

Exploration in Linear MDPs

Sham Kakade and Wen Sun

CS 6789: Foundations of Reinforcement Learning

Recap: linear MDP definition

Finite horizon time-dependent episodic MDP $\mathcal{M} = \{S, A, H, \{r\}_h, \{P\}_h, s_0\}$

Recap: linear MDP definition

Finite horizon time-dependent episodic MDP $\mathcal{M} = \{S, A, H, \{r\}_h, \{P\}_h, s_0\}$

Low-Rank
Decomposition:

$$\begin{array}{c} |S| \\ \boxed{P_h(s'|s, a)} \\ |S||A| \end{array} = \begin{array}{c} \boxed{\mu_h} \\ \boxed{\phi} \end{array}$$

Recap: linear MDP definition

Finite horizon time-dependent episodic MDP $\mathcal{M} = \{S, A, H, \{r\}_h, \{P\}_h, s_0\}$

$$P_h(s' | s, a) = \mu_h^\star(s') \cdot \phi(s, a), \quad \mu_h^\star \in S \mapsto \mathbb{R}^d, \phi \in S \times A \mapsto \mathbb{R}^d$$

$$r(s, a) = \theta_h^\star \cdot \phi(s, a), \quad \theta_h^\star \in \mathbb{R}^d$$

Low-Rank
Decomposition:

$$\begin{array}{c} |S| \\ \boxed{P_h(s' | s, a)} \\ |S||A| \end{array} = \begin{array}{c} \boxed{\mu_h} \\ \boxed{\phi} \end{array}$$

Recap:

1. Linear MDP's Q_h^\star & Q_h^π are all linear functions wrt ϕ

Recap:

1. Linear MDP's Q_h^\star & Q_h^π are all linear functions wrt ϕ
2. Since $P_h(\cdot | s, a) = \mu_h^\star \phi(s, a)$, we can estimate μ_h^\star by via regression

Recap:

1. Linear MDP's Q_h^\star & Q_h^π are all linear functions wrt ϕ

2. Since $P_h(\cdot | s, a) = \mu_h^\star \phi(s, a)$, we can estimate μ_h^\star by via regression

Since given s, a , $s' \sim P_h^\star(\cdot | s, a)$, we know $\mathbb{E}_{s' \sim P_h^\star(\cdot | s, a)} \delta(s') = P_h^\star(\cdot | s, a)$,

$$\delta \in \mathbb{R}^{|S|}$$

Recap:

1. Linear MDP's Q_h^\star & Q_h^π are all linear functions wrt ϕ

2. Since $P_h(\cdot | s, a) = \mu_h^\star \phi(s, a)$, we can estimate μ_h^\star by via regression

Since given s, a , $s' \sim P_h^\star(\cdot | s, a)$, we know $\mathbb{E}_{s' \sim P_h^\star(\cdot | s, a)} \delta(s') = P_h^\star(\cdot | s, a)$,

Given dataset $\{s_h^i, a_h^i, s_{h+1}^i\}_{i=1}^{n-1}$

Ridge Linear Regression:

$$\min_{\mu} \sum_{i=1}^{n-1} \|\mu \phi(s_h^i, a_h^i) - \delta(s_{h+1}^i)\|_2^2 + \lambda \|\mu\|_F^2$$

$$\mathbb{E}_{s, a} [\delta] = P_h^\star(\cdot | s, a)$$

Recap:

1. Linear MDP's Q_h^\star & Q_h^π are all linear functions wrt ϕ

2. Since $P_h(\cdot | s, a) = \mu_h^\star \phi(s, a)$, we can estimate μ_h^\star by via regression

Since given s, a , $s' \sim P_h^\star(\cdot | s, a)$, we know $\mathbb{E}_{s' \sim P_h^\star(\cdot | s, a)} \delta(s') = P_h^\star(\cdot | s, a)$,

Given dataset $\{s_h^i, a_h^i, s_{h+1}^i\}_{i=1}^{n-1}$

Ridge Linear Regression:

$$\min_{\mu} \sum_{i=1}^{n-1} \|\mu \phi(s_h^i, a_h^i) - \delta(s_{h+1}^i)\|_2^2 + \lambda \|\mu\|_F^2$$

$$\hat{\mu}_h^n = \sum_{i=1}^{n-1} \delta(s_{h+1}^i) \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}; \quad \widehat{P}_h^n(\cdot | s, a) = \hat{\mu}_h^n \phi(s, a)$$

Today:

Regret bound for the UCBVI algorithm for Linear MDP and its proof sketch

Outline:

1. Model fitting and its guarantee $(\hat{P}(\cdot | s, a) - P(\cdot | s, a))^{\top} V$, for some fixed V

2. Covering argument to bound $(\hat{P}(\cdot | s, a) - P(\cdot | s, a))^{\top} V$, for ALL $V \in \mathcal{F}$

3. UCBVI revisit and its guarantee

1. Model Learning in Linear MDPs

$$\widehat{P}_h^n(\cdot | s, a) = \widehat{\mu}_h^n \phi(s, a)$$

1. Model Learning in Linear MDPs

$$\widehat{P}_h^n(\cdot | s, a) = \widehat{\mu}_h^n \phi(s, a)$$

Cannot directly bound $\|\widehat{P}_h^n(\cdot | s, a) - P(\cdot | s, a)\|_1$ since $P(\cdot | s, a)$ is in $\mathbb{R}^{|S|}$

1. Model Learning in Linear MDPs

$$\widehat{P}_h^n(\cdot | s, a) = \widehat{\mu}_h^n \phi(s, a)$$

Cannot directly bound $\|\widehat{P}_h^n(\cdot | s, a) - P(\cdot | s, a)\|_1$ since $P(\cdot | s, a)$ is in $\mathbb{R}^{|S|}$

As in tabular-UCBVI and Generative Model, we care **average model error**:

1. Model Learning in Linear MDPs

$$\widehat{P}_h^n(\cdot | s, a) = \widehat{\mu}_h^n \phi(s, a)$$

Cannot directly bound $\|\widehat{P}_h^n(\cdot | s, a) - P(\cdot | s, a)\|_1$ since $P(\cdot | s, a)$ is in $\mathbb{R}^{|S|}$

As in tabular-UCBVI and Generative Model, we care **average model error**:

Consider a fixed function $V : S \mapsto [0, H]$, we can bound:

$$\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V \right|$$

$V(s)$

1. Model Learning in Linear MDPs

Ridge Linear Regression:
$$\min_{\mu} \sum_{i=1}^{n-1} \|\mu \phi(s_h^i, a_h^i) - \delta(s_{h+1}^i)\|_2^2 + \|\mu\|_F^2$$

$$\hat{\mu}_h^n = \sum_{i=1}^{n-1} \delta(s_{h+1}^i) \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}, \quad \Lambda_h^n = \sum_{i=1}^{n-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top + \lambda I$$

$$\hat{P}_h^n(\cdot | s, a) = \hat{\mu}_h^n \phi(s, a)$$

1. Model Learning in Linear MDPs

Ridge Linear Regression:
$$\min_{\mu} \sum_{i=1}^{n-1} \|\mu \phi(s_h^i, a_h^i) - \delta(s_{h+1}^i)\|_2^2 + \|\mu\|_F^2$$

$$\hat{\mu}_h^n = \sum_{i=1}^{n-1} \delta(s_{h+1}^i) \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}, \quad \Lambda_h^n = \sum_{i=1}^{n-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top + \lambda I$$

$$\hat{P}_h^n(\cdot | s, a) = \hat{\mu}_h^n \phi(s, a)$$

Lemma [Model Average Error under a fixed V]:

