

Exploration in Linear MDPs (Part II)

Recap: 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

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

$$P_h(s'|s, a) = \mu^\star \cdot \phi(s, a)$$

$$= \begin{bmatrix} | & | & \dots & | \\ \mu_1 & \mu_2 & \dots & \mu_d \\ | & | & & | \end{bmatrix} \phi(s, a)$$

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

Feature map ϕ is known to the learner!

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

Recap: Model Learning in Linear MDPs

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

Recap: Model Learning in Linear MDPs

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

Ridge Linear Regression:

one-hot vector

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

Recap: Model Learning in Linear MDPs

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

Ridge Linear Regression:

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

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

Recap: Model Learning in Linear MDPs

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

Ridge Linear Regression:

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

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

$$\hat{\mu}_h^n - \mu_h^* = -\lambda \mu_h^* (\Lambda_h^n)^{-1} + \underbrace{\sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}}_{\text{①}}, \quad \epsilon_h^i = \delta(s_{h+1}^i) - P_h(\cdot | s_h^i, a_h^i)$$

$$\mathbb{E}[\epsilon_h^i | s_h^i, a_h^i] = 0 \quad \|\epsilon_h^i\|_1 \leq 2$$

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:

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

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:

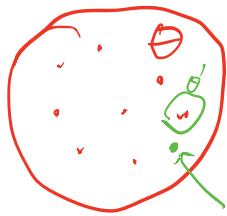
$$\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V \right| \lesssim \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \cdot H \left(\sqrt{\ln \frac{H}{\delta}} + \sqrt{d \ln \left(1 + \frac{N}{\lambda} \right)} \right)$$

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

$$|\mathcal{F}| = \infty$$

Covering

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

$$\|\theta' - \theta\|_2 \leq \epsilon$$

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 $|\mathcal{N}_\epsilon| \leq (1 + \underbrace{2R/\epsilon})^d$, 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

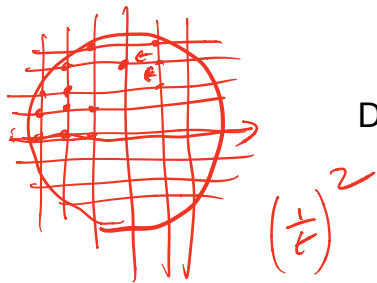
Lemma [Covering of Θ] We have $|\mathcal{N}_\epsilon| \leq (1 + 2R/\epsilon)^d$, and $\ln(|\mathcal{N}_\epsilon|) \leq d \ln(1 + 2R/\epsilon)$

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

$$\|f\|_\infty = \sup_{x \in \mathcal{X}} |f(x)|$$

$$f_{\theta_1}(x) = x^\top \cdot \theta_1$$
$$f_{\theta_2}(x) = x^\top \cdot \theta_2$$

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 $\forall \theta \in \Theta$:

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

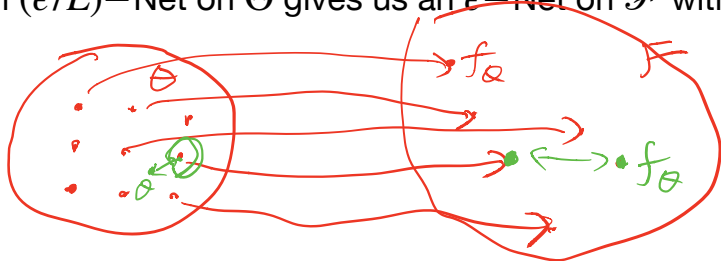
Denote ϵ -cover as the smallest \mathcal{N}_ϵ

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 $d(f_{\theta_1}, f_{\theta_2}) := \|f_{\theta_1} - f_{\theta_2}\|_\infty$

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



Lipshitz.

