

Exploration in Linear MDPs

Sham Kakade and Wen Sun

CS 6789: Foundations of Reinforcement Learning

Recap:

Recap:

Stochastic Linear Bandits

$$\mathcal{D} \subset \mathbb{R}^d \quad r(x) = \theta^\star \cdot x, \forall x$$

Recap:

Stochastic Linear Bandits

$$\mathcal{D} \subset \mathbb{R}^d \quad r(x) = \theta^\star \cdot x, \forall x$$

Every episode n , learner picks $x_n \in \mathcal{D}$

Recap:

Stochastic Linear Bandits

$$\mathcal{D} \subset \mathbb{R}^d \quad r(x) = \theta^\star \cdot x, \forall x$$

Every episode n , learner picks $x_n \in \mathcal{D}$

Learner receives a scalar $r_n = \theta^\star \cdot x_n + \epsilon_n$, $\mathbb{E}[\epsilon_n] = 0, |\epsilon_n| < \alpha$

Recap:

Stochastic Linear Bandits

$$\mathcal{D} \subset \mathbb{R}^d \quad r(x) = \theta^\star \cdot x, \forall x$$

Every episode n , learner picks $x_n \in \mathcal{D}$

Learner receives a scalar $r_n = \theta^\star \cdot x_n + \epsilon_n$, $\mathbb{E}[\epsilon_n] = 0$, $|\epsilon_n| < \alpha$

$$\text{Regret} = \mathbb{E} \left[\sum_{n=1}^N \theta^\star \cdot x^\star - \sum_{n=1}^N \theta^\star x_n \right]$$

Important Lemma:

$$E \left[\varepsilon_i \mid x_1, \dots, x_i \right] = 0$$

Lemma [Self Normalized Bound for Vector-Valued Martingales] Suppose $\{\varepsilon_n\}_{n=1}^{\infty}$ are mean zero random variables with $|\varepsilon_n| \leq \alpha$, for all n ; Let $\{x_i \in \mathbb{R}^d\}_{n=1}^{\infty}$ be some stochastic random process; Define $\Lambda^n = \lambda I + \sum_{i=1}^n x_i x_i^\top$, then with probability at least

$$1 - \delta, \text{ for all } n \geq 1: \left\| \underbrace{\sum_{i=1}^n x_i \varepsilon_i}_{\underbrace{\quad}} \right\|_{(\Lambda^n)^{-1}}^2 \leq 2\sigma^2 \ln \left(\frac{\det(\Lambda^n)^{1/2} \det(\lambda I)^{-1/2}}{\delta} \right)$$

$$\left(\sum_{i=1}^n x_i \varepsilon_i \right)^\top \Lambda_n^{-1} \left(\sum_{i=1}^n x_i \varepsilon_i \right)$$

Important Lemma:

$$\|x_i\|_2 \leq 1$$

Lemma [Self Normalized Bound for Vector-Valued Martingales] Suppose $\{\epsilon_n\}_{n=1}^{\infty}$ are mean zero random variables with $|\epsilon_n| \leq \alpha$, for all n ; Let $\{x_i \in \mathbb{R}^d\}_{n=1}^{\infty}$ be some stochastic random process; Define $\Lambda^n = \lambda I + \sum_{i=1}^n x_i x_i^\top$, then with probability at least

$$1 - \delta, \text{ for all } n \geq 1: \left\| \sum_{i=1}^n x_i \epsilon_i \right\|_{(\Lambda^n)^{-1}}^2 \leq 2\sigma^2 \ln \left(\frac{\det(\Lambda^n)^{1/2} \det(\lambda I)^{-1/2}}{\delta} \right)$$

$\approx \sigma^2 d \ln(n/\delta)$

$$\det(\Lambda^n) \leq (n + \lambda)^d \approx \sigma^2 d \ln(n) \checkmark$$

$$2\sigma^2 \ln(\det(\Lambda^n)^{1/2} \det(\lambda I)^{-1/2} / \delta) \leq \sigma^2 (d \ln(1 + n/\lambda) + 2 \ln(1/\delta))$$

Today's question

We extended MAB to linear bandit so that we can deal w/ infinitely many actions...