Consider a fixed $V : S \rightarrow [0, H]$. With probability at least $1 - \delta$, for any s, a, h, n , we have:

$$\left| \left(\hat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V \right| \leq \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \times \left(2H \sqrt{\ln \frac{H \det(\Lambda_h^n)^{1/2} \det(\lambda I)^{-1/2}}{\delta}} + H \sqrt{\lambda d} \right)$$

1. Model Learning in Linear MDPs: proof

$$\hat{\mu}_h^n - \mu_h^\star = -\lambda \mu_h^\star (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

1. Model Learning in Linear MDPs: proof

$$\hat{\mu}_h^n - \mu_h^\star = -\lambda \mu_h^\star (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

$$\hat{\mu}_h^n = \sum_{i=1}^{n-1} \delta(s_{h+1}^i) \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}, \quad \Lambda_h^n = \sum_{i=1}^{n-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top + \lambda I$$

1. Model Learning in Linear MDPs: proof

$$\mathbb{E}_{s^a} [\delta(s')] = P(\cdot | s^a)$$

$$\hat{\mu}_h^n - \mu_h^* = -\lambda \mu_h^* (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

$$\hat{\mu}_h^n = \sum_{i=1}^{n-1} \delta(s_{h+1}^i) \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}, \quad \Lambda_h^n = \sum_{i=1}^{n-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top + \lambda I$$

$$\hat{\mu}_h^n = \sum_{i=1}^{n-1} (P_h(\cdot | s_h^i, a_h^i) + \epsilon_h^i) \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

$$\Sigma = \delta(s') - P(\cdot | s^a)$$

$$\mathbb{E}[\Sigma] = 0$$

1. Model Learning in Linear MDPs: proof

$$\hat{\mu}_h^n - \mu_h^\star = -\lambda \mu_h^\star (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

$$\hat{\mu}_h^n = \sum_{i=1}^{n-1} \delta(s_{h+1}^i) \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}, \quad \Lambda_h^n = \sum_{i=1}^{n-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top + \lambda I$$

$$\hat{\mu}_h^n = \sum_{i=1}^{n-1} (P_h(\cdot | s_h^i, a_h^i) + \epsilon_h^i) \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1} = \sum_{i=1}^{n-1} (\underbrace{\mu_h^\star \phi(s_h^i, a_h^i)}_{= P(\cdot | s_h^i, a_h^i)} + \epsilon_h^i) \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

1. Model Learning in Linear MDPs: proof

$$\hat{\mu}_h^n - \mu_h^* = -\lambda \mu_h^* (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

$$\hat{\mu}_h^n = \sum_{i=1}^{n-1} \delta(s_{h+1}^i) \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}, \quad \Lambda_h^n = \sum_{i=1}^{n-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top + \lambda I$$

$$\hat{\mu}_h^n = \sum_{i=1}^{n-1} (P_h(\cdot | s_h^i, a_h^i) + \epsilon_h^i) \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1} = \sum_{i=1}^{n-1} (\mu_h^* \phi(s_h^i, a_h^i) + \epsilon_h^i \phi(s_h^i, a_h^i))^\top (\Lambda_h^n)^{-1}$$

$$= \mu_h^* \left(\sum_{i=1}^{n-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top \right) (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

+λI - λI

1. Model Learning in Linear MDPs: proof

$$\hat{\mu}_h^n - \mu_h^\star = -\lambda \mu_h^\star (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

$$\hat{\mu}_h^n = \sum_{i=1}^{n-1} \delta(s_{h+1}^i) \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}, \quad \Lambda_h^n = \sum_{i=1}^{n-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top + \lambda I$$

$$\begin{aligned} \hat{\mu}_h^n &= \sum_{i=1}^{n-1} (P_h(\cdot | s_h^i, a_h^i) + \epsilon_h^i) \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1} = \sum_{i=1}^{n-1} (\mu_h^\star \phi(s_h^i, a_h^i) + \epsilon_h^i) \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1} \\ &= \mu_h^\star \left(\sum_{i=1}^{n-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top \right) (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1} \\ &= \mu_h^\star - \lambda \mu_h^\star (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1} \end{aligned}$$

1. Model Learning in Linear MDPs

$$\hat{\rho}_{(s,a)} = \hat{\mu} \phi(s,a)$$

$$\hat{\mu}_h^n - \mu_h^* = -\lambda \mu_h^* (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

$$((\hat{\mu}_h^n - \mu_h^*) \phi(s, a))^\top V$$

1. Model Learning in Linear MDPs

$$\hat{\mu}_h^n - \mu_h^* = \left(\lambda \mu_h^* (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1} \right)^{-1}$$

$$((\hat{\mu}_h^n - \mu_h^*) \phi(s, a))^\top V$$

$$= -\lambda \phi(s, a)^\top (\Lambda_h^n)^{-1} (\mu_h^*)^\top V + \sum_{i=1}^{n-1} \phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V$$

1. Model Learning in Linear MDPs

$$\hat{\mu}_h^n - \mu_h^\star = -\lambda \mu_h^\star (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

$$((\hat{\mu}_h^n - \mu_h^\star) \phi(s, a))^\top V$$

$$= -\lambda \phi(s, a)^\top (\Lambda_h^n)^{-1} (\mu_h^\star)^\top V + \sum_{i=1}^{n-1} \phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V$$

$$a^\top b \leq \|a\| \|b\|$$

$$\left| \lambda \phi(s, a)^\top (\Lambda_h^n)^{-1} (\mu_h^\star)^\top V \right| \leq \lambda \|\phi(s, a)^\top (\Lambda_h^n)^{-1/2}\|_2 \|(\Lambda_h^n)^{-1/2} (\mu_h^\star)^\top V\|_2$$

$$\lambda \left(\phi(s, a)^\top (\Lambda_h^n)^{-1/2} \right) \left((\Lambda_h^n)^{-1/2} (\mu_h^\star)^\top V \right)$$

1. Model Learning in Linear MDPs

$$\hat{\mu}_h^n - \mu_h^\star = -\lambda \mu_h^\star (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

$$((\hat{\mu}_h^n - \mu_h^\star) \phi(s, a))^\top V$$

$$= -\lambda \phi(s, a)^\top (\Lambda_h^n)^{-1} (\mu_h^\star)^\top V + \sum_{i=1}^{n-1} \phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V$$

$$\left| \lambda \phi(s, a)^\top (\Lambda_h^n)^{-1} (\mu_h^\star)^\top V \right| \leq \lambda \|\phi(s, a)^\top (\Lambda_h^n)^{-1/2}\|_2 \|(\Lambda_h^n)^{-1/2} (\mu_h^\star)^\top V\|_2$$

$$= \lambda \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \|(\Lambda_h^n)^{-1/2} (\mu_h^\star)^\top V\|_2$$

1. Model Learning in Linear MDPs

$$\hat{\mu}_h^n - \mu_h^\star = -\lambda \mu_h^\star (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

$$((\hat{\mu}_h^n - \mu_h^\star) \phi(s, a))^\top V$$

$$= -\lambda \phi(s, a)^\top (\Lambda_h^n)^{-1} (\mu_h^\star)^\top V + \sum_{i=1}^{n-1} \phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V$$

$\text{Var}(\Lambda_h^n)$

$\Rightarrow \Delta$

$$\left| \lambda \phi(s, a)^\top (\Lambda_h^n)^{-1} (\mu_h^\star)^\top V \right| \leq \lambda \|\phi(s, a)^\top (\Lambda_h^n)^{-1/2}\|_2 \|(\Lambda_h^n)^{-1/2} (\mu_h^\star)^\top V\|_2$$

$$= \lambda \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \|(\Lambda_h^n)^{-1/2} (\mu_h^\star)^\top V\|_2$$

$$\leq \lambda \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \|(\Lambda_h^n)^{-1/2}\|_2 \|(\mu_h^\star)^\top V\|_2$$

