

Sharp Variance-Dependent Bounds in Reinforcement Learning: Best of Both Worlds in Stochastic and Deterministic Environments



Runlong Zhou¹ Ziha

Zihan Zhang

Simon S. Du¹

¹University of Washington

Preliminaries

Markov Decision Processes

- S states
- A actions
- Planning horizon H
- Reward function $R_h(s, a) \in \Delta([0, 1])$
- Transition probability $P_h(s'|s,a)$
- Maximum transition support $\Gamma = \max_{h,s,a} \|P_h(\cdot|s,a)\|_0$

Policy and value functions

- $\pi = {\{\pi_h\}_{h \in [H]}}$ where $\pi_h : \mathcal{S} \to \mathcal{A}$, optimal π^*
- Value and Q functions

$$V_h^{\pi}(s) := \mathbb{E}_{\pi} \left[\sum_{t=h}^{H} r_t \middle| s_h = s \right],$$

$$Q_h^{\pi}(s, a) := \mathbb{E}_{\pi} \left[\sum_{t=h}^{H} r_t \middle| (s_h, a_h) = (s, a) \right].$$

Episodic reinforcement learning

- K episodes, with policies $\pi^1, \pi^2, \dots, \pi^k$
- Performance measure

$$\mathsf{Regret}(K) := \sum_{k=1}^K (V_1^{\star}(s_1^k) - V_1^{\pi^k}(s_1^k)).$$

Conditions for MDPs

Condition 1: For any possible trajectory, its total reward in a single episode is upper-bounded by 1.

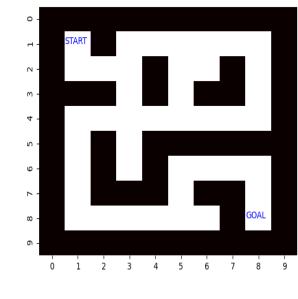
Condition 2: The MDP is time-homogeneous, i.e., the transition and reward are both independent on the timestep h.

Motivation

RL environments have different randomness:







Games with random environments

Robotics

Maze

Stochastic

High variance



Deterministic

Zero variance

Can we design an algorithm which automatically exploits randomness?

Total Multi-Step Conditional Variance

For trajectory $\tau = \{s_h, a_h\}_{h \in [H]}$, define

$$\operatorname{Var}_{\tau}^{\Sigma} := \sum_{h=1}^{H} (\mathbb{V}(R_h(s_h, a_h)) + \mathbb{V}(P_{s_h, a_h, h}, V_{h+1}^{\star})).$$

Let the trajectory of the k-th episode be τ^k , then we denote $\mathrm{Var}_{(k)}^\Sigma:=\mathrm{Var}_{\tau^k}^\Sigma$, and $\mathrm{Var}_K^\Sigma:=\sum_{k=1}^K\mathrm{Var}_{(k)}^\Sigma$.

Maximum Policy-Value Variance

For deterministic policy π , define

$$\mathsf{Var}_1^\pi(s) := \mathbb{E}_\pi \left[\left. \sum_{h=1}^H \left(\mathbb{V}(R_h(s_h, a_h)) + \mathbb{V}(P_{s_h, a_h, h}, V_{h+1}^\pi) \right) \, \right| \, s_1 = s \right]$$

and $\operatorname{Var}^{\star} := \max_{\pi \in \Pi, s \in \mathcal{S}} \operatorname{Var}_{1}^{\pi}(s)$.

Discussion on Variances

- Under Condition 1, we have $\operatorname{Var}_{\tau}^{\Sigma} \leqslant \widetilde{O}(1)$ with high probability if τ is sampled by any policy and $\operatorname{Var}^{\star} \leqslant 1$. Without Condition 1, $\operatorname{Var}_{\tau}^{\Sigma} \leqslant \widetilde{O}(H^2)$ and $\operatorname{Var}^{\star} \leqslant H^2$.
- For deterministic MDPs, $Var_{\tau}^{\Sigma} = Var^{\star} = 0$.
- $Var^* = 0 \implies Var^{\Sigma}_{\tau} = 0$, while reverse is not true.

Results of Model-Based Algorithms

The results are under Condition 1 and Condition 2

Algorithm	Regret	Variance- Dependent	Stochastic- Optimal	Deterministic Optimal
Euler	$\widetilde{O}(\sqrt{H\mathbb{Q}^{\star} \cdot SAK} + H^{5/2}S^2A)$	Yes	No	No
	$\widetilde{O}(\sqrt{SAK} + H^{5/2}S^2A)$	No	Yes	No
MVP	$\widetilde{O}(\sqrt{SAK} + S^2A)$	No	Yes	No
MVP-V	$\widetilde{O}(\sqrt{\min\{Var_K^\Sigma,Var^\star K\}SA}+\Gamma SA)$	Yes	Yes	Yes
This work	$O(\sqrt{\min\{\text{var}_K,\text{var}(K\}\mathcal{O}A+1\mathcal{O}A\}})$			
Lower bound	$\Omega(\sqrt{SAK})$ / $\Omega(SA)$	_	_	-

Remark:

- \mathbb{Q}^* could be as large as $\Omega(1)$ and $H\mathbb{Q}^* \geqslant \mathsf{Var}_{\tau}^{\Sigma}$.
- MVP-Vrecovers the optimal minimax rate because $Var^* \leq O(1)$, and is optimal for deterministic MDPs because variances are 0 and $\Gamma = 1$.

Results of Model-Free Algorithms

Algorithm	Regret	Variance- Dependent	Stochastic- Optimal
Q-learning (UCB-B)	$\widetilde{O}(\sqrt{H^4SAK} + H^{9/2}S^{3/2}A^{3/2})$	No	No
UCB-Advantage	$\widetilde{O}(\sqrt{H^3SAK} + \sqrt[4]{H^{33}S^8A^6K})$	No	Yes
Q-EarlySettled- Advantage	$\widetilde{O}(\sqrt{H^3SAK} + H^6SA)$	No	Yes
UCB-Advantage-V This work	$\widetilde{O}(\sqrt{\min\{Var_K^\Sigma,Var^{m{ imes}}K\}}HSA + \sqrt[4]{H^{15}S^5A^3K})$	Yes	Yes
Lower bound	$\Omega(\sqrt{H^3SAK})$ / $\Omega(H^2SA)$	_	_

Remark: Currently no generic model-free algorithm achieves a constant regret for deterministic MDPs, while ours achieves a $K^{1/4}$ rate.