

Oct. 5, 23

Lew8: LSVI-UCB in linear MDP

• RL problem

$$M = (S, A, \{r_h\}_{h \in [H]}, \{P_h\}_{h \in [H]})$$

interaction protocol

For episode $k = 1, \dots, K$:

– the adversary picks an initial state x_1^k

– the learner plays π^k (based on all information the learner has observed so far)

After K episodes, the regret of the learner:

$$\text{Regret}(K) = \sum_{k=1}^K [V_1^*(x_1^k) - V_1^{\pi^k}(x_1^k)]$$

recall:

$$V_1^{\pi}(x_1) = \mathbb{E}_{\pi} \left[\sum_{h=1}^H r_h(x_h, a_h) \right] \text{ where } (x_1, a_1, \dots, x_H, a_H) \sim (\{P_h\}_{h \in [H]}, \pi)$$

$$V_1^*(x_1) = \max_{\pi} V_1^{\pi}(x_1)$$

Linear MDP

$$\exists \phi: S \times A \rightarrow \mathbb{R}^d$$

$$\theta_h \in \mathbb{R}^d$$

$$\mu_h(s) = (\mu_h^{(1)}(s), \dots, \mu_h^{(d)}(s))$$

s.t.

$$\begin{cases} P_h(x'|x, a) = \phi(x, a)^T \mu_h(x') \\ r_h(x, a) = \phi(x, a)^T \theta_h \end{cases}$$

Chi Jin, Zhuoran Yang, Zhaoran Wang, Michael Jordan.

"Probably efficient RL with linear function approximation." COH'2020

Least-square Value Iteration w/ UCB CLSVI-UCB)

For episode $k=1, \dots, K$

- receive initial state x_1^k

- For step $h=H, \dots, 1$ do

$$\Lambda_h^k \leftarrow \sum_{\tau=1}^{k-1} \phi(x_h^\tau, a_h^\tau) \phi(x_h^\tau, a_h^\tau)^T + \lambda I$$

$$W_h^k \leftarrow (\Lambda_h^k)^{-1} \sum_{\tau=1}^{k-1} \phi(x_h^\tau, a_h^\tau) \left[r_h(x_h^\tau, a_h^\tau) + \max_a Q_{h+1}(x_{h+1}^\tau, a) \right]$$

$$= \underset{W \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{\tau=1}^{k-1} \left(\phi(x_h^\tau, a_h^\tau)^T W - r_h(x_h^\tau, a_h^\tau) - \max_a Q_{h+1}(x_{h+1}^\tau, a) \right)^2$$

$$Q_h^k(\cdot, \cdot) \leftarrow \min \left\{ \phi(\cdot, \cdot)^T W_h^k + \underbrace{\beta \|\phi(\cdot, \cdot)\|_{(\Lambda_h^k)^{-1}}}_\text{estimation error}, H \right\}$$

• encourage to take actions w/ large estimation error
→ exploration

- for step $h=1, \dots, H$:

take action $a_h^k \leftarrow \underset{a}{\operatorname{argmax}} Q_h(x_h^k, a)$ and observe $x_{h+1}^k \sim P_h(\cdot | x_h^k, a_h^k)$

• if all actions are well explored, the est. error is small → exploitation

Theorem: If choose $\beta = \tilde{O}(dH)$, then $\text{regret}_{\text{LSVI-UCB}}(k) = \tilde{O}(H^2 d^{3/2} \sqrt{k})$

| | |
|--------------------------------------|-------------------------------|
| LSVI-UCB (Jin et al. COLT'20) | $\tilde{O}(\sqrt{d^3 H^4 K})$ |
| Lower bound (Zhou et al. COLT'21) | $\Omega(d\sqrt{H^3 K})$ |
| LSVI-UCB++ (He et al. ICML'23) | $\tilde{O}(d\sqrt{H^3 K})$ |