Δ

1. Model Learning in Linear MDPs

$$\hat{\mu}_h^n - \mu_h^\star = -\lambda \mu_h^\star (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

$$((\hat{\mu}_h^n - \mu_h^\star) \phi(s, a))^\top V$$

$$= -\lambda \phi(s, a)^\top (\Lambda_h^n)^{-1} (\mu_h^\star)^\top V + \sum_{i=1}^{n-1} \phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V$$

$$\left| \lambda \phi(s, a)^\top (\Lambda_h^n)^{-1} (\mu_h^\star)^\top V \right| \leq \lambda \|\phi(s, a)^\top (\Lambda_h^n)^{-1/2}\|_2 \|(\Lambda_h^n)^{-1/2} (\mu_h^\star)^\top V\|_2$$

$$= \lambda \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \|(\Lambda_h^n)^{-1/2} (\mu_h^\star)^\top V\|_2$$

$$\leq \lambda \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \|(\Lambda_h^n)^{-1/2}\|_2 \|(\mu_h^\star)^\top V\|_2$$

Normalization
assumption on μ_h^\star



1. Model Learning in Linear MDPs

$$\hat{\mu}_h^n - \mu_h^\star = -\lambda \mu_h^\star (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

$$((\hat{\mu}_h^n - \mu_h^\star) \phi(s, a))^\top V$$

$$= -\lambda \phi(s, a)^\top (\Lambda_h^n)^{-1} (\mu_h^\star)^\top V + \sum_{i=1}^{n-1} \phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V$$

$$\left| \lambda \phi(s, a)^\top (\Lambda_h^n)^{-1} (\mu_h^\star)^\top V \right| \leq \lambda \|\phi(s, a)^\top (\Lambda_h^n)^{-1/2}\|_2 \|(\Lambda_h^n)^{-1/2} (\mu_h^\star)^\top V\|_2$$

$$= \lambda \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \|(\Lambda_h^n)^{-1/2} (\mu_h^\star)^\top V\|_2$$

$$\leq \lambda \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \|(\Lambda_h^n)^{-1/2}\|_2 \|(\mu_h^\star)^\top V\|_2 \leq \lambda \|\phi(s, a)^\top\|_{(\Lambda_h^n)^{-1}} \frac{H\sqrt{d}}{\sqrt{\lambda}}$$

Normalization
assumption on μ_h^\star

1. Model Learning in Linear MDPs

$$\hat{\mu}_h^n - \mu_h^\star = -\lambda \mu_h^\star (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

$$\left| ((\hat{\mu}_h^n - \mu_h^\star) \phi(s, a))^\top V \right| \leq H \sqrt{\lambda d} \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} + \left| \sum_{i=1}^{n-1} \phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V \right|$$

1. Model Learning in Linear MDPs

$$\hat{\mu}_h^n - \mu_h^* = -\lambda \mu_h^* (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

$$\left| ((\hat{\mu}_h^n - \mu_h^*) \phi(s, a))^\top V \right| \leq H \sqrt{\lambda d} \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} + \left| \sum_{i=1}^{n-1} \phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V \right|$$

$$\left| \sum_{i=1}^{n-1} \phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V \right| \leq \underbrace{\|\phi(s, a)^\top (\Lambda_h^n)^{-1/2}\|_2}_{\substack{(\Lambda_h^n)^{-1/2} \\ (\Lambda_h^n)^{-1/2}}} \left\| \sum_{i=1}^{n-1} (\Lambda_h^n)^{-1/2} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V \right\|_2$$

$$(\Lambda_h^n)^{-1/2} \quad (\Lambda_h^n)^{-1/2}$$

$$\|\phi(s, a)\|_{\Lambda_h^{n-1}}$$

1. Model Learning in Linear MDPs

$$\hat{\mu}_h^n - \mu_h^\star = -\lambda \mu_h^\star (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

$$\left| ((\hat{\mu}_h^n - \mu_h^\star) \phi(s, a))^\top V \right| \leq H \sqrt{\lambda d} \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} + \left| \sum_{i=1}^{n-1} \phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V \right|$$

$$\left| \sum_{i=1}^{n-1} \phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V \right| \leq \|\phi(s, a)^\top (\Lambda_h^n)^{-1/2}\|_2 \left\| \sum_{i=1}^{n-1} (\Lambda_h^n)^{-1/2} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V \right\|_2$$

$$= \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \left\| \sum_{i=1}^{n-1} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V \right\|_{(\Lambda_h^n)^{-1}}$$

$$E(\xi) = 0$$

$$E[\xi^\top V] = 0$$

1. Model Learning in Linear MDPs

$$\hat{\mu}_h^n - \mu_h^\star = -\lambda \mu_h^\star (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

$$\left| ((\hat{\mu}_h^n - \mu_h^\star) \phi(s, a))^\top V \right| \leq H \sqrt{\lambda d} \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} + \left| \sum_{i=1}^{n-1} \phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V \right|$$

$$\left| \sum_{i=1}^{n-1} \phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V \right| \leq \|\phi(s, a)^\top (\Lambda_h^n)^{-1/2}\|_2 \left\| \sum_{i=1}^{n-1} (\Lambda_h^n)^{-1/2} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V \right\|_2$$

$$= \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \left\| \sum_{i=1}^{n-1} \phi(s_h^i, a_h^i) ((\epsilon_h^i)^\top V) \right\|_{(\Lambda_h^n)^{-1}}$$

$$\mathbb{E} [(\epsilon_h^i)^\top V | \mathcal{H}_{i,h}] = 0, \quad |(\epsilon_h^i)^\top V| \leq \|\epsilon_h^i\|_1 \|V\|_\infty \leq 2H$$

1. Model Learning in Linear MDPs

$$\hat{\mu}_h^n - \mu_h^\star = -\lambda \mu_h^\star (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

$$\left| ((\hat{\mu}_h^n - \mu_h^\star) \phi(s, a))^\top V \right| \leq H \sqrt{\lambda d} \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} + \left| \sum_{i=1}^{n-1} \phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V \right|$$

$$\left| \sum_{i=1}^{n-1} \phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V \right| \leq \|\phi(s, a)^\top (\Lambda_h^n)^{-1/2}\|_2 \left\| \sum_{i=1}^{n-1} (\Lambda_h^n)^{-1/2} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V \right\|_2$$

$$= \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \left\| \sum_{i=1}^{n-1} \phi(s_h^i, a_h^i) ((\epsilon_h^i)^\top V) \right\|_{(\Lambda_h^n)^{-1}} \leq \|\phi(s, a)^\top\|_{(\Lambda_h^n)^{-1}} \times 2H \sqrt{\ln \frac{\det(\Lambda_h^n)^{1/2} \det(\lambda I)^{-1/2}}{\delta}}$$

$$\mathbb{E} [(\epsilon_h^i)^\top V | \mathcal{H}_{i,h}] = 0, \quad |(\epsilon_h^i)^\top V| \leq \|\epsilon_h^i\|_1 \|V\|_\infty \leq 2H$$

