

- HW2 (due 2 weeks)
- presentation = critique form
- please send me your slides ^{notes} after your presentation.

Linear MDPs: $M = (\Gamma_h, P_h)_{h=0}^{H-1}$. feature map $\phi: S \times A \rightarrow \mathbb{R}^d$.

$$\exists (\alpha_h^*)_{h=0}^{H-1} \subseteq \mathbb{R}^d.$$

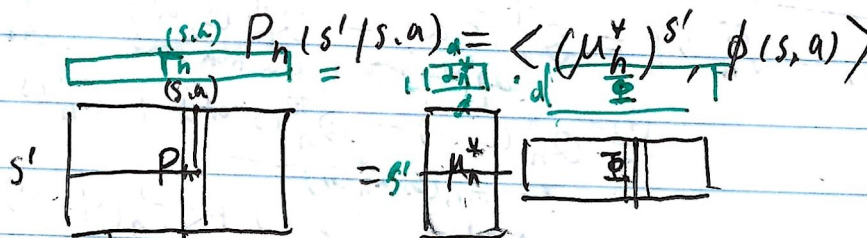
$$(\mu_h^*)_{h=0}^{H-1} \subseteq \mathbb{R}^{S \times d}, \text{ such that:}$$

(1) $\forall h, \forall s, a, s': \Gamma_h(s, a) = \langle \alpha_h^*, \phi(s, a) \rangle$

(2) $\max_{s, a} \|\phi(s, a)\|_2 \leq V$

(3) $\forall V: \|V\|_{\infty} \leq V \Rightarrow \|V^T \mu_h^*\|_2 \leq \sqrt{d} V$

(4) $\forall h, \|\alpha_h^*\|_2 \leq V$



(exercise)

Recall: Linear MDPs satisfies that: $\forall V: S \rightarrow \mathbb{R}$.

$$L_h: \mathbb{R}^S \rightarrow \mathbb{R}^{S \times A}$$

$$(L_h V)(s, a) = \Gamma_h(s, a) + \langle P_h(\cdot | s, a), V \rangle$$

[L_h links V_{h+1}^π to Q_h^π
 V_{h+1}^* to Q_h^*]

$$= \langle \alpha_h^*, \phi(s, a) \rangle + \langle \mu_h^*, \phi(s, a), V \rangle$$

$$= (\alpha_h^*)^T \phi(s, a) + V^T \cdot \mu_h^* \phi(s, a)$$

important notation.

$$= \langle \alpha_h^* + (\mu_h^*)^T V, \phi(s, a) \rangle = \langle \Theta_h^*(V), \phi(s, a) \rangle$$

(linear in $\phi(s, a)$.)

How to design an online RL algorithm for linear MDPs?

Idea: combine LSVI with optimism. (useful in online RL)

For episodes $k = 0, 1, \dots, K$:

$$\text{// Data: } D_h = \left\{ (S_h^i, a_h^i, r_h^i, S_{h+1}^i) \right\}_{i=0}^{R-1}, h = 0, 1, \dots, H-1.$$

~~For~~

// Idea 1: construct optimistic estimates of Q_h^*, V_h^*

$$Q_h^k \leq 0.$$

For $h = H-1, H-2, \dots, 0$:

|| given $\hat{V}_{n+1}^k \geq V_{n+1}^*$. want to construct $\hat{Q}_n^k \geq Q_n^*$.

write on another sheet of board) Recall: $Q_n^*(s,a) = (L_h V_{n+1}^*)(s,a)$
 can we build estimate of $L_h V_{n+1}^*$?

although we don't have $L_h V_{n+1}^*$ as $L_h \hat{V}_{n+1}^k \geq L_h V_{n+1}^*$, (exercise).
 we have \hat{V}_{n+1}^k as an upper bound of it. it suffices to build tight upper bound estimates of

$$L_h \hat{V}_{n+1}^k$$

Idea: $L_h \hat{V}_{n+1}^k = \langle \theta_h^*(\hat{V}_{n+1}^k), \phi(s,a) \rangle$. so it suffices to estimate $\theta_h(\hat{V}_{n+1}^k)$; specifically.

Similar to UCB. we will use a confidence set (generalizing confidence interval) to estimate it.

