## Linear MDP w/ a simulator

Recap:

$$\#\text{samples} = O\left(\frac{pdeg(H) \cdot |S||A|}{\varepsilon^2}\right)$$

→ linear in #numbers of states

What if: #states are exponentially large? (e.g. Atari games)

↓

we almost cannot visit one state twice

Question: How to generalize from observed states to unobserved states?

↑

similarity btw states

↑

function approx + (assumptions)

↑

simplest model: linear MDP

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

time-inhomogeneous

$$r_h: S \times A \to \mathbb{R} \quad, \quad P_h(\cdot | s, a) \in \triangle(S)$$

Linear MDP:- $\phi: S \times A \longrightarrow \mathbb{R}^d$ is a known feature map

$$\forall h \in [H], \exists w_h \in \mathbb{R}^d: \quad r_h(s, a) = \langle \phi(s, a), \theta_h \rangle \quad \forall (s, a)$$

$$\exists \mu_h \in \{S \to \mathbb{R}^{d}\}: \quad P_h(s' | s, a) = \langle \phi(s, a), \mu_h(s') \rangle$$

e.g. tabular MDP: $d = |S| \cdot |A|$, $\phi(s, a) = e_{(s,a)}$

Lemma: $Q_h^\pi$ is linear in $\phi$ $\quad \forall h, \pi$

$$\forall \pi, \exists \quad w(\pi) \in (\mathbb{R}^d)^H \quad Q_h^\pi(s, a) = \langle \phi(s, a), w_h(\pi) \rangle$$

Linear MDP w/ simulator

$(s, a, h) \longrightarrow \boxed{\text{sim}} \longrightarrow s' \sim P_h(\cdot | s, a)$

assume $P_h$ is known for simplicity

how to estimate the optimal policy?

Least-square value iteration (LSVI):

- Construct a "core set" $K \in S \times A$ of state-action pairs

- Collect data w/ the simulator:

for $h = 1:H$ do
For each $(s, a) \in K$,
query $(s, a, h)$ $n$ times $\longrightarrow$ $s_1', \ldots, s_n' \overset{iid}{\sim} P_h(\cdot | s, a)$
Add $\{(s, a, s_i')\}_{i \in [n]}$ to $D_h$

- Backup estimation:

$$\hat{V}_{H+1}(s) = 0 \quad \forall_s$$

for $h = H, H-1, \ldots, 1$ :

$$\hat{w}_h = \underset{w \in \mathbb{R}^d}{\arg\min} \sum_{(s,a,s') \in D_h} \left( \phi(s,a)^T w - \hat{r}_h(s,a) - \hat{V}_{h+1}(s') \right)^2$$

$$\hat{Q}_h = \phi^T \hat{w}_h$$

$$\hat{V}_h(s) = \max_a \hat{Q}_h(s,a)$$

$$\hat{\pi}_h(s) \in \arg\max_a \hat{Q}_h(s,a)$$

Return: $\hat{\pi} = \{\hat{\pi}_h\}_{h \in [H]}$

Assume: $\text{span} \{ \phi(s,a) \mid (s,a) \in S \times A \} = \mathbb{R}^d$

core set

$$K = \{ (\bar{s}_i, \bar{a}_i) \}_{i \in [d]} \qquad \text{span} \{ \phi(\bar{s}_i, \bar{a}_i) \mid (\bar{s}_i, \bar{a}_i) \in K \} = \mathbb{R}^d$$

<span style="color:orange">D-optimal design</span>

$$\# \text{samples} = d \cdot n \cdot H$$

Goal: Bound $V_1^{\pi^*}(s) - V_1^{\hat{\pi}}(s)$

Define: $\Lambda = \sum_{i \in [d]} \phi(\bar{s}_i, \bar{a}_i) \phi^T(\bar{s}_i, \bar{a}_i)$

<span style="color:orange">Dirac distribution</span>

$$= \frac{1}{n} \sum_{\substack{(\bar{s}, \bar{a}) \\ \in D_n}} \phi(\bar{s}, \bar{a}) \phi^T(\bar{s}, \bar{a})$$

$$\Psi = ( \phi(\bar{s}_i, \bar{a}_i) )_{i \in [d]} \in \mathbb{R}^{d \times d}$$

investigate the least-square solution:

$$\hat{w}_h = \underset{w \in \mathbb{R}^d}{\arg\min} \sum_{(\bar{s},\bar{a},\bar{s}') \in D_h} \left( \phi(\bar{s},\bar{a})^\top w - r_h(\bar{s},\bar{a}) - \hat{V}_{h+1}(\bar{s}') \right)^2$$

$$= \frac{1}{n} \bar{\Lambda}^{-1} \sum_{(\bar{s},\bar{a},\bar{s}') \in D_h} \phi(\bar{s},\bar{a}) \left( r_h(\bar{s},\bar{a}) + \hat{V}_{h+1}(\bar{s}') \right)$$