Detour: Covering Number and An Example

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

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

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

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

$$f_{w, \beta, \Lambda} : S \mapsto \mathbb{R}^+$$

Detour: Covering Number and An Example

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

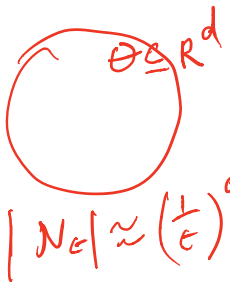
Define the function $f_{w, \beta, \Lambda} : S \rightarrow [0, H]$

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

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_∞**

Detour: Covering Number and Example

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



Detour: Covering Number and Example

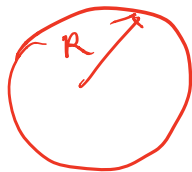
$$f_{w,\beta,\Lambda} \text{ : } f_{w,\beta,\Lambda}(s) := \min \left\{ \max_a \left(w^\top \phi(s,a) + \beta \sqrt{\phi(s,a)^\top \Lambda^{-1} \phi(s,a)} \right), H \right\}$$

$d \quad 1 \quad d^2$

$\Lambda \in \mathbb{R}^{d \times d}$

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

Detour: Covering Number and Example



$$|\mathcal{N}_d = \left(1 + \frac{R}{\epsilon}\right)^d$$

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

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

$\forall s \in \mathcal{S}$

Key step in the proof:

$$|f_\theta(s) - f_{\hat{\theta}}(s)| \leq \|w - \hat{w}\|_2 + |\beta - \hat{\beta}|/\sqrt{\lambda} + B\sqrt{\|\Lambda^{-1} - \hat{\Lambda}^{-1}\|_F}$$

(w, β, Λ) $(\hat{w}, \hat{\beta}, \hat{\Lambda})$ $\frac{\epsilon}{3} - \text{Net}$ $\sqrt{\lambda} \epsilon$ $\frac{\epsilon^2}{B^2} - \text{Net} \left\{ A \in \mathbb{R}^{d \times d} : \sigma_{\max}(A) \leq \frac{\epsilon}{\lambda} \right\}$

Detour: Uniform Convergence using Covering

Define $f_{w,\beta,\Lambda} : S \rightarrow [0,H]$, $f_{w,\beta,\Lambda}(s) := \min \left\{ \max_a \left(w^\top \phi(s, a) + \beta \sqrt{\phi(s, a)^\top \Lambda^{-1} \phi(s, a)} \right), H \right\}$

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

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

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

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

$$|\mathcal{F}| = \infty$$

$$\| \widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \|_1 \approx \sqrt{\frac{5 \ln(1/\delta)}{\dots}}$$

Detour: Uniform Convergence using Covering

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

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

Detour: Uniform Convergence using Covering

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: let's start with \mathcal{N}_ϵ over \mathcal{F}

$$\ln |\mathcal{N}_\epsilon| \approx d^2$$

Detour: Uniform Convergence using Covering

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: let's start with \mathcal{N}_ϵ over \mathcal{F}

1. For a fixed $\hat{f} \in \mathcal{N}_\epsilon$: $\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot \hat{f} \right| \lesssim \left\| \phi(s, a) \right\|_{(\Lambda_h^n)^{-1}} \cdot H \left(\sqrt{\ln \frac{H}{\delta}} + \sqrt{d \ln \left(1 + \frac{N}{\lambda} \right)} \right)$

Detour: Uniform Convergence using Covering

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

$$\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot f \right| = \underbrace{\widetilde{O}(Hd)}_{AA} \cdot \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \quad \ln |\mathcal{N}_\epsilon| \leq d^2$$

Proof Sketch: let's start with \mathcal{N}_ϵ over \mathcal{F}

1. For a fixed $\hat{f} \in \mathcal{N}_\epsilon$: $\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot \hat{f} \right| \lesssim \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \cdot H \left(\sqrt{\ln \frac{H}{\delta}} + \sqrt{d \ln \left(1 + \frac{N}{\lambda} \right)} \right)$