Can we extend discrete MDPs to some kind linear MDPs?

Outline for this lecture:

1. Introduction of low-rank MDP
2. Planning in low-rank MDP (i.e., DP) and UCBVI algorithm
3. Non-parametric model learning in linear MDPs

Notations and Useful Inequalities

For real-value matrix A :

$$\|A\|_F^2 = \sum_{i,j} A_{i,j}^2 \qquad \|A\|_2 = \sup_{x:\|x\|_2 \leq 1} \|Ax\|_2$$

Notations and Useful Inequalities

For real-value matrix A :

$$\|A\|_F^2 = \sum_{i,j} A_{i,j}^2$$

$$\|A\|_2 = \sup_{x:\|x\|_2 \leq 1} \|Ax\|_2$$

$$\|Ax\|_2 \leq \|A\|_2 \|x\|_2$$

Notations and Useful Inequalities

For real-value matrix A :

$$\|A\|_F^2 = \sum_{i,j} A_{i,j}^2 \qquad \|A\|_2 = \sup_{x:\|x\|_2 \leq 1} \|Ax\|_2$$

$$\|Ax\|_2 \leq \|A\|_2 \|x\|_2$$

For Positive Definite matrix Λ :

Notations and Useful Inequalities

For real-value matrix A :

$$\|A\|_F^2 = \sum_{i,j} A_{i,j}^2 \qquad \|A\|_2 = \sup_{x:\|x\|_2 \leq 1} \|Ax\|_2$$

$$\|Ax\|_2 \leq \|A\|_2 \|x\|_2$$

For Positive Definite matrix Λ :

$$\sigma_{\min}(\Lambda), \sigma_{\max}(\Lambda)$$

Notations and Useful Inequalities

For real-value matrix A :

$$\|A\|_F^2 = \sum_{i,j} A_{i,j}^2 \quad \|A\|_2 = \sup_{x:\|x\|_2 \leq 1} \|Ax\|_2$$

$$\|Ax\|_2 \leq \|A\|_2 \|x\|_2$$

For Positive Definite matrix Λ :

$$\sigma_{\min}(\Lambda), \sigma_{\max}(\Lambda) \quad \det(\Lambda) = \prod_{i=1}^d \sigma_i$$

Notations and Useful Inequalities

For real-value matrix A :

$$\|A\|_F^2 = \sum_{i,j} A_{i,j}^2 \quad \|A\|_2 = \sup_{x:\|x\|_2 \leq 1} \|Ax\|_2$$

$$\|Ax\|_2 \leq \|A\|_2 \|x\|_2$$

For Positive Definite matrix Λ :

$$\sigma_{\min}(\Lambda), \sigma_{\max}(\Lambda) \quad \det(\Lambda) = \prod_{i=1}^d \sigma_i$$

$$\|x\|_{\Lambda}^2 = x^{\top} \Lambda x$$

Notations and Useful Inequalities

For real-value matrix A :

$$\|A\|_F^2 = \sum_{i,j} A_{i,j}^2 \quad \|A\|_2 = \sup_{x:\|x\|_2 \leq 1} \|Ax\|_2$$

$$\|Ax\|_2 \leq \|A\|_2 \|x\|_2$$

For Positive Definite matrix Λ :

$$\sigma_{\min}(\Lambda), \sigma_{\max}(\Lambda) \quad \det(\Lambda) = \prod_{i=1}^d \sigma_i$$

$$\|x\|_{\Lambda}^2 = x^{\top} \Lambda x$$

$$\mathbb{E}_{s' \sim P_h(\cdot | s, a)} f(s') = P_h(\cdot | s, a) \cdot f$$

Low-Rank MDP Definition

Low-Rank MDP Definition

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

Low-Rank MDP Definition

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

S & A could be large or even continuous, hence $\text{poly}(S, A)$ is not acceptable

Low-Rank MDP Definition

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

S & A could be large or even continuous, hence $\text{poly}(S, A)$ is not acceptable

