

So far we have assumed $|S| < \infty$, $|A| < \infty$

And our regret scaled w/ $\text{poly}(|S|, |A|) \sqrt{T}$

Today: Two ways to handle infinite state+action spaces.

Let $\phi: S \times A \rightarrow \mathbb{R}^d$ be a feature map.

Def] Linear MDP. Assume $\exists \{\theta_h\}_{h=1}^H \subset \mathbb{R}^d$, $\{\mu_h\} \subset \mathbb{R}^{|S| \times d}$

$$r_h(s, a) = \theta_h^\top \phi(s, a) \quad P_h(s' | s, a) = \mu_h(s') \phi(s, a).$$

Assume $\{\theta_h\}_h$ are known, and $\{\mu_h\}$ are unknown.

Things to note

• P_h has $|S|d$ parameters, and thus could be arbitrarily large.

• If we think of $\underbrace{P_h}_{\mathbb{R}^{|S|d} \times \mathbb{R}^{|S|d}} \in \mathbb{R}^{|S| \times |S| \times d}$ as a matrix, it has rank d .

$$P_h = \mu \Phi \quad \Phi = [\underbrace{\phi(s_1, a_1) \dots \phi(s_{|S|}, a_{|S|})}_{\mathbb{R}^d \times SA} \dots]$$

$\mathbb{R}^{|S|d}$ $\mathbb{R}^d \times SA$

$\rightarrow P_h$ must be at most rank d .

• $Q_h(s, a)$ is linear in $\phi(s, a)$.

$$\begin{aligned}
 Q_h^*(s, a) &= r_h(s, a) + \mathbb{E}_{S' \sim P(\cdot | s, a)} [V_{h+1}^*(s')] \\
 &= \theta_h^T \phi(s, a) + P(\cdot | s, a)^T V_{h+1}^* \\
 &= \theta_h^T \phi(s, a) + (V_{h+1}^*)^T \mu_h \phi(s, a) \\
 &= (\theta_h^T + (V_{h+1}^*)^T \mu_h) \phi(s, a) \\
 &= W_h^T \phi(s, a)
 \end{aligned}$$

Suppose we have data of form

$$\left\{ (s_h^i, a_h^i, s_{h+1}^i) \right\}_{i=1}^{k-1}$$

Want to estimate μ_h

$$\begin{aligned}
 \hat{\mu}_h^k &= \underset{\mu \in \mathbb{R}^{S \times d}}{\operatorname{argmin}} \sum_{i=1}^{k-1} \left\| \delta(s_{h+1}^i) - \mu \phi(s_h^i, a_h^i) \right\|_2^2 + \lambda \|\mu\|_F^2 \\
 &= \sum_{i=1}^{k-1} \delta(s_{h+1}^i) \phi(s_h^i, a_h^i)^T (\Lambda_h^k)^{-1}
 \end{aligned}$$

where $\Lambda_h^k = \sum_{i=1}^{k-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^T + \lambda I$

$$\mathbb{E}[\delta(s_{h+1}^i) | s_h^i, a_h^i] = \mu_h(s_{h+1}^i) \phi(s_h^i, a_h^i)$$

$$\mathbb{E}[\hat{\mu}_h^k] \approx \mu_h^k$$

$$\text{Cov}(\hat{\mu}_h^k(s^i)) = (\Lambda_h^k)^{-1}$$

Key piece of analysis in finite MPS was bounding

$$\begin{aligned} |V^T(\hat{P}_{h+1}(\cdot | s, a) - P_{h+1}(\cdot | s, a))| &= |V^T(\hat{\mu}_{h+1}^k - \mu_{h+1}) \phi(s, a)| \\ &= |V^T(\hat{\mu}_{h+1}^k - \mu_{h+1})(\Lambda_h^k)^{1/2} (\Lambda_h^k)^{-1/2} \phi(s, a)| \\ &\leq \|(\hat{\mu}_{h+1}^k - \mu_{h+1})^T V\|_{\Lambda_h^k} \cdot \|\phi(s, a)\|_{(\Lambda_h^k)^{-1}} \\ &\leq \beta \cdot \|\phi(s, a)\|_{(\Lambda_h^k)^{-1}} \end{aligned}$$

↑
hides factors of H, d .

But no dependence on $|S|$ or $|A|$

Algorithm

for $k=1, 2, \dots$

$$\hat{\mu}_k = \underset{\mu \in \mathbb{R}^{S \times d}}{\operatorname{argmin}} \sum_{i=1}^{k-1} \left\| \delta(s_{h+1}^i) - \mu \phi(s_h^i, a_h^i) \right\|_2^2 + \lambda \|\mu\|_F^2 \quad \forall h$$

$$V_{H+1} = 0$$

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

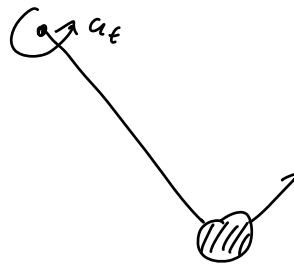
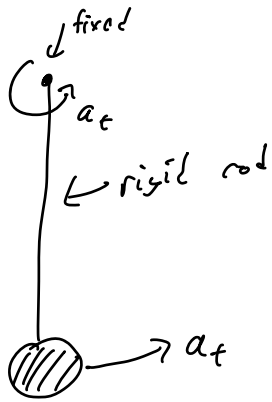
$$\hat{Q}_h^k(s, a) = \min \left\{ H, \beta \|\phi(s, a)\|_{(\hat{\mu}_h^k)^{-1}} + \theta_h^T \phi(s, a) + \hat{V}_{h+1}^k \cdot \hat{\mu}_{h+1}^k \phi(s, a) \right\}$$

$$= \min \left\{ H, \beta \|\phi(s, a)\|_{(\hat{\mu}_h^k)^{-1}} + r_h(s, a) + \hat{P}_{h+1}(s'|s, a)^T \hat{V}_{h+1}^k \right\}$$

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

Roll-out collect $\left\{ (s_h^k, a_h^k, s_{h+1}^k) \right\}_{k=1}^{H+1}$

Final Regret $\leq \sqrt{d^3 K} \operatorname{poly}(H)$.



Goal state

Linear dynamical systems.

State $x_t \in \mathbb{R}^d$

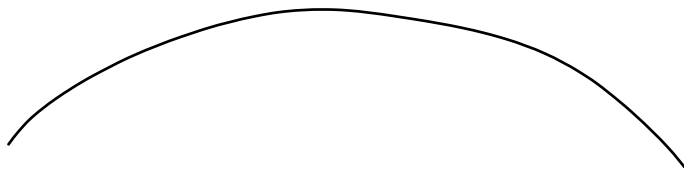
Action (input) $u_t \in \mathbb{R}^p$, $w_t \sim \mathcal{N}(0, I)$ noise

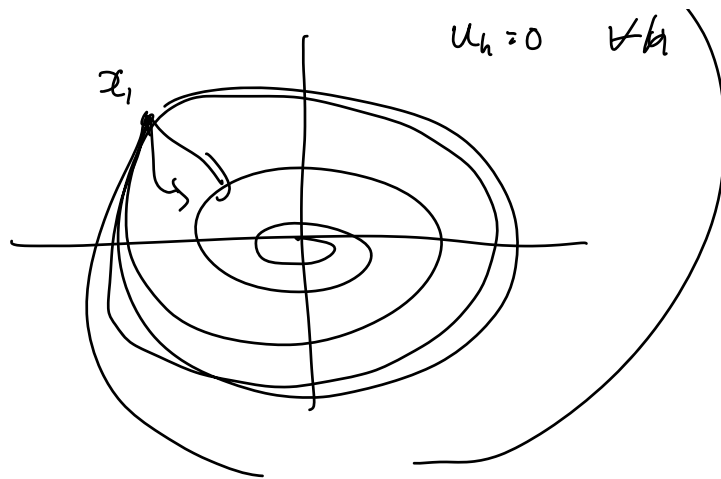
$$x_{t+1} = Ax_t + Bu_t + w_t$$

$$P_{h+1}(\cdot | x_h, u_h) \sim \mathcal{N}(Ax_h + Bu_h, I)$$

For known matrices $Q \in \mathbb{R}^{d \times d}$, $R \in \mathbb{R}^{p \times p}$

loss function $\mathbb{E}\left[x_H^T Q x_H + \sum_{h=1}^{H-1} x_h^T Q x_h + u_h^T R u_h\right]$





$$Q_n^{\pi}(x, u) = \mathbb{E} \left[x_H^T Q x_H + \sum_{t=h}^{H-1} (x_t^T Q x_t + u_t^T R u_t) \mid x_h = x, u_h = u \right]$$

$$Q_H^{\pi}(x, u) = \mathbb{E} [x_H^T Q x_H]$$

$$x_H = A x_{H-1} + B u_{H-1} + w_{H-1}$$

$$Q_{H-1}^{\star}(x, u) = \mathbb{E} \left[(A x + B u + w_{H-1})^T Q (A x + B u + w_{H-1}) + x^T Q x + u^T R u \right]$$

$$\arg \min_u Q_{H-1}^{\star}(x, u) = -K_{H-1} x_{H-1}$$

Value Iteration says in general $\exists K_t \forall t$

$$\text{opt. } u_t = -K_t x_t$$