2. union bound: $\forall \hat{f} \in \mathcal{N}_\epsilon$: $\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V \right| \lesssim \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \cdot H \left(\sqrt{\ln \frac{H |\mathcal{N}_\epsilon|}{\delta}} + \sqrt{d \ln \left(1 + \frac{N}{\lambda} \right)} \right)$

failure $\frac{\delta}{|\mathcal{N}_\epsilon|}$

$\ln \frac{|\mathcal{N}_\epsilon|}{\delta} = \left(d^2 + \ln \frac{1}{\delta} \right) \approx o(d)$

Detour: Uniform Convergence using Covering

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: let's start with \mathcal{N}_ϵ over \mathcal{F}

1. For a fixed $\hat{f} \in \mathcal{N}_\epsilon$: $\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot \hat{f} \right| \lesssim \left\| \phi(s, a) \right\|_{(\Lambda_h^n)^{-1}} \cdot H \left(\sqrt{\ln \frac{H}{\delta}} + \sqrt{d \ln \left(1 + \frac{N}{\lambda} \right)} \right)$

2. union bound: $\forall \hat{f} \in \mathcal{N}_\epsilon$: $\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V \right| \lesssim \left\| \phi(s, a) \right\|_{(\Lambda_h^n)^{-1}} \cdot H \left(\sqrt{\ln \frac{H |\mathcal{N}_\epsilon|}{\delta}} + \sqrt{d \ln \left(1 + \frac{N}{\lambda} \right)} \right)$



3. Consider any $f \in \mathcal{F}$:

$$\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot f \right| \leq \underbrace{\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot \hat{f} \right|}_{\text{Step 2}} + \underbrace{\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot (f - \hat{f}) \right|}_{\epsilon}$$

Handwritten annotations: A red arrow points from $\hat{f} \in \mathcal{N}_\epsilon$ to the first term. A red arrow points from ϵ to the second term.

Detour: Uniform Convergence using Covering

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

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

Proof Sketch: let's start with \mathcal{N}_ϵ over \mathcal{F}

1. For a fixed $\hat{f} \in \mathcal{N}_\epsilon$: $\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot \hat{f} \right| \lesssim \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \cdot H \left(\sqrt{\ln \frac{H}{\delta}} + \sqrt{d \ln \left(1 + \frac{N}{\lambda} \right)} \right)$

2. union bound: $\forall \hat{f} \in \mathcal{N}_\epsilon$: $\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V \right| \lesssim \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \cdot H \left(\sqrt{\ln \frac{H|\mathcal{N}_\epsilon|}{\delta}} + \sqrt{d \ln \left(1 + \frac{N}{\lambda} \right)} \right)$

3. Consider any $f \in \mathcal{F}$:

$$\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot f \right| \leq \underbrace{\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot \hat{f} \right|}_{Hd \sqrt{\ln(\frac{1}{\epsilon})}} + \underbrace{\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot (f - \hat{f}) \right|}_{\leftarrow}$$

4. Tune parameter ϵ

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)

$$\ln(N_{\epsilon}) \approx O(d^2)$$

$$\begin{array}{c} f_{w,p,\eta} \\ \uparrow \quad \uparrow \quad \uparrow \\ d \quad 1 \quad d \times d \end{array}$$

Summary of Covering Argument

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

Let's get back to Linear MDPs again!

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

$$\hat{P}_h^n(\cdot | s, a) \notin \Delta(S)$$

$$\hat{P}_h^n(s' | s, a) = \frac{N(s, a, s')}{N(s, a)}$$

Algorithm: UCBVI in Linear MDPs

Algorithm: UCBVI in Linear MDPs

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

Algorithm: UCBVI in Linear MDPs

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

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

Algorithm: UCBVI in Linear MDPs

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

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

Covering

$\beta = \widetilde{O}(Hd)$, 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

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

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

$\beta = \widetilde{O}(Hd)$, 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: $\pi^{n+1} = \text{Value-Iter} \left(\left(\underbrace{\{\widehat{P}_h^n\}_h}_{\Delta}, \underbrace{\{r_h + b_h^n\}}_{\sim} \right) \right)$