Define $\hat{\theta}_h^k = \underset{\theta}{\operatorname{argmin}} \left\{ \frac{1}{n} \|\theta\|^2 + \sum_{i=0}^{k-1} \frac{\langle \theta, \phi(S_h^i, a_h^i) \rangle - (r_h^i + \hat{V}_{n+1}^k(S_{n+1}^i))}{(r_h^i + \hat{V}_{n+1}^k(S_{n+1}^i))} \right\}$
 ridge regression

Denote by $\phi_n^i = \phi(S_h^i, a_h^i)$.

$$\hat{\theta}_h^k = (L_h)^{-1} \left(\sum_{i=0}^{k-1} \phi_n^i (r_h^i + \hat{V}_{n+1}^k(S_{n+1}^i)) \right)$$

$$\text{where } \Lambda_h^k = kI + \sum_{i=0}^{k-1} \phi_n^i (\phi_n^i)^T$$

set $B_h^k = \{ \theta : \|\theta - \hat{\theta}_h^k\|_{\Lambda_h^k} \leq \beta \}$, where

(with prob $1 - \frac{1}{k}$: $\theta_h^*(\hat{V}_{n+1}^k) \in B_h^k$)

$$\beta = \text{const.} \cdot \sqrt{\ln \frac{nk}{\delta}}$$

(so w.h.p. $\geq \langle \theta_h^*(\hat{V}_{n+1}^k), \phi(s,a) \rangle$)

Define $\hat{Q}_n^k(s,a) = \max_{\theta \in B_h^k} \langle \theta, \phi(s,a) \rangle = (L_h \hat{V}_{n+1}^k)(s,a)$

$$= \max_{\|\Delta\theta\|_{\Lambda_h^k} \leq \beta} \langle \hat{\theta}_h^k, \phi(s,a) \rangle + \langle \Delta\theta, \phi(s,a) \rangle$$

(see textbook for another "model-based" ^{explanation} ~~description~~ of the ^{same} algorithm)

$$M^k = (\hat{P}_h^k, \hat{r}_h^k)_{h=0}^H$$

$$= \langle \hat{\Theta}_h^k, \phi(s, a) \rangle + \beta \cdot \|\phi(s, a)\| (\Lambda_h^k)^{-1}$$

— Define $\hat{V}_h^k(s) = \min_{a \in A} (\max_{a \in A} \hat{Q}_h^k(s, a), H)$

— Define $\hat{\pi}_h^k(s) = \operatorname{argmax}_{a \in A} \hat{Q}_h^k(s, a)$

// execute policy $\hat{\pi}^k$.

For $h = 0, 1, \dots, H-1$:

see s_h^k . take action $a_h^k = \hat{\pi}_h^k(s_h^k)$,

receive $r_h^k = r_h(s_h^k, a_h^k)$, transition to S_{h+1}^k .

LSVI-UCB algorithm.

Analysis:

Theorem: LSVI-UCB has regret

$$\operatorname{Reg}(K) \leq \tilde{O}(H^2 \sqrt{d^3 K})$$

$\tilde{O}(\cdot)$ ignore log factors

pf structure: ① ~~optimism: $V_h^k \geq V_h^*$~~ ~~concentration: $\forall k, E_k = \{ \forall h, \Theta_h^*(\hat{V}_{h+1}^k) \in \mathcal{B}_h^k \}$~~ whp

② $P(E_k) \geq 1 - \frac{1}{k}$

③ optimism: $\forall k$ on E_k , $\hat{V}_h^k \geq V_h^*$, $\hat{Q}_h^k \geq Q_h^*$

$$\textcircled{3} \operatorname{Reg}(K) \leq \mathbb{E} \left[\sum_{k=0}^{K-1} (V_0^*(s_0^k) - V_0^{\pi^k}(s_0^k)) \mathbb{I}(E_k) \right]$$

$$\stackrel{\text{optimism}}{\leq} \mathbb{E} \left[\sum_{k=0}^{K-1} (\hat{V}_0^k(s_0^k) - V_0^{\pi^k}(s_0^k)) \mathbb{I}(E_k) \right]$$

$$\leq \mathbb{E} \left[\sum_{k=0}^{K-1} \sum_{h=0}^{H-1} 2\beta \cdot \|\phi(s_h^k, a_h^k)\| (\Lambda_h^k)^{-1} \mathbb{I}(E_k) \right]$$

bonus / uncertainty of experienced state-action pairs.