1. Model Learning in Linear MDPs

$$\hat{\mu}_h^n - \mu_h^\star = -\lambda \mu_h^\star (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

$$\left| ((\hat{\mu}_h^n - \mu_h^\star) \phi(s, a))^\top V \right| \leq H \sqrt{\lambda d} \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} + \left| \sum_{i=1}^{n-1} \phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V \right|$$

$$\left| \sum_{i=1}^{n-1} \phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V \right| \leq \|\phi(s, a)^\top (\Lambda_h^n)^{-1/2}\|_2 \left\| \sum_{i=1}^{n-1} (\Lambda_h^n)^{-1/2} \phi(s_h^i, a_h^i) (\epsilon_h^i)^\top V \right\|_2$$

$$= \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \left\| \sum_{i=1}^{n-1} \phi(s_h^i, a_h^i) ((\epsilon_h^i)^\top V) \right\|_{(\Lambda_h^n)^{-1}} \leq \|\phi(s, a)^\top\|_{(\Lambda_h^n)^{-1}} \times 2H \sqrt{\ln \frac{\det(\Lambda_h^n)^{1/2} \det(\lambda I)^{-1/2}}{\delta}}$$

With prob $1 - \delta$, $\forall n$

$$\mathbb{E} [(\epsilon_h^i)^\top V | \mathcal{H}_{i,h}] = 0, \quad |(\epsilon_h^i)^\top V| \leq \|\epsilon_h^i\|_1 \|V\|_\infty \leq 2H$$

1. Model Learning in Linear MDPs

Lemma [Model Average Error under a fixed V]:

Consider a fixed $V : S \rightarrow [0, H]$. With probability at least $1 - \delta$, for all s, a, n, h , we have:

$$\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V \right| \leq \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \times \left(2H \sqrt{\ln \frac{H \det(\Lambda_h^n)^{1/2} \det(\lambda I)^{-1/2}}{\delta}} + H \sqrt{\lambda d} \right)$$

1. Model Learning in Linear MDPs

Lemma [Model Average Error under a fixed V]:

Consider a fixed $V : S \rightarrow [0, H]$. With probability at least $1 - \delta$, for all s, a, n, h , we have:

$$\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V \right| \leq \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \times \left(2H \sqrt{\ln \frac{H \det(\Lambda_h^n)^{1/2} \det(\lambda I)^{-1/2}}{\delta}} + H \sqrt{\lambda d} \right)$$

$$\det(\Lambda_h^n)^{1/2} \det(\lambda I)^{-1/2} \leq (n + \lambda)^{d/2} \lambda^{-d/2} = (N/\lambda + 1)^{d/2}$$

$$\sigma_{\max}(\Lambda_h^n)$$

$$\Lambda_h^n = \sum_{i=1}^n \phi_i \phi_i^T + \lambda I$$

$$\det(\Lambda_h^n) \leq (N + \lambda)^d$$

1. Model Learning in Linear MDPs

Lemma [Model Average Error under a fixed V]:

Consider a fixed $V : S \rightarrow [0, H]$. With probability at least $1 - \delta$, for all s, a, n, h , we have:

$$\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V \right| \leq \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \times \left(2H \sqrt{\ln \frac{H \det(\Lambda_h^n)^{1/2} \det(\lambda I)^{-1/2}}{\delta}} + H \sqrt{\lambda d} \right)$$

$$\det(\Lambda_h^n)^{1/2} \det(\lambda I)^{-1/2} \leq (n + \lambda)^{d/2} \lambda^{-d/2} = (N/\lambda + 1)^{d/2}$$

$$\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V \right| \leq \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \left(2H \sqrt{d \ln \left(\frac{NH}{\lambda} + 1 \right) + \ln \left(\frac{1}{\delta} \right)} + H \sqrt{\lambda d} \right)$$

1. Model Learning in Linear MDPs

Lemma [Model Average Error under a fixed V]:

Consider a fixed $V : S \rightarrow [0, H]$. With probability at least $1 - \delta$, for all s, a, n, h , we have:

$$\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V \right| \leq \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \times \left(2H \sqrt{\ln \frac{H \det(\Lambda_h^n)^{1/2} \det(\lambda I)^{-1/2}}{\delta}} + H \sqrt{\lambda d} \right)$$

$$\det(\Lambda_h^n)^{1/2} \det(\lambda I)^{-1/2} \leq (n + \lambda)^{d/2} \lambda^{-d/2} = (N/\lambda + 1)^{d/2}$$

ASA

$$\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V \right| \leq \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \left(2H \sqrt{d \ln \left(\frac{NH}{\lambda} + 1 \right) + \ln \left(\frac{1}{\delta} \right)} + H \sqrt{\lambda d} \right)$$

$$= \widetilde{O} \left(H \sqrt{d} + H \sqrt{\ln(1/\delta)} \right) \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}}$$

1. Model learning: summary

Lemma [Model Average Error under a fixed V]:

$$\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V \right| = \widetilde{O} \left(H\sqrt{d} + H\sqrt{\ln(1/\delta)} \right) \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}}$$

1. Model learning: summary

Lemma [Model Average Error under a fixed V]:

$$\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V \right| = \widetilde{O} \left(H\sqrt{d} + H\sqrt{\ln(1/\delta)} \right) \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}}$$

Q: Can we get a uniform convergence argument for a function class \mathcal{F} ?

$V \in \mathcal{F}$

Outline:



1. Model fitting and its guarantee $(\hat{P}(\cdot | s, a) - P(\cdot | s, a))^{\top} V$, for some fixed V

2. Covering argument to bound $(\hat{P}(\cdot | s, a) - P(\cdot | s, a))^{\top} V$, for ALL $V \in \mathcal{F}$

3. UCBVI revisit and its guarantee

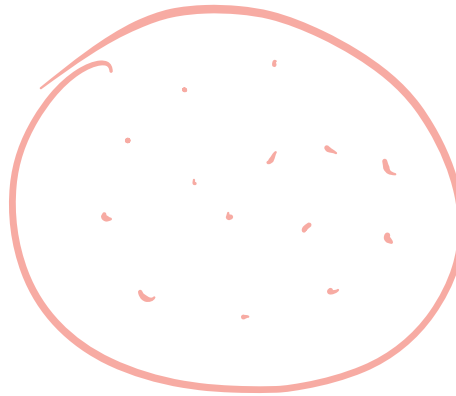
Detour: Covering Number

Consider the ball $\Theta = \{\theta : \theta \in \mathbb{R}^d, \|\theta\|_2 \leq R\}$.

Denote ϵ -**Net** as a subset $\mathcal{N}_\epsilon \subseteq \Theta$, such that $\forall \theta \in \Theta$:

$$\exists \theta' \in \mathcal{N}_\epsilon, \text{ s.t. } \|\theta' - \theta\|_2 \leq \epsilon.$$

Denote ϵ -**cover** as the smallest \mathcal{N}_ϵ



Detour: Covering Number

Consider the ball $\Theta = \{\theta : \theta \in \mathbb{R}^d, \|\theta\|_2 \leq R\}$.

Denote ϵ -**Net** as a subset $\mathcal{N}_\epsilon \subseteq \Theta$, such that $\forall \theta \in \Theta$:

$$\exists \theta' \in \mathcal{N}_\epsilon, \text{ s.t. } \|\theta' - \theta\|_2 \leq \epsilon.$$

Denote ϵ -**cover** as the smallest \mathcal{N}_ϵ

Lemma [Covering of Θ] We have $|\mathcal{N}_\epsilon| \leq (1 + 2R/\epsilon)^d$, and $\ln(|\mathcal{N}_\epsilon|) \leq d \ln(1 + 2R/\epsilon)$

Detour: Covering Number

Consider the ball $\Theta = \{\theta : \theta \in \mathbb{R}^d, \|\theta\|_2 \leq R\}$.