Low-Rank
Decomposition:

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

$d \ll |S||A|$

Low-Rank MDP Definition

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

S & A could be large or even continuous, hence $\text{poly}(S, A)$ is not acceptable

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

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

Low-Rank Decomposition:

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

Low-Rank MDP Definition

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

S & A could be large or even continuous, hence $\text{poly}(S, A)$ is not acceptable

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

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

Low-Rank
Decomposition:

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

$\text{poly}(d)$ rather than $\text{poly}(SA)$

Linear MDP Definition

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

S & A could be large or even continuous, hence $\text{poly}(S, A)$ is not acceptable

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

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

is known

Feature map ϕ is known to the learner!

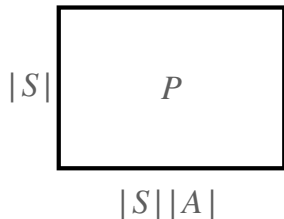
(We assume reward is known, i.e., θ^* is known)

Linear MDP Example

It generalizes tabular MDPs: $\phi(s, a)$ one-hot vector $\mathbb{R}^{|S||A|}$

$$P(\cdot | s, a) = P\phi(s, a)$$

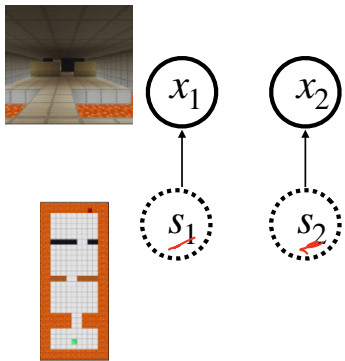
where $P \in \mathbb{R}^{|S| \times |S||A|}$ is the transition matrix



$$\text{Rank}(P) \leq \sum_{i=1}^{|S|} 1$$

Low-Rank Example

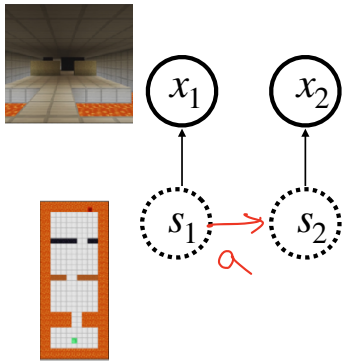
Can encode latent variables: block-MDPs



Low-Rank Example

Can encode latent variables: block-MDPs

Discrete latent state space S : $|S|$ is small, transition $T : S \times A \mapsto S$

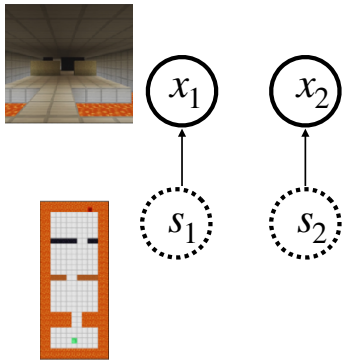


Low-Rank Example

Can encode latent variables: block-MDPs

Discrete latent state space S : $|S|$ is small, transition $T : S \times A \mapsto S$

Large observation space X (hence any poly dependency on $|X|$ is bad)



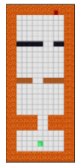
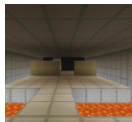
Low-Rank Example

Can encode latent variables: block-MDPs

Discrete latent state space S : $|S|$ is small, transition $T : S \times A \mapsto S$

Large observation space X (hence any poly dependency on $|X|$ is bad)

Each state s has an emission distribution $\mu_s \in \Delta(X)$, also μ_s and $\mu_{s'}$ have **disjoint support** for any $s \neq s'$ (i.e, latent state is decodable)



Low-Rank Example

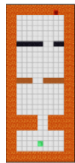
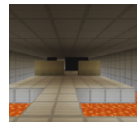
Can encode latent variables: block-MDPs

Discrete latent state space S : $|S|$ is small, transition $T : S \times A \mapsto S$

Large observation space X (hence any poly dependency on $|X|$ is bad)