$O(\sqrt{dH})$ gap

variance-aware LSVI to shave off \sqrt{H} factor

reference-advantage decomposition
 $(V_{hH}^K = V_{hH}^* + V_{hH}^K - V_{hH}^*)$

to shave off \sqrt{d} factor

Zhou et al. COLT'20: Nearly minimax optimal RL for linear mixture MDP

He et al., ICML'23: Nearly minimax optimal RL for linear MDP

Note: Q-learning paper of Chi Jin only gives lower bound of order $\Omega(\sqrt{dH^3 K})$ if we set $d = SA$

Sketch proof of the theorem

$$\bullet E_1 = \left\{ Q_h^k(x, a) \geq Q_h^*(x, a), \quad \forall (x, a, h, k) \in S \times A \times [H] \times [K] \right\} \quad (\text{optimism})$$

$$\bullet \delta_h^k = V_h^k(x_h^k) - V_h^{\pi_k}(x_h^k)$$

(computation of V_h^k, π_k do not use (x_h^k, a_h^k))

$$S_{h+1}^k = \mathbb{E}[\delta_{h+1}^k \mid x_h^k, a_h^k] - \delta_{h+1}^k$$

$$E_2 = \left\{ \delta_h^k \leq \delta_{h+1}^k + S_{h+1}^k + 2\beta \|\phi_h^k\|_{(\Lambda_h^k)^{-1}} : \forall (k, h) \right\}$$

$$\text{Regret}(K) = \sum_{k=1}^K [V_1^k(x_1^k) - V_1^{\pi_k}(x_1^k)]$$

$$\stackrel{\mathbb{E}_1}{\leq} \sum_{k=1}^K [V_1^k(x_1^k) - V_1^{\pi_k}(x_1^k)]$$

$$= \sum_{k=1}^K \delta_1^k \stackrel{\mathbb{E}_2}{\leq} \sum_{k,h} S_h^k + 2\beta \sum_{k,h} \|\Phi_h^k\|_{(\Lambda_h^k)^{-1}}$$

• $\{S_h^k\}_{k,h}$ is a martingale difference sequence. and $|S_h^k| \leq 2H$

Azuma-Hoeffding's inequality

$$\sum_{k,h} S_h^k \leq \sqrt{KH^2 \log(1/\delta)} \quad \text{w.p. } 1-\delta$$

$$\sum_{k=1}^K \|\Phi_h^k\|_{(\Lambda_h^k)^{-1}}^2 \leq 2 \log \frac{\det(\Lambda_h^{K+1})}{\det(\Lambda_h^1)} = O(d \log K)$$

recall: $\Lambda_h^k = \sum_{\tau=1}^{k-1} \Phi_h^\tau (\Phi_h^\tau)^\top + \lambda I$

$$\sum_{k,h} \|\Phi_h^k\|_{(\Lambda_h^k)^{-1}} \leq \sum_h \sqrt{K \sum \|\Phi_h^k\|_{(\Lambda_h^k)^{-1}}^2} \leq H \sqrt{Kd \log K}$$

Sketch proof for $\Pr(E_2)$ is large

• Concentration lemma

$$\forall (k, h): \left\| \sum_{\tau=1}^{k-1} \phi(x_h^\tau, a_h^\tau) [V_{h+1}^k(x_{h+1}^\tau) - P_h V_{h+1}^k(x_h^\tau, a_h^\tau)] \right\|_{(\Lambda_h^k)^{-1}} \leq \tilde{O}(dH)$$

Fix π ,

$$\forall (x, a, h, k): \phi(x, a)^T W_h^k - Q_h^\pi(x, a) \leq P_h [V_{h+1}^k - V_{h+1}^\pi](x, a) + \underbrace{dH}_{\beta} \|\phi(x, a)\|_{(\Lambda_h^k)^{-1}}$$

E_2

E_1

Sketch proof for $\Pr(E_1)$ is large

$$\cdot \quad \phi(x, a)^T W_H^k - Q_H^*(x, a) \leq \beta \|\phi(x, a)\| (\Lambda_H^k)^{-1}$$

$$\cdot \quad \phi(x, a)^T W_h^k - Q_h^*(x, a) - \underbrace{P_h(V_{hH}^k - V_{hH}^*)}_{\leq 0 \text{ by induction}}(x, a) \leq \beta \|\phi(x, a)\| (\Lambda_h^k)^{-1}$$