$$\textcircled{4} \quad \sum_{k=0}^{K-1} \sum_{h=0}^{H-1} \|\phi(s_h^k, a_h^k)\| (\Lambda_h^k)^{-1} \leq \tilde{O}\left(\sqrt{dK}\right)$$

using elliptic potential lemma (will introduce next)

In summary, $\textcircled{1} - \textcircled{4}$.

$$\Rightarrow \text{Reg}(K) \leq \tilde{O}\left(\beta \cdot H \sqrt{dK}\right) = \tilde{O}\left(H^2 \sqrt{d^3 K}\right)$$

Pf of $\textcircled{1}$: Let $f = \hat{V}_{n+1}^k$.

$$\begin{aligned} (\mathcal{L}_h f)_{(s_h^i, a_h^i)} &= \langle \theta_h^y(f), \phi(s_h^i, a_h^i) \rangle \stackrel{\text{Expectation}}{\leftarrow} r_h^i + f(s_{n+1}^i) \\ &= \langle \theta_h^y(f), \phi_h^i \rangle + \varepsilon_h^i(f) \\ \hat{\theta}_h^k(f) &= (\Lambda_h^k)^{-1} \left(\sum_{i=1}^{k-1} \phi_h^i (\phi_h^{iT} \theta_h^y(f) + \varepsilon_h^i(f)) \right) \end{aligned}$$

$$\Rightarrow \hat{\theta}_h^k(f) - \theta_h^y(f)$$

$$= (\Lambda_h^k)^{-1} \left[\left((\Lambda_h^k) - \lambda \mathbf{I} \right) \theta_h^y(f) + \sum_{i=0}^{k-1} \phi_h^i \varepsilon_h^i(f) \right] - \theta_h^y(f)$$

$$= -\lambda^{-1} \theta_h^y(f) + (\Lambda_h^k)^{-1} \left(\sum_{i=0}^{k-1} \phi_h^i \varepsilon_h^i(f) \right)$$

$$\Rightarrow \|\hat{\theta}_h^k(f) - \theta_h^y(f)\|_{\Lambda_h^k}$$

$$\leq \sqrt{\lambda} \|\theta_h^y(f)\| + \left\| \sum_{i=0}^{k-1} \phi_h^i \varepsilon_h^i(f) \right\|_{(\Lambda_h^k)^{-1}}$$

$\hat{\theta}_h^k(f)$ is the ridge regression soln of

$$\left(\phi(s_h^i, a_h^i), (\mathcal{L}_h f)(s_h^i, a_h^i) + \varepsilon_i(f) \right)_{i=0}^{k-1}$$

regression quantiles $\langle \theta_h^y(f), \phi(s_h^i, a_h^i) \rangle$

$$\begin{aligned} \|\hat{\theta}_h^k(f) - \theta_h^y(f)\|_{(\Lambda_h^k)^{-1}} &\leq \|\theta_h^y(f)\|_2 + H \sqrt{d + \ln \frac{1}{\delta}} \\ &\leq \tilde{O}\left(H \sqrt{d \ln \frac{1}{\delta}}\right) \quad (*) \end{aligned}$$

However: $f = \hat{V}_{h+1}^k$ depends on previous data collected.
 so conditioned on $\hat{V}_{h+1}^k (S_h^i, a_h^i)_{i=0}^{k-1}$, $(S_{h+1}^i)_{i=0}^{k-1}$ may not be independent.
 So (*) may not be true.

Fix: ~~the~~ key observation: \hat{V}_{h+1}^k come from a "simple" fn class.
 argue that the concentration holds for all $f \in \mathcal{F}_{L, \beta}$.

$$f \in \mathcal{F}_{L, \beta} \iff f \text{ w.p. } \beta. \text{1} : \|W\|_2 \leq L, 0 \leq \beta \leq \beta.$$

with $L = 2HK$
 $\beta = \delta(Hd)$.

$$\text{where } f \text{ w.p. } \beta. \text{1} = \min(H, \max_a \langle W, \phi(s, a) \rangle + \beta \|\phi(s, a)\|_{\Lambda^{-1}})$$

we show w.p. $1 - \delta$. $\forall f \in \mathcal{F}_{L, \beta. \text{1}}$.