$$= \frac{1}{n} \bar{\Lambda}^{-1} \sum_{(\bar{s},\bar{a},\bar{s}') \in D_h} \phi(\bar{s},\bar{a}) \left( \phi(\bar{s},\bar{a})^\top \theta_h + \hat{V}_{h+1}(\bar{s}') \right)$$

$$= \theta_h + \frac{1}{n} \bar{\Lambda}^{-1} \sum_{(\bar{s},\bar{a},\bar{s}') \in D_h} \phi(\bar{s},\bar{a}) \hat{V}_{h+1}(\bar{s}')$$

$$\phi(s,a)^\top \hat{w}_h = \phi(s,a)^\top \theta_h + \frac{1}{n} \phi(s,a)^\top \bar{\Lambda}^{-1} \sum_{(\bar{s},\bar{a},\bar{s}') \in D_h} \phi(\bar{s},\bar{a}) \hat{V}_{h+1}(\bar{s}')$$

$$= r_h(s,a) + \left[ \hat{P}_h \hat{V}_{h+1} \right](s,a)$$

Where $\quad \hat{P}_h(s'|s,a) = \frac{1}{n} \phi(s,a)^\top \bar{\Lambda}^{-1} \sum_{(\bar{s},\bar{a},\bar{s}') \in D_h} \phi(\bar{s},\bar{a}) \, \delta_{\bar{s}'}(s')$

$$= \phi(s,a)^T \bar{\Lambda}^{-1} \sum_{j=1}^{d} \phi(\bar{s}_j, \bar{a}_j) \frac{1}{n} \sum_{i=1}^{n} \delta_{\bar{s}'_{j,i}}(s')$$

Lemma (Error decomposition):

$$[P_h V](s,a) = \underset{s' \sim P_h(\cdot | s,a)}{\mathbb{E}}[V(s')]$$

$$V_1^*(s_1) - V_1^{\hat{\pi}}(s_1) = \mathbb{E}_{\pi^*}\left[\sum_{h=1}^{H}[(P_h - \hat{P}_h)\hat{V}_{h+1}](s_h, a_h)\right]$$

$$- \mathbb{E}_{\hat{\pi}}\left[\sum_{h=1}^{H}[(P_h - \hat{P}_h)\hat{V}_{h+1}](s_h, a_h)\right]$$

$$+ \mathbb{E}_{\hat{\pi}}\left[\sum_{h=1}^{H}\underbrace{\langle \hat{Q}_h(s_h, \cdot), \pi^*(\cdot | s_h) - \hat{\pi}(\cdot | s_h)\rangle_A}_{\leq 0}\right]$$

We only need to bound:

$$[(P_h - \hat{P}_h)\hat{V}_{h+1}](s,a)$$

$$\left[ (P_h - \hat{P}_h) \hat{V}_{h+1} \right] (s,a) = \phi(s,a)^T \sum_{s' \in S} \mu_n(s') \hat{V}_{h+1}(s')$$

$$- \phi(s,a)^T \bar{\Lambda}^{-1} \sum_{j=1}^{d} \phi(\bar{s}_j, \bar{a}_j) \frac{1}{n} \sum_{i=1}^{n} \hat{V}_{n+1}(\bar{s}_j^i)$$

$$= \phi(s,a)^T \bar{\Lambda}^{-1} \left[ \sum_{j=1}^{d} \phi(\bar{s}_j, \bar{a}_j) \phi^T(\bar{s}_j, \hat{a}_j) \sum_{s' \in S} \mu_n(s') \hat{V}_{n+1}(s') \right.$$

$$\left. - \sum_{j=1}^{d} \phi(\bar{s}_j, \hat{a}_j) \frac{1}{n} \sum_{i=1}^{n} \hat{V}_{n+1}(\bar{s}'_{j,i}) \right]$$

$$= \phi(s,a)^{T} \bar{\Lambda}^{-1} \left[ \sum_{j=1}^{d} \phi(\bar{s}_j, \bar{a}_j) \left[ \underset{s \sim P_h(\cdot | \bar{s}_j, \bar{a}_j)}{\mathbb{E}} \left[ \hat{V}_{h+1}(s) \right] - \frac{1}{n} \sum_{i=1}^{n} \hat{V}_{n+1}(\bar{s}'_{j,i}) \right] \right]$$

$$\underbrace{\hspace{6cm}}_{\displaystyle H \sqrt{\frac{\log(Hd/\delta)}{n}}}$$

$$= \phi(s,a)^T \left( \psi \psi^T \right)^{-1} \psi \, \varepsilon$$

$$= \phi(s,a)^T \left( \psi^T \right)^{-1} \psi^{-1} \psi \, \varepsilon$$

$$= \phi(s,a)^T (\Psi T)^{-1} \varepsilon$$

$$\leq \| \phi(s,a)^T (\Psi T)^{-1} \|_1 \cdot \| \varepsilon \|_\infty$$

$$\leq L \cdot H \sqrt{\frac{\log(H d / \delta)}{n}}$$