Each state s has an emission distribution $\mu_s \in \Delta(X)$, also μ_s and $\mu_{s'}$ have **disjoint support** for any $s \neq s'$ (i.e, latent state is decodable)

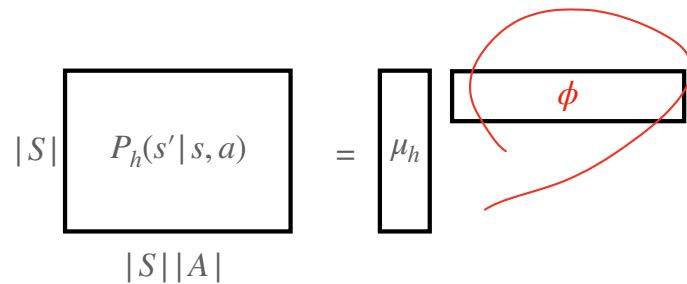
$$P(x'|x, a) = \sum_{s' \in \{s_1, s_2, s_3\}} T(s' | \omega(x), a) \mu_{s'}(x') = [\mu_{s_1}(x'), \mu_{s_2}(x'), \mu_{s_3}(x')] \begin{bmatrix} T(s_1 | \omega(x), a) \\ T(s_2 | \omega(x), a) \\ T(s_3 | \omega(x), a) \end{bmatrix} \in \mathbb{R}^3$$



ω

We only study Linear MDPs here (i.e., low-rank + known ϕ).
Learning in Low-rank MDP is much harder (coming later!)

Low-Rank
Decomposition:

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


Outline for this lecture:



1. Introduction of low-rank MDP

2. Planning in low-rank MDP (i.e., DP) and the UCBVI algorithm

3. Non-parametric model learning in linear MDPs

Planning in Linear MDP: Value Iteration

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

$$\mu_h^* \in \mathbb{R}^{|S| \times d} \quad r_h(s, a) = (\theta_h^*)^\top \phi(s, a), \quad \theta_h^* \in \mathbb{R}^d$$

Planning in Linear MDP: Value Iteration

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

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

$$V_H^\star(s) = 0, \forall s,$$

Planning in Linear MDP: Value Iteration

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

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

$$V_H^\star(s) = 0, \forall s,$$

$$Q_h^\star(s, a) = r_h(s, a) + \mathbb{E}_{s' \sim P_h(\cdot | s, a)} V_{h+1}^\star(s')$$

$$\underbrace{P_h(\cdot | sa)}^\top V_{h+1}^\star$$

$$= (\mu_h^\star \phi(s, a))^\top V_{h+1}^\star = \phi(s, a)^\top (\mu_h^\star{}^\top V_{h+1}^\star)$$

Planning in Linear MDP: Value Iteration

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

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

$$V_H^\star(s) = 0, \forall s,$$

$$\begin{aligned} Q_h^\star(s, a) &= r_h(s, a) + \mathbb{E}_{s' \sim P_h(\cdot | s, a)} V_{h+1}^\star(s') \\ &= \theta_h^\star \cdot \phi(s, a) + (\mu_h^\star \phi(s, a))^\top V_{h+1}^\star \end{aligned}$$

Planning in Linear MDP: Value Iteration

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

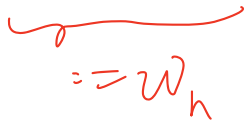
$$r_h(s, a) = (\theta_h^\star)^\top \phi(s, a), \quad \theta_h^\star \in \mathbb{R}^d$$

$$V_H^\star(s) = 0, \forall s,$$

$$Q_h^\star(s, a) = r_h(s, a) + \mathbb{E}_{s' \sim P_h(\cdot | s, a)} V_{h+1}^\star(s')$$

$$= \theta_h^\star \cdot \phi(s, a) + (\mu_h^\star \phi(s, a))^\top V_{h+1}^\star$$

$$= \phi(s, a)^\top (\theta_h^\star + (\mu_h^\star)^\top V_{h+1}^\star)$$



A red handwritten bracket underlines the term $(\theta_h^\star + (\mu_h^\star)^\top V_{h+1}^\star)$ in the previous equation. Below the bracket, the text $:= \mathcal{W}_h$ is written in red, defining the bracketed term as \mathcal{W}_h .