$$\|\hat{\Theta}_h^k(f) - \Theta_h^*(f)\|_{(\Lambda_h^k)^{-1}} \leq \delta(H\sqrt{d^2 + \ln \frac{1}{\delta}})$$

via a covering argument. (AJKS Sect. 8.4)

pf of ②: $\forall k$. suppose E_k happens. $\textcircled{1}$ $\forall s$ $\hat{Q}_h^k(s) = \min(H, \max_a \hat{Q}_h^k(s, a))$
 Backward induction. $\textcircled{2}$ $\forall h, s, a$. $\hat{Q}_h^k(s, a) = \max_{\theta \in \mathcal{D}_h^k} \langle \theta, \phi(s, a) \rangle$

$h=H$. trivial.

suppose $\hat{Q}_{h+1}^k \geq Q_{h+1}^*$. $\hat{V}_{h+1}^k \geq V_{h+1}^*$. (inductive hypothesis)

as $\forall s, a$. $\hat{Q}_h^k(s, a) = \max_{\theta \in \mathcal{B}_h^k} \langle \theta, \phi(s, a) \rangle$ $\textcircled{1}$ $\textcircled{2}$ \Rightarrow Strong optimism

$$\geq \langle \Theta_h^*(\hat{V}_{h+1}^k), \phi(s, a) \rangle$$

$$= (\hat{L}_h \hat{V}_{h+1}^k)(s, a)$$

(almost satisfy Bellman opt. eqn w/ equality replaced w/ ineq.)

$$\geq (L_h V_{h+1}^*)(s, a) = Q_h^*(s, a).$$

$$\Rightarrow \hat{Q}_h^k \geq Q_h^*$$

$$\max_a \hat{Q}_h^k(s, a) \geq \max_a Q_h^*(s, a) = V_h^*(s) \text{ and } H \geq V_h^*(s).$$

$$\Rightarrow \hat{V}_h^k(s) \geq V_h^*(s).$$

leave this as an exercise

$\forall s$.

pf of ③ $\hat{V}_0^k(s_0^k) - V_0^{\pi^k}(s_0^k)$

Let $a_0^k = \pi_0^k(s_0^k)$
 $= \min(\hat{Q}_0^k(s_0^k, a_0^k), H) - Q_0^{\pi^k}(s_0^k, a_0^k)$

$\leq \hat{Q}_0^k(s_0^k, a_0^k) - Q_0^{\pi^k}(s_0^k, a_0^k)$

Generalized Simulation Lemma:

given episodic MDP. $M, \hat{Q} = (\hat{Q}_h)_{h=0}^{H-1} : S \times A \rightarrow \mathbb{R}$, $\pi = (\pi_h)_{h=0}^{H-1}$ history independent policy. then

$\hat{Q}_h(s_h, a_h) - Q_h^\pi(s_h, a_h)$

$= \mathbb{E} \left[\sum_{i=h}^{H-1} \text{Bell}_h^\pi(s_h, a_h) \mid s_h, a_h, M, \pi \right]$, where

$\text{Bell}_h^\pi(s, a) = \hat{Q}_h(s, a) - (J_h^\pi \hat{Q}_{h+1})(s, a)$ is the Bellman error of \hat{Q} at step h .
(how much does \hat{Q}_h violate Bellman consistency equation of π)
 $(J_h^\pi f)(s, a) = r_h(s, a) + \sum_{s'} P_h(s'/s, a) f(s, \pi_{h+1}(s))$.

① it generalizes simulation lemma b/c $(\hat{Q}_h)_{h=0}^{H-1}$ may not be value fn of some policy.

② intuitively, if Bellman error of \hat{Q} is small at all ~~steps~~ steps, $\hat{Q} \approx Q^\pi$.

Applying to our setting, at episode k .

$\text{Bell}_h^{\pi^k}(s, a) = \hat{Q}_h^k(s, a) - \left(r_h(s, a) + \sum_{s'} P_h(s'/s, a) \hat{Q}_{h+1}^k(s, \pi_{h+1}^k(s)) \right)$

$\leq \hat{Q}_h^k(s, a) - (L_h \hat{V}_{h+1}^k)(s, a)$

$= \max_{\theta \in B_h^k} \langle \theta, \phi(s, a) \rangle - \langle \theta_h^*(\hat{V}_{h+1}^k), \phi(s, a) \rangle$

$\|\theta - \theta_h^*\|_{\Lambda_h^k} \leq 2\beta \cdot \|\phi(s, a)\|_{(\Lambda_h^k)^{-1}}$

max $\hat{Q}_{h+1}^k(s, a)$
 $\approx \hat{V}_{h+1}^k(s)$

in other words.

