

Linear MDP $(S, A, \{P_h\}, \{R_h\}, H, d_0)$.

$$P_h(s'|s, a) = \phi_h(s, a)^T \psi(s') \rightarrow \mathbb{R}^d$$

$$R_h(s, a) = \phi_h(s, a)^T \theta_R$$

$$\|\phi\|_2 \leq 1$$

$$R_h \in [0, 1]$$

$\Rightarrow \forall f_{h+1}, \mathcal{T}f_{h+1}$ is linear in ϕ_h . $V_{max} = H$.

Alg (UCB-LSVI): for $t=1, 2, \dots, T$

for $h = \underline{H}, H-1, \dots, 1$.

1. Define: $\Lambda_h^t = I + \sum_{i=1}^{t-1} \phi_h^i (\phi_h^i)^T$

where $\phi_h^i := \phi_h(s_h^i, a_h^i)$ \rightarrow episode index.

2. $\tilde{Q}_h^t(s, a) = \phi_h(s, a)^T (\Lambda_h^t)^{-1} \sum_{i=1}^{t-1} \phi_h^i (r_h^i + \tilde{V}_{h+1}^t(s_{h+1}^i))$

3. $\hat{Q}_h^t(s, a) = \tilde{Q}_h^t(s, a) + \beta \sqrt{\phi_h(s, a)^T (\Lambda_h^t)^{-1} \phi_h(s, a)}$
 (optimism)

$\hat{V}_h^t(s, a) = \min \{ H, \max_a \hat{Q}_h^t(s, a) \}$

4. interact w/ MDP using π^t greedy w.r.t. \hat{Q}_h^t .

$$\text{Regret}_T := \sum_{t=1}^T (J(\pi^*) - J(\pi^t)).$$

Lemma. $\forall h, t$. Define $b_h^t(s, a) :=$

$$\left| \underbrace{\bar{Q}_h^t(s, a)}_{(\text{w.h.p.}) \Delta} - \underbrace{\left(R_h(s, a) + P_h(s, a)^T \hat{V}_{h+1}^t \right)}_{\beta} \right|$$

$$b_h^t(s, a) \leq \underbrace{\|\phi(s, a)\|_{(L_h^t)^{-1}}}_{\Delta} \tilde{O}(H(d, t, \sqrt{\log \frac{1}{\delta}})).$$

$$\Rightarrow \hat{Q}_h^t(s, a) \geq R_h(s, a) + P_h(s, a)^T \hat{V}_{h+1}^t.$$

$\bar{X}_1, \dots, \bar{X}_n, \Delta.$

$$\hat{\mathbb{E}}[f(x_i)] \rightarrow \mathbb{E}[f(x_i)]$$

$$f = f(x_1, \dots, x_n) \in \mathcal{F}.$$

$$x \in \mathbb{R}^d. \quad \mathcal{F} = \left\{ \overbrace{f(x) = x^T A x} \right.$$

$$x^T \theta. \quad \left. \forall A \in \mathbb{R}^{d \times d} \right\}$$

$$x^T A x = \sum_{i,j} x_i x_j A_{ij}$$

Lemma: Optimism: $\forall h, t \quad J(\pi^*) = \mathbb{E}_{d_0}[V_1^*]$

$$\rightarrow \boxed{\hat{Q}_h^t \geq Q_h^*, \quad \hat{V}_h^t \geq V_h^*} \quad \checkmark$$

Proof: $h = H+1 \quad \checkmark$

$$\boxed{\hat{Q}_h^t(s, a) \geq R_h(s, a) + P_h(s, a)^T \hat{V}_{h+1}^t}$$

$$\geq R_h + P_h^T V_{h+1}^*$$

$$= Q_h^*(s, a)$$

$$\hat{V}_h(s) = \min(H, \max_a \hat{Q}_h^t(s, a))$$

$$\downarrow$$

$$V_h^*(s)$$

$$\max_a \hat{Q}_h^t(s, a)$$

$$\parallel$$

$$V_h^*(s)$$

$$\geq V_h^*(s)$$

$$\text{Regret}_T = \sum_{t=1}^T \underbrace{J(\pi^*) - J(\pi^t)}$$

$$\leq \sum_{t=1}^T \mathbb{E}_{d_0} [\hat{Q}_1^t(s_1, \pi^t)] - J(\pi^t)$$

$$\leq \sum_{t=1}^T \sum_{h=1}^H \mathbb{E}_{d_n^{\pi^t}} [\hat{Q}_n^t(s, a) - R_n(s, a) - P_n(s, a)^T [\max_a \hat{Q}_{n+1}^t(\cdot, a)]]$$

$\min(t)$

$$\leq \sum_{t=1}^T \sum_{h=1}^H \mathbb{E}_{d_n^{\pi^t}} [\hat{Q}_n^t(s, a) + \beta \cdot \|\phi(s, a)\|_{(V_h^t)^{-1}} - R_n - P_n^T \hat{V}_{n+1}]$$

$$\leq \sum_{t=1}^T \sum_{h=1}^H \mathbb{E}_{d_n^{\pi^t}} [2\beta \|\phi(s, a)\|_{(V_h^t)^{-1}}]$$

$$= 2\beta \sum_{h=1}^H \sum_{f=1}^T \mathbb{E}_{d_h^{\pi^+}} \left[\underbrace{\|\phi(s, a)\|}_{\Delta} \underbrace{(\Delta_h^+)^{-1}}_{\Delta_h^+} \right].$$

$$\leq 2\beta \sum_{h=1}^H \left(\underbrace{\sum_{f=1}^T \|\phi_n^+(s_n^+, a_n^+)\| (\Delta_h^+)^{-1}}_{\text{Azuma's.}} \right) \underbrace{\mathbb{E}[\dots]}_{\substack{\text{elliptical} \\ \text{potential.}}} \left| \begin{matrix} s_n^+, a_n^+ \\ \vdots \\ s_n^+, a_n^+ \end{matrix} \right.$$