$$:= \mathcal{W}_h$$

Planning in Linear MDP: Value Iteration

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

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

$$V_H^\star(s) = 0, \forall s,$$

$$\begin{aligned} Q_h^\star(s, a) &= r_h(s, a) + \mathbb{E}_{s' \sim P_h(\cdot | s, a)} V_{h+1}^\star(s') \\ &= \theta_h^\star \cdot \phi(s, a) + (\mu_h^\star \phi(s, a))^\top V_{h+1}^\star \\ &= \phi(s, a)^\top (\theta_h^\star + (\mu_h^\star)^\top V_{h+1}^\star) \\ &= \phi(s, a)^\top w_h \end{aligned}$$

Planning in Linear MDP: Value Iteration

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

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

$$V_H^\star(s) = 0, \forall s,$$

$$\begin{aligned} Q_h^\star(s, a) &= r_h(s, a) + \mathbb{E}_{s' \sim P_h(\cdot | s, a)} V_{h+1}^\star(s') \\ &= \theta_h^\star \cdot \phi(s, a) + (\mu_h^\star \phi(s, a))^\top V_{h+1}^\star \\ &= \phi(s, a)^\top (\theta_h^\star + (\mu_h^\star)^\top V_{h+1}^\star) \\ &= \phi(s, a)^\top w_h \end{aligned}$$

$$V_h^\star(s) = \max_a \phi(s, a)^\top w_h, \quad \pi_h^\star(s) = \arg \max_a \phi(s, a)^\top w_h$$

Planning in Linear MDP: Value Iteration

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

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

$$V_H^\star(s) = 0, \forall s,$$

$$\begin{aligned} Q_h^\star(s, a) &= r_h(s, a) + \mathbb{E}_{s' \sim P_h(\cdot | s, a)} V_{h+1}^\star(s') \\ &= \theta_h^\star \cdot \phi(s, a) + \left(\mu_h^\star \phi(s, a) \right)^\top V_{h+1}^\star \\ &= \phi(s, a)^\top \left(\theta_h^\star + (\mu_h^\star)^\top V_{h+1}^\star \right) \\ &= \phi(s, a)^\top w_h \end{aligned}$$

Indeed we can show that $Q_h^\pi(\cdot, \cdot)$

Is linear with respect to ϕ as well, for any π, h

$$V_h^\star(s) = \max_a \phi(s, a)^\top w_h, \quad \pi_h^\star(s) = \arg \max_a \phi(s, a)^\top w_h$$

UCBVI in Linear MDPs

At the beginning of iteration n :

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 $\{s_h^i, a_h^i, s_{h+1}^i\}_{i=0}^{n-1}$



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 $\{s_h^i, a_h^i, s_{h+1}^i\}_{i=0}^{n-1}$

2. Design reward bonus $b_h^n(s, a), \forall s, a$

τ

$\otimes \text{poly}(H, d)$

set ϕ_{r_0} ←

be one-hot encoding

$\Rightarrow \sqrt{\frac{1}{N(s, a)}}$

$\sqrt{\frac{\phi^\top \Sigma^{-1} \phi}{\Delta}}$

$\Sigma = \sum_{i=1}^{n-1} \phi_i \phi_i^\top$

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 $\{s_h^i, a_h^i, s_{h+1}^i\}_{i=0}^{n-1}$

2. Design reward bonus $b_h^n(s, a), \forall s, a$

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

Outline for this lecture:



1. Introduction of low-rank MDP



2. Planning in low-rank MDP (i.e., DP) and the UCBVI algorithm

3. Non-parametric model learning in linear MDPs

Additional Assumptions in Linear MDPs to permit linear regression analysis

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

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

Additional Assumptions in Linear MDPs to permit linear regression analysis

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

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

Norm bounds:

$$\sup_{s,a} \|\phi(s, a)\|_2 \leq 1, \quad \|\theta_h^\star\|_2 \leq W, \quad \|v^\top \mu_h^\star\|_2 \leq \sqrt{d}, \quad \forall v \text{ s.t. } \|v\|_\infty \leq 1$$

1. Model Learning in Linear MDPs

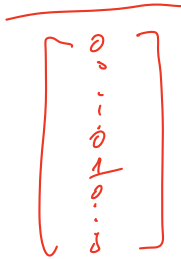
$$\mathcal{D}_h^n = \left\{ s_h^i, a_h^i, s_{h+1}^i \right\}_{i=1}^{n-1} \quad \Lambda_h^n = \sum_{i=1}^{n-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top + \lambda I$$

△

1. Model Learning in Linear MDPs

$$\mathcal{D}_h^n = \left\{ s_h^i, a_h^i, s_{h+1}^i \right\}_{i=1}^{n-1} \quad \Lambda_h^n = \sum_{i=1}^{n-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top + \lambda I$$

Denote $\delta(s) \in \mathbb{R}^{|S|}$ with zero everywhere except the entry corresponding to s




A hand-drawn red vector $\begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$ enclosed in red square brackets. The vector has a '1' in the middle and '0's and vertical ellipses above and below it, representing a unit vector in a specific dimension.

1. Model Learning in Linear MDPs

$$\mathcal{D}_h^n = \left\{ s_h^i, a_h^i, s_{h+1}^i \right\}_{i=1}^{n-1} \quad \Lambda_h^n = \sum_{i=1}^{n-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top + \lambda I$$

Denote $\delta(s) \in \mathbb{R}^{|S|}$ with zero everywhere except the entry corresponding to s

Given s, a , note that $\mathbb{E}_{s' \sim P_h(\cdot | s, a)} [\delta(s')] = P_h(\cdot | s, a) = \mu_h^* \phi(s, a)$



1. Model Learning in Linear MDPs

$$\mathcal{D}_h^n = \left\{ s_h^i, a_h^i, s_{h+1}^i \right\}_{i=1}^{n-1} \quad \Lambda_h^n = \sum_{i=1}^{n-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top + \lambda I$$

Denote $\delta(s) \in \mathbb{R}^{|S|}$ with zero everywhere except the entry corresponding to s

Given s, a , note that $\mathbb{E}_{s' \sim P_h(\cdot | s, a)} [\delta(s')] = P_h(\cdot | s, a) = \mu_h^* \phi(s, a)$

Denote $\epsilon_{s,a} = \delta(s') - P_h(\cdot | s, a)$, we have $\mathbb{E}_{s' \sim P_h(\cdot | s, a)} [\epsilon_{s,a}] = 0$, and $\|\epsilon_{s,a}\|_1 \leq 2$

$$\epsilon_p |s|$$

$$\|\epsilon_{s,a}\|_2 \leq \|\delta(s')\|_2 + \|P_h(\cdot | s, a)\|_2$$

1. Model Learning in Linear MDPs

$$\mathcal{D}_h^n = \left\{ s_h^i, a_h^i, s_{h+1}^i \right\}_{i=1}^{n-1} \quad \Lambda_h^n = \sum_{i=1}^{n-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top + \lambda I$$

Denote $\delta(s) \in \mathbb{R}^{|S|}$ with zero everywhere except the entry corresponding to s

Given s, a , note that $\mathbb{E}_{s' \sim P_h(\cdot | s, a)} [\delta(s')] = P_h(\cdot | s, a) = \mu_h^* \phi(s, a)$

Denote $\epsilon_{s,a} = \delta(s') - P_h(\cdot | s, a)$, we have $\mathbb{E}_{s' \sim P_h(\cdot | s, a)} [\epsilon_{s,a}] = 0$, and $\|\epsilon_{s,a}\|_1 \leq 2$

Ridge Linear Regression:

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

1. Model Learning in Linear MDPs

$$\mathcal{D}_h^n = \left\{ s_h^i, a_h^i, s_{h+1}^i \right\}_{i=1}^{n-1} \quad \Lambda_h^n = \sum_{i=1}^{n-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top + \lambda I$$