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

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

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

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

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

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

Δ

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

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

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

$$\mathbb{E}_{s' \sim \widehat{P}^n(\cdot | s, a)} \widehat{V}_{n+1}^n(s')$$

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

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

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

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

$$\rightarrow \widehat{\mu}_h^n \phi(s, a)$$

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

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

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

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

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

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

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

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

$$\widehat{\omega}_h^n \in \mathbb{R}^d$$

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

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

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

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

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

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

$$= \beta \sqrt{\phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s, a)} + \phi(s, a)^\top \widehat{w}_h^n \quad \leftarrow \begin{matrix} \text{:} \\ \Rightarrow \end{matrix} \widehat{w}_h^n$$

~~_____~~

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

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

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

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

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

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

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

$$V_h^*(s) = \min \left\{ \max_a \left(\underbrace{\phi(s, a)^\top w_h + \beta \sqrt{\phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s, a)}}_{\widehat{Q}_h^n(s, a)}, H \right), H \right\}, \quad \pi_h^*(s) = \arg \max_a \widehat{Q}_h^n(s, a)$$

$$\mathcal{F}_{w, \beta, \Lambda}$$

$$= \min \left\{ \max_a w^\top \phi + \beta \sqrt{\dots}, H \right\}$$

$$\widehat{V}_h^n \in \mathcal{F}$$

$$\Delta$$

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

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: $\widehat{V}_{h+1}^n(s) \geq V_{h+1}^*(s), \forall s$

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

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

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: $\widehat{V}_{h+1}^n(s) \geq V_{h+1}^*(s), \forall s$

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

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: $\widehat{V}_{h+1}^n(s) \geq V_{h+1}^*(s), \forall s$

$$\begin{aligned} \widehat{Q}_h^n(s, a) - Q_h^*(s, a) &= b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n - \underbrace{P_h(\cdot | s, a)}_{\text{Valid Distribution}} \cdot V_{h+1}^* \\ &\geq \cancel{b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n - P_h(\cdot | s, a) \cdot \widehat{V}_{h+1}^n} \\ &\geq b_h^n(s, a) + \underbrace{\widehat{P}_h^n(\cdot | s, a)}_{\text{Valid Distribution}} \cdot V_{h+1}^* - P_h(\cdot | s, a) \cdot V_{h+1}^* \quad \checkmark \end{aligned}$$

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

$$\| \widehat{P}(\cdot | s, a) - P(\cdot | s, a) \|_1 \leq \sqrt{\frac{\text{slu}(s)}{N(s, a)}}$$

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: $\widehat{V}_{h+1}^n(s) \geq V_{h+1}^*(s), \forall s$

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

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

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

$$\begin{aligned} &\uparrow \\ &(Hd) \cdot \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \\ &= 0 \end{aligned}$$

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

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

$$\widehat{V}_{h+1}^n(s) = \min \left\{ \max_a w^T \phi + \beta \sqrt{\lambda^{-1} \phi}, \tau \right\}$$

covering argument

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!!!

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

$$\geq 0$$

4. Regret Decomposition

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

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

Regret
at ep- n

π^n under
 $(\widehat{P}_n, r + b_n^n)$

π^n under (P, r)

4. Regret Decomposition

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

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

(apply Simulation Lemma here)

4. Regret Decomposition

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

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

d_h^π (s.a.)

