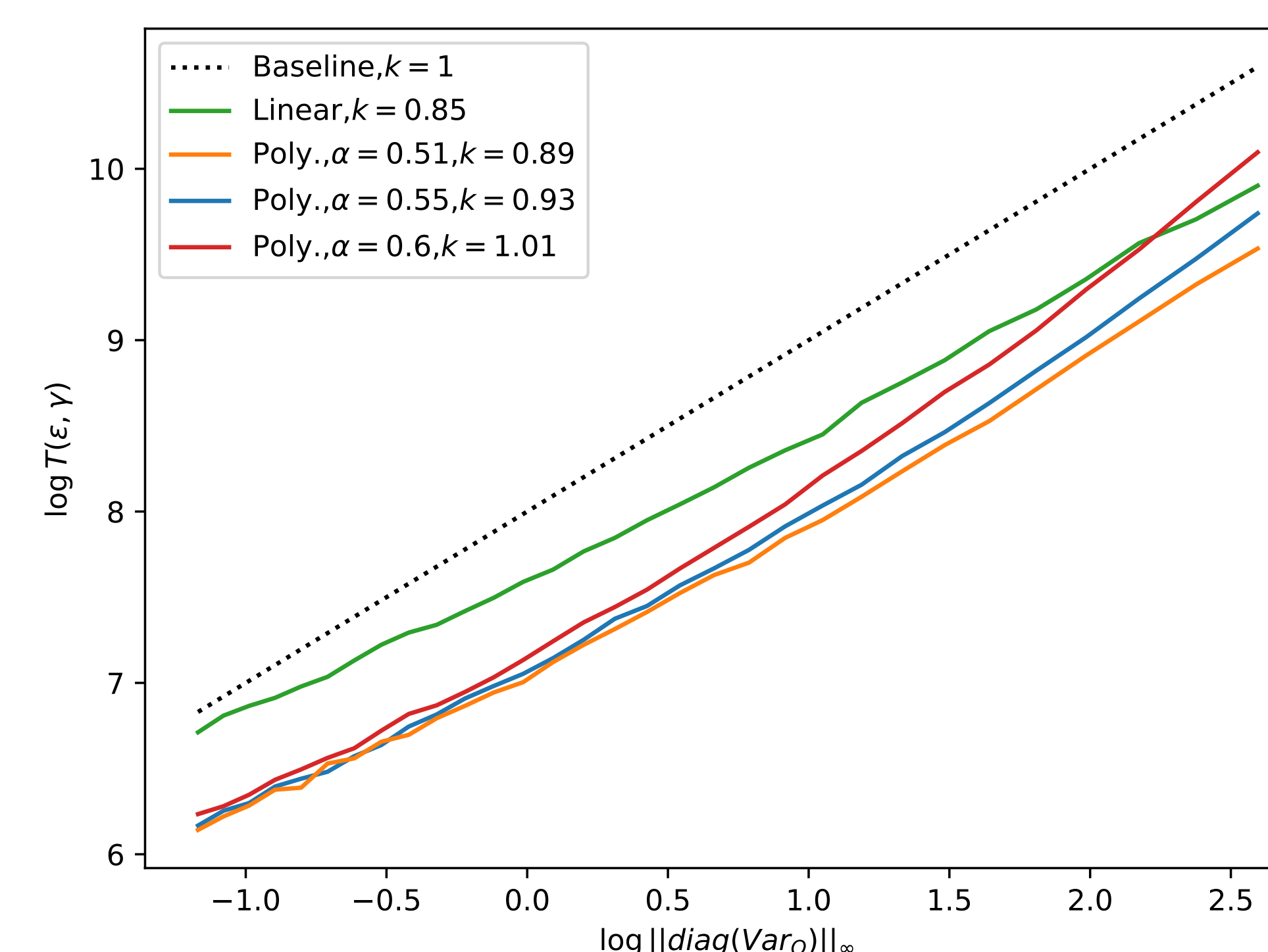


Figure 1: Log-log plots of the sample complexity  $T(\epsilon, \gamma)$  versus the asymptotic variance  $\|\text{diag}(\text{Var}_Q)\|_{\infty}$  where  $T(\epsilon, \gamma) = \inf\{T : \mathbb{E} \|\bar{Q}_T - Q^*\|_{\infty} \leq \epsilon\}$  in a  $\gamma$ -discounted MDP.