Denote ϵ -**Net** as a subset $\mathcal{N}_\epsilon \subseteq \Theta$, such that $\forall \theta \in \Theta$:

$$\exists \theta' \in \mathcal{N}_\epsilon, \text{ s.t. } \|\theta' - \theta\|_2 \leq \epsilon.$$

Denote ϵ -**cover** as the smallest \mathcal{N}_ϵ

Lemma [Covering of Θ] We have $|\mathcal{N}_\epsilon| \leq (1 + 2R/\epsilon)^d$, and $\ln(|\mathcal{N}_\epsilon|) \leq d \ln(1 + 2R/\epsilon)$

Now consider a function class $\mathcal{F} = \{f_\theta : \theta \in \Theta\}$,
and for any $f_{\theta_1}, f_{\theta_2} \in \mathcal{F}$, $\|f_{\theta_1} - f_{\theta_2}\|_\infty \leq L\|\theta_1 - \theta_2\|_2$



Detour: Covering Number

Consider the ball $\Theta = \{\theta : \theta \in \mathbb{R}^d, \|\theta\|_2 \leq R\}$.

Denote ϵ -**Net** as a subset $\mathcal{N}_\epsilon \subseteq \Theta$, such that $\forall \theta \in \Theta$:

$$\exists \theta' \in \mathcal{N}_\epsilon, \text{ s.t. } \|\theta' - \theta\|_2 \leq \epsilon.$$

Denote ϵ -**cover** as the smallest \mathcal{N}_ϵ

Lemma [Covering of Θ] We have $|\mathcal{N}_\epsilon| \leq (1 + 2R/\epsilon)^d$, and $\ln(|\mathcal{N}_\epsilon|) \leq d \ln(1 + 2R/\epsilon)$

Now consider a function class $\mathcal{F} = \{f_\theta : \theta \in \Theta\}$,

and for any $f_{\theta_1}, f_{\theta_2} \in \mathcal{F}$, $\|f_{\theta_1} - f_{\theta_2}\|_\infty \leq L \|\theta_1 - \theta_2\|_2$

Then (ϵ/L) -Net on Θ gives us an ϵ -Net on \mathcal{F} with $d(f_{\theta_1}, f_{\theta_2}) := \|f_{\theta_1} - f_{\theta_2}\|_\infty$

$$L \cdot \frac{\epsilon}{L}$$

Detour: Covering Number and An Example

Consider a specific parameterization $\theta = (w, \beta, \Lambda)$,

$$\Theta = \{(w, \beta, \Lambda) : \|w\|_2 \leq L, \beta \in [0, B], \sigma_{\min}(\Lambda) \geq \lambda\}$$



Detour: Covering Number and An Example

Consider a specific parameterization $\theta = (w, \beta, \Lambda)$,
 $\Theta = \{(w, \beta, \Lambda) : \|w\|_2 \leq L, \beta \in [0, B], \sigma_{\min}(\Lambda) \geq \lambda\}$

Define the function $f_{w, \beta, \Lambda} : S \rightarrow [0, H]$

$$f_{w, \beta, \Lambda}(s) := \min \left\{ \max_a \left(\underbrace{w^\top \phi(s, a)}_{Q^*} + \beta \underbrace{\sqrt{\phi(s, a)^\top \Lambda^{-1} \phi(s, a)}}_{\text{bonus}}, H \right) \right\}$$

Detour: Covering Number and An Example

Consider a specific parameterization $\theta = (w, \beta, \Lambda)$,
 $\Theta = \{(w, \beta, \Lambda) : \|w\|_2 \leq L, \beta \in [0, B], \sigma_{\min}(\Lambda) \geq \lambda\}$

Define the function $f_{w, \beta, \Lambda} : S \rightarrow [0, H]$

$$f_{w, \beta, \Lambda}(s) := \min \left\{ \max_a \left(w^\top \phi(s, a) + \beta \sqrt{\phi(s, a)^\top \Lambda^{-1} \phi(s, a)} \right), H \right\}$$

Denote $\mathcal{F} = \{f_{w, \beta, \Lambda} : \|w\|_2 \leq L, \beta \in [0, B], \sigma_{\min}(\Lambda) \geq \lambda\}$, **what's the covering number of \mathcal{F} under ℓ_∞**

$$\|f_\theta - f_{\theta'}\|_\infty \leq L \cdot \|\theta - \theta'\|$$

Detour: Covering Number

$$f_{w,\beta,\Lambda} \text{ :}, f_{w,\beta,\Lambda}(s) := \min \left\{ \max_a \left(w^\top \phi(s, a) + \beta \sqrt{\phi(s, a)^\top \Lambda^{-1} \phi(s, a)} \right), H \right\}$$

Detour: Covering Number

$$f_{w,\beta,\Lambda} := f_{w,\beta,\Lambda}(s) := \min \left\{ \max_a \left(w^\top \phi(s, a) + \beta \sqrt{\phi(s, a)^\top \Lambda^{-1} \phi(s, a)} \right), H \right\}$$

Lemma: Denote $\mathcal{F} = \{f_{w,\beta,\Lambda} : \|w\|_2 \leq L, \beta \in [0, B], \sigma_{\min}(\Lambda) \geq \lambda\}$,
under ℓ_∞ we have: $\ln |\mathcal{N}_\epsilon| \leq d \ln(1 + 6L/\epsilon) + 2d^2 \ln(1 + 18B^2\sqrt{d}/(\lambda\epsilon^2)) = \widetilde{O}(d^2)$

Detour: Covering Number

$$f_{w,\beta,\Lambda} \text{ :}, f_{w,\beta,\Lambda}(s) := \min \left\{ \max_a \left(w^\top \phi(s, a) + \beta \sqrt{\phi(s, a)^\top \Lambda^{-1} \phi(s, a)} \right), H \right\}$$

Lemma: Denote $\mathcal{F} = \{f_{w,\beta,\Lambda} : \|w\|_2 \leq L, \beta \in [0, B], \sigma_{\min}(\Lambda) \geq \lambda\}$,
 under ℓ_∞ we have: $\ln |\mathcal{N}_\epsilon| \leq d \ln(1 + 6L/\epsilon) + 2d^2 \ln(1 + 18B^2\sqrt{d}/(\lambda\epsilon^2)) = \tilde{O}(d^2)$

Key step in the proof:

$$\underbrace{|f_\theta(s) - f_{\hat{\theta}}(s)|}_{\Delta \quad \Delta} \leq \underbrace{\|w - \hat{w}\|_2}_{\Sigma\text{-Net}} + \underbrace{|\beta - \hat{\beta}|/\sqrt{\lambda}}_{\cancel{\Sigma} \quad \Sigma \cdot \sqrt{\lambda} \text{-Net}} + B \underbrace{\sqrt{\|\Lambda^{-1} - \hat{\Lambda}^{-1}\|_F}}_{\frac{\sqrt{\Sigma}}{B} \text{-Net}}$$

Detour: Uniform Convergence

$$f_{w,\beta,\Lambda} \text{ : } f_{w,\beta,\Lambda}(s) := \min \left\{ \max_a \left(w^\top \phi(s, a) + \beta \sqrt{\phi(s, a)^\top \Lambda^{-1} \phi(s, a)} \right), H \right\}$$

Lemma: Denote $\mathcal{F} = \{f_{w,\beta,\Lambda} : \|w\|_2 \leq L, \beta \in [0, B], \sigma_{\min}(\Lambda) \geq \lambda\}$,

under ℓ_∞ we have: $\ln |\mathcal{N}_\epsilon| \leq d \ln(1 + 6L/\epsilon) + 2d^2 \ln(1 + 18B^2 \sqrt{d}/(\lambda\epsilon^2)) = \tilde{O}(d^2)$

Lemma [uniform convergence]: With probability at least $1 - \delta$, for all s, a, h, n , and **ALL** $f \in \mathcal{F}$:

$$\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot f \right| = \tilde{O}(Hd) \cdot \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}}$$

$f \in \mathcal{F}$

(fixed ν :
 $H\sqrt{d} \|\phi\|_{\Lambda^{-1}}$)