(apply Simulation Lemma here)

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

$=$ ← Tabular: Holder $\|\widehat{P} - P\|_1 \cdot \|\widehat{V}\|_\infty$

$$\leq H \cdot \sqrt{\frac{\sum \ln(1/\delta)}{N(s,a)}}$$

$$\underbrace{\left(\widehat{P} - P \right) \cdot V^*}_{\text{green}} + \underbrace{\left(\widehat{P} - P \right) \left(\widehat{V} - V^* \right)}_{\text{green}} \leftarrow$$

$$\leq \sqrt{\frac{\ln(1/\delta)}{N(s,a)}}$$

$$\begin{aligned} \widehat{P} &\rightarrow P \\ \widehat{V} &\rightarrow V^* \end{aligned}$$

$$\frac{1}{N(s,a)} \rightarrow \ln(N)$$

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)

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s_h, a_h \sim d_h^{\pi^n}} \left[b_h^n(s_h, a_h) + \underbrace{\left(\widehat{P}_h^n(\cdot | s_h, a_h) - P_h(\cdot | s_h, a_h) \right) \cdot \widehat{V}_{h+1}^n}_{\text{OUD} \cdot \|\phi(s_{-g})\| (\Lambda_h^n)^{-1} \approx b_h^n(s_e)}$$

$$\lesssim \sum_{h=0}^{H-1} \mathbb{E}_{s_h, a_h \sim d_h^{\pi^n}} \left[\underbrace{b_h^n(s_h, a_h)}_{\text{OUD}} \right]$$

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)

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

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

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

$$\left| \left(\hat{P}_n^n(c|\cdot|s_a) - P_h^n(c|\cdot|s_a) \right) \cdot f \right| \leq H d \|\phi(\cdot, a)\| \frac{1}{(\lambda_h^n)^T}, \forall s_a, n, h, f \in \mathcal{F}$$

4. Concluding the Regret Computation

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

4. Concluding the Regret Computation

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

4. Concluding the Regret Computation

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

\uparrow
 CH-inequality

4. Concluding the Regret Computation

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

$$\lesssim \beta \mathbb{E} \left[\sum_{n=1}^N \sum_{h=0}^{H-1} \sqrt{\phi(s_h^n, a_h^n)^\top (\Lambda_h^n)^{-1} \phi(s_h^n, a_h^n)} \right] + \delta NH$$

$$\lesssim \beta \mathbb{E} \left[\sum_{h=0}^{H-1} \sqrt{N} \sqrt{\sum_{n=1}^N \phi(s_h^n, a_h^n)^\top (\Lambda_h^n)^{-1} \phi(s_h^n, a_h^n)} \right] + \delta NH$$

$$\lesssim \widetilde{O}(H^2 d^{1.5} \sqrt{N})$$

$$O(d \ln(N))$$

$$\phi(s, a) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \in \mathbb{R}^{SA}$$

$$\in \mathbb{R}^{SA}$$

$$d = SA$$

$$H^2 (SA)^{1.5} \sqrt{N}$$

$$\text{set } \delta = \frac{1}{NH}$$

Linear Bandit $H=1$

$$d^{1.5} \sqrt{N}$$

Linear Bandit $\rightarrow d \sqrt{N}$

Tabular:

$\hat{P}(\cdot | S_n) \leftarrow$ linear regression on

$$\left| \left(\hat{P}(\cdot | S_n) - P(\cdot | S_n) \right) \cdot f \right| \quad \begin{array}{l} \nwarrow \text{pre-fixed} \\ f \quad \boxed{*} \end{array}$$

$\mathcal{F}, |\mathcal{F}| = \infty$

- ① covering ϵ -Net for \mathcal{F}
 - ② $\boxed{*}$ + Union Bound on ϵ -Net
 - ③ def ϵ -Net + Triangle Inequality
-

④ optimization

② simulation lemma

$$\textcircled{3} \sum_{n=1}^N \phi(S_n, a_n^T) \left(\mathbb{E}_n^n \right)^T \phi(S_n, a_n^a) \leq d \cdot \ln(N)$$

$$|(\hat{\theta} - \theta^*)^T x| \leq \underbrace{\|x\|_{\Lambda^{-1}}}_{\leq \sqrt{\alpha}} \sigma \sqrt{\alpha}$$

$$\Lambda = \sum_{i=1}^n x_i x_i^T$$

$$\Lambda = U \Sigma U^T$$

$\phi(s, \alpha)$