## Optimal Variance

We prove the following results.

- RAL = *regular and asymptotically linear*.
- Given the dataset  $\mathcal{D}$ , the asymptotic variance matrix of any RAL estimator  $\widehat{Q}_T$  of  $Q^*$  computed from  $\mathcal{D}$  satisfying

$$\lim_{T \rightarrow \infty} T \mathbb{E}(\widehat{Q}_T - Q^*)(\widehat{Q}_T - Q^*)^{\top} \succeq \text{Var}_Q,$$

where  $\mathbf{A} \succeq \mathbf{B}$  means  $\mathbf{A} - \mathbf{B}$  is positive semidefinite.

- The averaged iterate  $\bar{Q}_T$  is the **optimal RAL** estimator for  $Q^*$  due to

$$\sqrt{T}(\bar{Q}_T - Q^*) = \frac{1}{\sqrt{T}} \sum_{t=1}^T (\mathbf{I} - \gamma \mathbf{P}^{\pi^*})^{-1} \mathbf{Z}_t + o_{\mathbb{P}}(1),$$

where  $\mathbf{Z}_t = (\mathbf{r}_t - r) + (\mathbf{P}_t - \mathbf{P})\mathbf{V}^*$ .

Averaged Q-Learning iterates has the optimal asymptotic variance matrix.

## Functional Central Limit Theorem

We **can do** statistical inference **without** estimating  $\text{Var}_Q$  by using the FCLT.

- Given  $\{Q_t\}_{t \in [T]}$ , its partial-sum processes is

$$\phi_T(r) := \frac{1}{\sqrt{T}} \sum_{t=1}^{\lfloor Tr \rfloor} (Q_t - Q^*), r \in [0, 1]. \quad (3)$$

- We assume that (i)  $\sup_{s,a} \mathbb{E}R^4(s, a) < \infty$ ; (ii)  $\pi^*$  is unique; and (iii)  $\eta_t = t^{-\alpha}$  ( $0.5 < \alpha < 1$ ).
- We show the following weak convergence

$$\phi_T(\cdot) \xrightarrow{w} \text{Var}_Q^{1/2} \mathbf{B}_D(\cdot), \quad (4)$$

where  $\text{Var}_Q$  is defined in (1) and  $\mathbf{B}_D$  is the standard  $D$ -dim Brownian motion on  $[0, 1]$ .

## Online Statistics Inference

- By continuous mapping theorem, for any continuous functional  $f: \mathbb{D}[0, 1] \rightarrow \mathbb{R}$ ,

$$f(\phi_T) \xrightarrow{d} f(\text{Var}_Q^{1/2} \mathbf{B}_D).$$

- Once  $f$  is scale-invariant, we have  $f(\text{Var}_Q^{1/2} \mathbf{B}_D) = f(\mathbf{B}_D)$  and thus

$$f(\phi_T) \xrightarrow{d} f(\mathbf{B}_D).$$

- (1) a **pivotal** statistic relying on  $Q^*$  and  $\mathcal{D}$  **known** distribution
- (2) can be computed **online efficiently**

- Confidence intervals for  $Q^*$  can be obtained by inverting some constraint regime on  $f(\phi_T)$ .

- We find a  $f$  following the spirit of t-statistics:

$$f(\mathbf{B}) := \mathbf{B}(1)^{\top} \left( \int_0^1 \bar{\mathbf{B}}(r) \bar{\mathbf{B}}(r)^{\top} dr \right)^{-1} \mathbf{B}(1)$$

with  $\bar{\mathbf{B}}(r) = \mathbf{B}(r) - r\mathbf{B}(1)$ .

- A numerical illustration:

