

# CS292F StatRL Lecture 10

## Exploration in Linear MDPs

Instructor: Yu-Xiang Wang

Spring 2021

UC Santa Barbara

# Recap: Lecture 9

- Exploration in Tabular MDPs
  - Problem setup: Episodic Finite-H MDP with non-stationary transitions.
  - Regret definition
- UCB-VI ([Azar et al., 2017](#))
  - A model-based approach; requires estimating  $P$ .
- Ideas in the proof
  - Concentration
  - Optimism via exploration bonus in value iterations
  - A few other tricks

# This lecture: Exploration in Reinforcement Learning

- Why is it challenging?
  - The reward depends on both  $s$ ,  $a$
  - Unlike the generative model setting, we cannot just choose any  $s$  to explore.
  - The data needs to be actively collected
- We will study
  - Tabular MDP
  - **Linear MDPs**
  - Both in the finite horizon episodic setting.

# Recap: episodic finite horizon MDPs with non-stationary transitions

- Problem setup / notations

- MDP:  $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \{r_h\}_h, \{P_h\}_h, H, s_0\}$

- Policy depends on time step

$$\pi = \{\pi_0, \dots, \pi_{H-1}\}$$

$\pi_h: \mathcal{S} \rightarrow \mathcal{A}$   
 $\Delta(\mathcal{A})$

- Performance measure

$$\text{Regret} := \mathbb{E} \left[ \sum_{k=0}^{K-1} \left( \underset{\uparrow}{V^*} - \underset{\uparrow}{V^{\pi^k}} \right) \right]$$

$V^* := V_0^{\pi^*}(s_0)$   
 $V^{\pi^k} := V_0^{\pi^k}(s_0)$

# Recap: Linear function approximation in TD-learning

- Why do we need it?
  - Recall the PACMAN example.

Let's say we discover through experience that this state is bad:



In naïve q-learning, we know nothing about this state:



Or even this one!



- Describe the state by its features, and value functions linear in features
  - $V_w(s) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_n f_n(s)$
  - $Q_w(s,a) = w_1 f_1(s,a) + w_2 f_2(s,a) + \dots + w_n f_n(s,a)$

# How do we formally analyze this?

- What are the assumptions to make?
  - $Q^*(s,a)$  approximately linear?
  - $Q^\pi(s,a)$  is approximately linear for all  $\pi$ ?
  - $Q^*(s,a)$  is exactly linear?
  - $Q^\pi(s,a)$  is exactly linear for all  $\pi$ ?

$$\|f - Q^*\|_\infty \leq \epsilon$$

↓

$$\pi_f = \text{argmax}_\pi f(s,a)$$

$$|V^* - V^{\pi_f}| \leq \frac{\epsilon}{1-\gamma}$$

$$\phi(s,a)$$

$$Q^*(s,a) = (\Theta^*)^T \phi(s,a)$$

Exponential sample complexity / regret lower bounds for the approximate case...

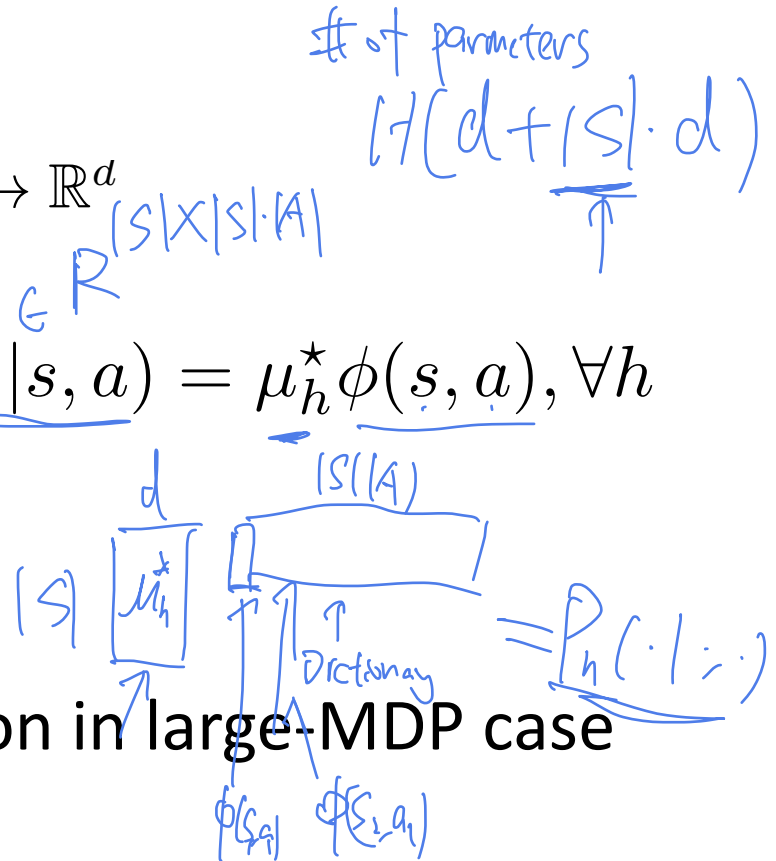
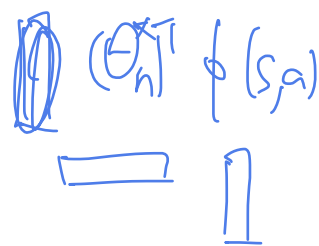
(Du, Kakade, Wang, Yang, 2019) Is a good representation sufficient for sample efficient reinforcement learning?