Detour: Uniform Convergence

$$f_{w,\beta,\Lambda} \text{ :}, f_{w,\beta,\Lambda}(s) := \min \left\{ \max_a \left(w^\top \phi(s, a) + \beta \sqrt{\phi(s, a)^\top \Lambda^{-1} \phi(s, a)} \right), H \right\}$$

Lemma: Denote $\mathcal{F} = \{f_{w,\beta,\Lambda} : \|w\|_2 \leq L, \beta \in [0, B], \sigma_{\min}(\Lambda) \geq \lambda\}$,

under ℓ_∞ we have: $\ln |\mathcal{N}_\epsilon| \leq d \ln(1 + 6L/\epsilon) + 2d^2 \ln(1 + 18B^2\sqrt{d}/(\lambda\epsilon^2)) = \widetilde{O}(d^2)$

Lemma [uniform convergence]: With probability at least $1 - \delta$, for all s, a, h, n , and **ALL** $f \in \mathcal{F}$:

$$\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot f \right| = \widetilde{O}(Hd) \cdot \left\| \phi(s, a) \right\|_{(\Lambda_h^n)^{-1}}$$

Proof sketch: the model error we had for a fixed $V + \epsilon$ -Net argument
(Same high level steps as the ones we used for HW1's last question)

Detour: Uniform Convergence

$$f_{w,\beta,\Lambda} \text{ :}, f_{w,\beta,\Lambda}(s) := \min \left\{ \max_a \left(w^\top \phi(s, a) + \beta \sqrt{\phi(s, a)^\top \Lambda^{-1} \phi(s, a)} \right), H \right\}$$

Lemma: Denote $\mathcal{F} = \{f_{w,\beta,\Lambda} : \|w\|_2 \leq L, \beta \in [0, B], \sigma_{\min}(\Lambda) \geq \lambda\}$,

under ℓ_∞ we have: $\ln |\mathcal{N}_\epsilon| \leq d \ln(1 + 6L/\epsilon) + 2d^2 \ln(1 + 18B^2 \sqrt{d}/(\lambda\epsilon^2)) = \tilde{O}(d^2)$

Lemma [uniform convergence]: With probability at least $1 - \delta$, for all s, a, h, n , and **ALL** $f \in \mathcal{F}$:

$$\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot f \right| = \tilde{O}(Hd) \cdot \left\| \phi(s, a) \right\|_{(\Lambda_h^n)^{-1}}$$

This will be our bonus term

Proof sketch: the model error we had for a fixed $V + \epsilon$ -Net argument
(Same high level steps as the ones we used for HW1's last question)

Summary of Covering Argument

Covering allows us to build a uniform convergence result (i.e., $\forall f \in \mathcal{F}$)
over a infinite hypothesis class
(Intuitively, log of covering number scales w.r.t to the # of parameters)

Outline:



1. Model fitting and its guarantee $(\hat{P}(\cdot | s, a) - P(\cdot | s, a))^{\top} V$, for some fixed V



2. Covering argument to bound $(\hat{P}(\cdot | s, a) - P(\cdot | s, a))^{\top} V$, for ALL $V \in \mathcal{F}$

3. UCBVI revisit and its guarantee

Algorithm: UCBVI in Linear MDPs

At the beginning of iteration n:

1. Learn transition model $\{\hat{P}_h^n\}_{h=0}^{H-1}$ from all previous data via Ridge linear regression

$$\hat{\mu} \cdot \phi$$

$$\# \text{ comment: } \min_{\mu} \sum_{i=1}^{n-1} \|\mu \phi(s_h^i, a_h^i) - \delta(s_{h+1}^i)\|_2^2 + \lambda \|\mu\|_F^2$$

Algorithm: UCBVI in Linear MDPs

At the beginning of iteration n :

1. Learn transition model $\{\widehat{P}_h^n\}_{h=0}^{H-1}$ from all previous data via Ridge linear regression

comment: $\min_{\mu} \sum_{i=1}^{n-1} \|\mu \phi(s_h^i, a_h^i) - \delta(s_{h+1}^i)\|_2^2 + \lambda \|\mu\|_F^2$

2. Design reward bonus $b_h^n(s, a) = \beta \sqrt{\phi(s, a) (\Lambda_h^n)^{-1} \phi(s, a)}$

$\|\phi(s, a)\|_{(\Lambda_h^n)^{-1}}$

Algorithm: UCBVI in Linear MDPs

At the beginning of iteration n :

1. Learn transition model $\{\widehat{P}_h^n\}_{h=0}^{H-1}$ from all previous data via Ridge linear regression

$$\# \text{ comment: } \min_{\mu} \sum_{i=1}^{n-1} \|\mu \phi(s_h^i, a_h^i) - \delta(s_{h+1}^i)\|_2^2 + \lambda \|\mu\|_F^2$$

2. Design reward bonus $b_h^n(s, a) = \beta \sqrt{\phi(s, a) (\Lambda_h^n)^{-1} \phi(s, a)}$

#Comment: $\beta = \widetilde{O}(Hd)$, reward bonus upper bounds $\left| (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot f \right|, \forall f \in \mathcal{F}$

Algorithm: UCBVI in Linear MDPs

At the beginning of iteration n:

1. Learn transition model $\{\widehat{P}_h^n\}_{h=0}^{H-1}$ from all previous data via Ridge linear regression

$$\# \text{ comment: } \min_{\mu} \sum_{i=1}^{n-1} \|\mu \phi(s_h^i, a_h^i) - \delta(s_{h+1}^i)\|_2^2 + \lambda \|\mu\|_F^2$$

2. Design reward bonus $b_h^n(s, a) = \beta \sqrt{\phi(s, a) (\Lambda_h^n)^{-1} \phi(s, a)}$

#Comment: $\beta = \widetilde{O}(Hd)$, reward bonus upper bounds $\left| (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot f \right|, \forall f \in \mathcal{F}$

3. Plan: $\pi^{n+1} = \text{Value-Iter} \left(\{\widehat{P}^n\}_h, \{r_h + b_h^n\} \right)$

Algorithm: UCBVI in Linear MDPs

At the beginning of iteration n:

1. Learn transition model $\{\widehat{P}_h^n\}_{h=0}^{H-1}$ from all previous data via Ridge linear regression

$$\# \text{ comment: } \min_{\mu} \sum_{i=1}^{n-1} \|\mu \phi(s_h^i, a_h^i) - \delta(s_{h+1}^i)\|_2^2 + \lambda \|\mu\|_F^2$$

2. Design reward bonus $b_h^n(s, a) = \beta \sqrt{\phi(s, a)(\Lambda_h^n)^{-1}\phi(s, a)}$

#Comment: $\beta = \widetilde{O}(Hd)$, reward bonus upper bounds $\left| (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot f \right|, \forall f \in \mathcal{F}$

3. Plan: $\pi^{n+1} = \text{Value-Iter} \left(\{\widehat{P}_h^n\}_h, \{r_h + b_h^n\} \right)$

4. Execute π^{n+1} for H steps

Regret bound for UCBVI in linear MDP

$$\mathbb{E} \left[\sum_{n=1}^N (V^{\star} - V^{\pi^n}) \right] \leq \widetilde{O}(H^2 d^{1.5} \sqrt{N})$$

No S or A polynomial dependence!

2. Value Iteration in the Learned Model w/ Reward Bonus

$$\pi^{n+1} = \text{Value-Iter} \left(\{ \widehat{P}^n \}_h, \{ r_h + b_h^n \} \right)$$

2. Value Iteration in the Learned Model w/ Reward Bonus

$$\pi^{n+1} = \text{Value-Iter} \left(\{ \widehat{P}^n \}_h, \{ r_h + b_h^n \} \right)$$

$$\widehat{V}_H^n(s) = 0, \forall s,$$



2. Value Iteration in the Learned Model w/ Reward Bonus

$$\pi^{n+1} = \text{Value-Iter} \left(\{ \widehat{P}^n \}_h, \{ r_h + b_h^n \} \right)$$

$$\widehat{V}_H^n(s) = 0, \forall s,$$

$H \rightarrow 0$

$$\widehat{Q}_h^n(s, a) = r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n$$

$\underline{\quad}$ \triangle

2. Value Iteration in the Learned Model w/ Reward Bonus

$$\pi^{n+1} = \text{Value-Iter} \left(\{ \widehat{P}^n \}_h, \{ r_h + b_h^n \} \right)$$

$$\widehat{V}_H^n(s) = 0, \forall s,$$

$$\widehat{Q}_h^n(s, a) = r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n$$

$$= \theta_h^* \cdot \phi(s, a) + \beta \sqrt{\phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s, a)} + \underbrace{(\widehat{\mu}_h^n \phi(s, a))^\top}_{\text{red underline}} \widehat{V}_{h+1}^n$$

2. Value Iteration in the Learned Model w/ Reward Bonus

$$\pi^{n+1} = \text{Value-Iter} \left(\{ \widehat{P}^n \}_h, \{ r_h + b_h^n \} \right)$$

$$\widehat{V}_H^n(s) = 0, \forall s,$$

$$\widehat{Q}_h^n(s, a) = r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n$$

$$= \theta_h^* \cdot \phi(s, a) + \beta \sqrt{\phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s, a)} + (\widehat{\mu}_h^n \phi(s, a))^\top \widehat{V}_{h+1}^n$$

linear in ϕ

$$= \beta \sqrt{\phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s, a)} + \left(\theta_h^* + (\widehat{\mu}_h^n)^\top \widehat{V}_{h+1}^n(s') \right)^\top \phi(s, a)$$

2. Value Iteration in the Learned Model w/ Reward Bonus

$$\pi^{n+1} = \text{Value-Iter} \left(\{ \widehat{P}^n \}_h, \{ r_h + b_h^n \} \right)$$

$$\widehat{V}_H^n(s) = 0, \forall s,$$

$$\widehat{Q}_h^n(s, a) = r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n$$

$$= \theta_h^* \cdot \phi(s, a) + \beta \sqrt{\phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s, a)} + (\widehat{\mu}_h^n \phi(s, a))^\top \widehat{V}_{h+1}^n$$

$$= \beta \sqrt{\phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s, a)} + \left(\theta_h^* + (\widehat{\mu}_h^n)^\top \widehat{V}_{h+1}^n(s') \right)^\top \phi(s, a)$$

$$= \beta \sqrt{\phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s, a)} + \phi(s, a)^\top \widehat{w}_h^n$$

\widehat{w}_h^n
 \hat{w}_h^n

2. Value Iteration in the Learned Model w/ Reward Bonus

$$\pi^{n+1} = \text{Value-Iter} \left(\{ \widehat{P}^n \}_h, \{ r_h + b_h^n \} \right)$$

$$\widehat{V}_H^n(s) = 0, \forall s,$$

$$\widehat{Q}_h^n(s, a) = r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n$$

$$= \theta_h^* \cdot \phi(s, a) + \beta \sqrt{\phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s, a)} + (\widehat{\mu}_h^n \phi(s, a))^\top \widehat{V}_{h+1}^n$$

$$= \beta \sqrt{\phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s, a)} + \left(\theta_h^* + (\widehat{\mu}_h^n)^\top \widehat{V}_{h+1}^n(s') \right)^\top \phi(s, a)$$

$$= \beta \sqrt{\phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s, a)} + \phi(s, a)^\top \widehat{w}_h^n$$

$$\widehat{V}_h^n(s) = \min \left\{ \max_a \left(\phi(s, a)^\top \widehat{w}_h^n + \beta \sqrt{\phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s, a)} \right), H \right\}, \left(\pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a) \right)$$

$f \in \mathcal{F}$

3. Prove Optimism

Lemma [Optimism]: with high probability, for all n, h, s :

$$\widehat{V}_h^n(s) \geq V_h^*(s)$$

3. Prove Optimism

Lemma [Optimism]: with high probability, for all n, h, s :

$$\widehat{V}_h^n(s) \geq V_h^*(s)$$

Proof Sketch: let's do induction here with Inductive hypothesis: $\widehat{V}_{h+1}^n(s) \geq V_{h+1}^*(s), \forall s$

3. Prove Optimism

Lemma [Optimism]: with high probability, for all n, h, s :

$$\widehat{V}_h^n(s) \geq V_h^*(s)$$

Proof Sketch: let's do induction here with Inductive hypothesis: $\widehat{V}_{h+1}^n(s) \geq V_{h+1}^*(s), \forall s$

$$\widehat{Q}_h^n(s, a) - Q_h^*(s, a) = b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n - P_h(\cdot | s, a) \cdot V_{h+1}^* \sim \widehat{V}_{h+1}^n$$

Δ Δ \times

\uparrow

$\widehat{r} + b + \widehat{P}(s, a) \cdot \widehat{V}_{h+1}$ $\widehat{r} + P^*(s, a) \cdot V_{h+1}^*$

(Note: In the original image, \widehat{V}_{h+1}^n is circled in green and labeled V_{h+1}^* , and V_{h+1}^* is circled in red. A red arrow points from the difference between the two terms to the inductive hypothesis \widehat{V}_{h+1}^n .)

3. Prove Optimism

Lemma [Optimism]: with high probability, for all n, h, s :

$$\widehat{V}_h^n(s) \geq V_h^*(s)$$

Proof Sketch: let's do induction here with Inductive hypothesis: $\widehat{V}_{h+1}^n(s) \geq V_{h+1}^*(s), \forall s$

$$\widehat{Q}_h^n(s, a) - Q_h^*(s, a) = b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n - P_h(\cdot | s, a) \cdot V_{h+1}^*$$

$$\geq b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n - P_h(\cdot | s, a) \cdot \widehat{V}_{h+1}^n$$

Valid Transition

3. Prove Optimism

Lemma [Optimism]: with high probability, for all n, h, s :

$$\widehat{V}_h^n(s) \geq V_h^*(s)$$

Proof Sketch: let's do induction here with Inductive hypothesis: $\widehat{V}_{h+1}^n(s) \geq V_{h+1}^*(s), \forall s$

$$\widehat{Q}_h^n(s, a) - Q_h^*(s, a) = b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n - P_h(\cdot | s, a) \cdot V_{h+1}^*$$

$$\geq b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n - P_h(\cdot | s, a) \cdot \widehat{V}_{h+1}^n$$

NOTE this is different from what we did in tabular MDP!!!

\widehat{P} is not a valid
Transition

3. Prove Optimism

Lemma [Optimism]: with high probability, for all n, h, s :

$$\widehat{V}_h^n(s) \geq V_h^*(s)$$

Proof Sketch: let's do induction here with Inductive hypothesis: $\widehat{V}_{h+1}^n(s) \geq V_{h+1}^*(s), \forall s$

$$\widehat{Q}_h^n(s, a) - Q_h^*(s, a) = b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n - P_h(\cdot | s, a) \cdot V_{h+1}^*$$

$$\geq b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n - P_h(\cdot | s, a) \cdot \widehat{V}_{h+1}^n$$

NOTE this is different from what we did in tabular MDP!!!

$$\geq b_h^n(s, a) - \left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot \widehat{V}_{h+1}^n \right|$$



3. Prove Optimism

Lemma [Optimism]: with high probability, for all n, h, s :

$$\widehat{V}_h^n(s) \geq V_h^*(s)$$

Proof Sketch: let's do induction here with Inductive hypothesis: $\widehat{V}_{h+1}^n(s) \geq V_{h+1}^*(s), \forall s$

$$\widehat{Q}_h^n(s, a) - Q_h^*(s, a) = b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n - P_h(\cdot | s, a) \cdot V_{h+1}^*$$

$$\geq b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n - P_h(\cdot | s, a) \cdot \widehat{V}_{h+1}^n$$

$$\geq b_h^n(s, a) - \left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot \widehat{V}_{h+1}^n \right|$$

$$\geq 0$$

NOTE this is different from what we did in tabular MDP!!!

$$\leq \text{hd} \|\Phi\|_{\infty} - 1$$

4. Regret Decomposition

Conditioned on history up to the end of episode $n-1$:

$$V_0^*(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$$

value of π^n in (P, r, γ, b)

optimism

$$\pi_n = VI(\widehat{P}, r, \gamma, b)$$

value of π_n in $(\widehat{P}, r, \gamma, b)$

4. Regret Decomposition

Conditioned on history up to the end of episode $n-1$:

$$V_0^*(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$$

(apply Simulation Lemma here)

4. Regret Decomposition

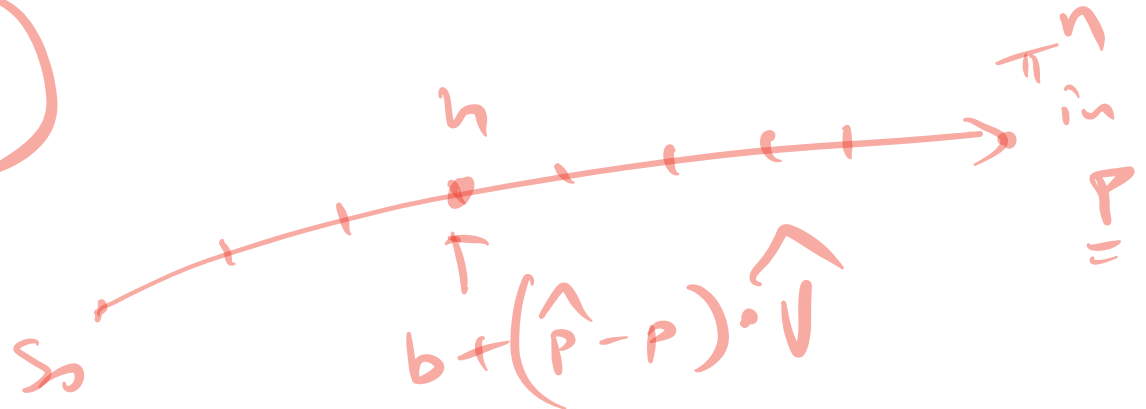
Conditioned on history up to the end of episode $n-1$:

$$V_0^*(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$$

(apply Simulation Lemma here)

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s_h, a_h \sim d_h^{\pi^n}} \left[\underbrace{b_h^n(s_h, a_h)}_{\gamma + b - \tau} + \underbrace{\left(\widehat{P}_h^n(\cdot | s_h, a_h) - P_h(\cdot | s_h, a_h) \right) \cdot \widehat{V}_{h+1}^n}_{b + (\widehat{P} - P) \cdot \widehat{V}} \right]$$

$$(\gamma + b - \tau)$$



4. Regret Decomposition

Conditioned on history up to the end of episode $n-1$:

$$V_0^*(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$$

(apply Simulation Lemma here)

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s_h, a_h \sim d_h^{\pi^n}} \left[b_h^n(s_h, a_h) + \left(\widehat{P}_h^n(\cdot | s_h, a_h) - P_h(\cdot | s_h, a_h) \right) \cdot \widehat{V}_{h+1}^n \right]$$

$$\stackrel{2}{\leq} \sum_{h=0}^{H-1} \mathbb{E}_{s_h, a_h \sim d_h^{\pi^n}} [b_h^n(s_h, a_h)]$$

$$= b_h^n(s_h, a_h)$$

4. Regret Decomposition

Conditioned on history up to the end of episode $n-1$:

$$V_0^*(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$$

(apply Simulation Lemma here)

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s_h, a_h \sim d_h^{\pi^n}} \left[b_h^n(s_h, a_h) + \left(\widehat{P}_h^n(\cdot | s_h, a_h) - P_h(\cdot | s_h, a_h) \right) \cdot \widehat{V}_{h+1}^n \right]$$

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s_h, a_h \sim d_h^{\pi^n}} [b_h^n(s_h, a_h)]$$

being loges \rightarrow

$$= \sum_{h=0}^{H-1} \mathbb{E}_{s_h, a_h \sim d_h^{\pi^n}} \left[\beta \sqrt{\phi(s_h, a_h)^\top (\Lambda_h^n)^{-1} \phi(s_h, a_h)} \right]$$

4. Concluding the Regret Computation

$$\mathbb{E} \left[\sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right] = \mathbb{E} \left[\mathbf{1}[\text{good event holds}] \sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right] + \mathbb{E} \left[\mathbf{1}[\text{good event doesn't hold}] \sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right]$$

4. Concluding the Regret Computation

$$\begin{aligned} \mathbb{E} \left[\sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right] &= \mathbb{E} \left[\mathbf{1}[\text{good event holds}] \sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right] + \mathbb{E} \left[\mathbf{1}[\text{good event doesn't hold}] \sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right] \\ &\lesssim \beta \mathbb{E} \left[\sum_{n=1}^N \sum_{h=0}^{H-1} \sqrt{\phi(s_h^n, a_h^n)^\top (\Lambda_h^n)^{-1} \phi(s_h^n, a_h^n)} \right] + \delta NH \end{aligned}$$

4. Concluding the Regret Computation

$$\begin{aligned} \mathbb{E} \left[\sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right] &= \mathbb{E} \left[\mathbf{1}[\text{good event holds}] \sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right] + \mathbb{E} \left[\mathbf{1}[\text{good event doesn't hold}] \sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right] \\ &\lesssim \beta \mathbb{E} \left[\sum_{n=1}^N \sum_{h=0}^{H-1} \sqrt{\phi(s_h^n, a_h^n)^\top (\Lambda_h^n)^{-1} \phi(s_h^n, a_h^n)} \right] + \delta NH \\ &\lesssim \beta \mathbb{E} \left[\sum_{h=0}^{H-1} \sqrt{N} \sqrt{\sum_{n=1}^N \phi(s_h^n, a_h^n)^\top (\Lambda_h^n)^{-1} \phi(s_h^n, a_h^n)} \right] + \delta NH \end{aligned}$$

4. Concluding the Regret Computation

$$\begin{aligned}
\mathbb{E} \left[\sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right] &= \mathbb{E} \left[\mathbf{1}[\text{good event holds}] \sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right] + \mathbb{E} \left[\mathbf{1}[\text{good event doesn't hold}] \sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right] \\
&\lesssim \beta \mathbb{E} \left[\sum_{n=1}^N \sum_{h=0}^{H-1} \sqrt{\phi(s_h^n, a_h^n)^\top (\Lambda_h^n)^{-1} \phi(s_h^n, a_h^n)} \right] + \delta NH \\
&\lesssim \beta \mathbb{E} \left[\sum_{h=0}^{H-1} \sqrt{N} \sqrt{\sum_{n=1}^N \phi(s_h^n, a_h^n)^\top (\Lambda_h^n)^{-1} \phi(s_h^n, a_h^n)} \right] + \delta NH \\
&\lesssim \widetilde{O}(H^2 d^{1.5} \sqrt{N})
\end{aligned}$$