$$\left(\Lambda_0^k (s_0^k, a_0^k) - Q_0^{\pi^k} (s_0^k, a_0^k) \right) I(E_k)$$

$$\leq \mathbb{E} \left[\sum_{h=0}^{H-1} 2\beta \cdot \|\phi(s_h^k, a_h^k)\| (\Lambda_h^k)^{-1} \left| \mathcal{H}^{k-1} \right. \right] I(E_k)$$

$$\left\{ (s_h^{k'}, a_h^{k'})_{h=0}^{H-1}, s_0^k, a_0^k \right\}$$

Pf of ④: $\sum_{k=0}^{K-1} \sum_{h=0}^{H-1} \|\phi(s_h^k, a_h^k)\| (\Lambda_h^k)^{-1}$

$$\Lambda_h^k = I + \sum_{k'=0}^{k-1} \phi_h^{k'} (\phi_h^{k'})^T$$

so we just need to bound

(*) $\sum_{k=0}^{K-1} \|\phi_h^k\| (\Lambda_h^k)^{-1}$ for each h separately

The Elliptic potential lemma (EPL)

suppose $u_1, \dots, u_T \in \mathbb{R}^d$. $V_t = \mu I + \sum_{s=1}^t u_s u_s^T$. then

in 1-d:

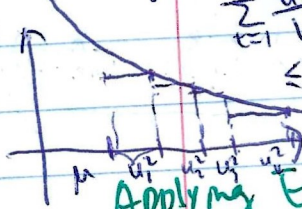
$$V_t = \mu + \sum_{s=1}^t u_s^2$$

$$\sum_{t=1}^T \|u_t\|_{V_t^{-1}}^2 \leq \ln \frac{\det(V_T)}{\det(V_0)}$$

If $\forall t \|u_t\| \leq D$.

$$\leq \int_{\mu}^{\mu + \frac{D^2 T}{2}} \frac{1}{x} dx$$

$$\leq d \ln \left(1 + \frac{D^2 T}{d\mu} \right) = \delta(d)$$



Applying EPL to (*):

$$(*) \leq \sqrt{K} \cdot \sum_{k=0}^{K-1} \underbrace{\|\phi_h^k\|}_{u_t} \underbrace{(\Lambda_h^k)^{-1}}_{V_{t+1}, \text{ not } V_t!}$$

$$\Lambda_h^{k+1} = \Lambda_h^k + \phi_h^k (\phi_h^k)^T$$

$$\leq \Lambda_h^k + I$$

$$\leq 2\Lambda_h^k \quad \text{(Exercise)}$$

$$B \Rightarrow A^T \geq B^T$$

$$\Rightarrow (\Lambda_h^k)^{-1} \leq 2(\Lambda_h^{k+1})^{-1}$$

$$\leq \sqrt{K \sum_{k=0}^{K-1} \|\phi_h^k\|^2 (\Lambda_h^{k+1})^{-1}}$$

$$\leq \tilde{O}(\sqrt{Kd})$$

pf of generalized simulation lemma:

$$h, s_n, a_n: \hat{Q}_h(s_n, a_n) - Q_h^\pi(s_n, a_n)$$

$$\text{Introduce Bellman error}$$

$$= \hat{Q}_h(s_n, a_n) - (J_h^\pi \hat{Q}_{h+1})(s_n, a_n) + (J_h^\pi \hat{Q}_{h+1})(s_n, a_n) - (J_h^\pi Q_{h+1}^\pi)(s_n, a_n)$$

$$= \text{Bell}_h^\pi(s_n, a_n) + \left[r_h(s_n, a_n) + \underbrace{\langle \hat{Q}_{h+1}, \hat{Q}_{h+1} \rangle}_{\sum_{s'} P_h(s'|s_n, a_n) \cdot \hat{Q}_{h+1}(s', \pi_{h+1}(s'))} \right]$$

$$- \left[r_h(s_n, a_n) + \sum_{s'} P_h(s'|s_n, a_n) \cdot Q_{h+1}^\pi(s', \pi_{h+1}(s')) \right]$$

$$= \text{Bell}_h^\pi(s_n, a_n) + \mathbb{E} \left[\hat{Q}_{h+1}(s_{n+1}, a_{n+1}) - Q_{h+1}^\pi(s_{n+1}, a_{n+1}) \mid s_n, a_n, \pi, M \right]$$

keep expanding