Denote $\delta(s) \in \mathbb{R}^{|S|}$ with zero everywhere except the entry corresponding to s

Given s, a , note that $\mathbb{E}_{s' \sim P_h(\cdot | s, a)} [\delta(s')] = P_h(\cdot | s, a) = \mu_h^* \phi(s, a)$

Denote $\epsilon_{s,a} = \delta(s') - P_h(\cdot | s, a)$, we have $\mathbb{E}_{s' \sim P_h(\cdot | s, a)} [\epsilon_{s,a}] = 0$, and $\|\epsilon_{s,a}\|_1 \leq 2$

Ridge Linear Regression:

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

$$\hat{\mu}_h^n = \sum_{i=1}^{n-1} \underbrace{\delta(s_{h+1}^i)}_{\mathbb{R}^{|S| \times 2}} \underbrace{\phi(s_h^i, a_h^i)^\top}_{\mathbb{R}^{1 \times d}} (\underbrace{\Lambda_h^n^{-1}}_{\mathbb{R}^{d \times d}})$$

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

1. Model Learning in Linear MDPs

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

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

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

1. Model Learning in Linear MDPs

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

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

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

Can we bound the ℓ_1 error on distributions, i.e., $\|\hat{P}_h^n(\cdot | s, a) - P(\cdot | s, a)\|_1$?

$$\leq \sqrt{\frac{|S|}{N(s,a)}}$$

1. Model Learning in Linear MDPs

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

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

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

Can we bound the ℓ_1 error on distributions, i.e., $\|\hat{P}_h^n(\cdot | s, a) - P(\cdot | s, a)\|_1$?

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

1. Model Learning in Linear MDPs

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

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

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

Can we bound the ℓ_1 error on distributions, i.e., $\|\widehat{P}_h^n(\cdot | s, a) - P(\cdot | s, a)\|_1$?

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

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

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

1. Model Learning in Linear MDPs

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

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

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

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:

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

1. Model Learning in Linear MDPs

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

1. Model Learning in Linear MDPs

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

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

1. Model Learning in Linear MDPs

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

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

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

1. Model Learning in Linear MDPs

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

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

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

1. Model Learning in Linear MDPs

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

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

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

1. Model Learning in Linear MDPs

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

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

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

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

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

1. Model Learning in Linear MDPs

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

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

1. Model Learning in Linear MDPs

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

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

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

1. Model Learning in Linear MDPs

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

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

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

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

1. Model Learning in Linear MDPs

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

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

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

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

1. Model Learning in Linear MDPs

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

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

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

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

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

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

1. Model Learning in Linear MDPs

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

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

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

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

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

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

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

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

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

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

1. Model Learning in Linear MDPs

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

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

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

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

1. Model Learning in Linear MDPs

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

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

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

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

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

1. Model Learning in Linear MDPs

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

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

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

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

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

1. Model Learning in Linear MDPs

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

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

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

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

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:

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

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:

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

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

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:

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

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

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

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:

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

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

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

2. Reward Bonus Design

Lemma [Model Average Error under a fixed V]:

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

$$b_h^n(s, a) = \beta \sqrt{\phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s, a)}, \quad \beta = \widetilde{O}(dH)$$

Next lecture: reward bonus design + regret bound

Summary for today:

1. Introduction of low-rank / Linear MDPs (linear Q^* , Q^π in feature ϕ)

2. Model-fitting in low-rank MDP (non-parametric regression)

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

3. Average model-error (over a fixed function V):

Summary for today:

1. Introduction of low-rank / Linear MDPs (linear Q^* , Q^π in feature ϕ)

2. Model-fitting in low-rank MDP (non-parametric regression)

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

3. Average model-error (over a fixed function V):

$$\left| \left(\hat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V \right| \leq \tilde{O} \left(H\sqrt{d} \right) \cdot \left\| \phi(s, a) \right\|_{(\Lambda_h^n)^{-1}}$$