$\mathcal{N}(e^H)$

# Linear MDPs

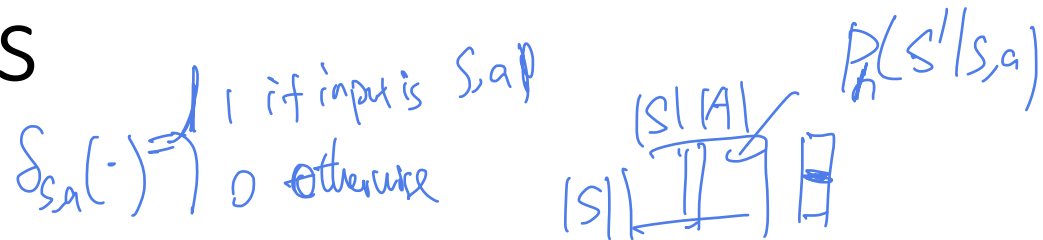
- Exists feature map  $\phi : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}^d$ 
  - Such that:

$$r_h(s, a) = \theta_h^* \cdot \phi(s, a), \quad \underline{P_h(\cdot | s, a)} = \underline{\mu_h^* \phi(s, a)}, \quad \forall h$$



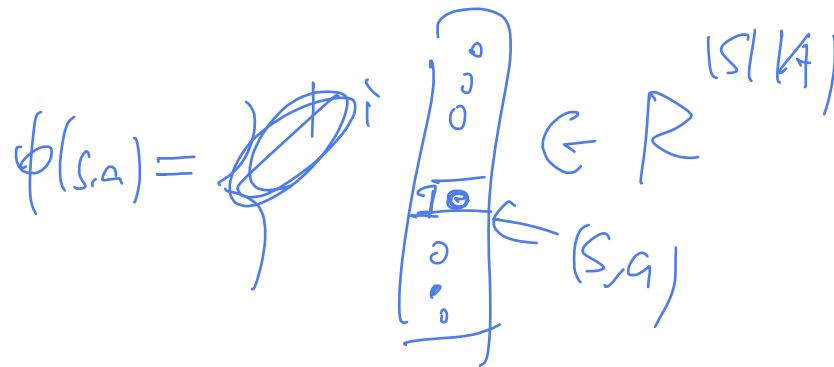
- Implies a low-rank assumption in large-MDP case

# Tabular MDPs are instances of linear MDPs



$$\underline{r_h(s, a)} = \theta_h^* \cdot \phi(s, a), \quad \underline{P_h(\cdot | s, a)} = \underline{\mu_h^* \phi(s, a)}, \forall h$$

- Choose  $d = |S| |A|$
- Feature map to be:





# Linear MDPs imply that the Q function for any policy is linear

- Optimal Q\* function:

$$\begin{aligned}
 Q_h^*(S,a) &= r(S,a) + P_n(\cdot|S,a) \cdot V_{h+1}^*(\cdot) \\
 &= \theta_h^* \phi(S,a) + (U_h^* \phi(S,a))^T V_{h+1}^*(\cdot) \\
 &= (\theta_h^* + U_h^{*T} V_{h+1}^*)^T \phi(S,a) = w^{*T} \phi(S,a)
 \end{aligned}$$

with any  
 $f: S \rightarrow \mathbb{R}$   
 $f: V_{h+1}^*(\cdot)$

- Claim 7.2 (AJKS) For any function of the state, the Bellman backup is a linear in the feature.

# Recap: LinUCB in linear bandits

- In every round:

1. Use ridge regression to estimate the model parameters

$$\hat{\mu}_t = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^{t-1} \underbrace{(x_i^\top \theta - r_i)^2}_{\text{residual}} + \lambda \underbrace{\|\theta\|^2}_{\text{regularizer}}$$

2. Construct an ellipsoidal confidence set of the parameters

$$\underline{\text{Ball}}_t = \left\{ \mu \mid (\mu - \hat{\mu}_t)^\top \Sigma_t (\mu - \hat{\mu}_t) \leq \beta_t \right\}$$

where  $\underline{\Sigma}_t = \sum_{i=1}^{t-1} x_i x_i^\top + \lambda I_d$



3. Choose actions that maximize UCB

$$X_t = \underset{x \in \mathcal{D}}{\operatorname{argmax}} \max_{\mu \in \underline{\text{Ball}}_t} \langle x, \mu \rangle$$

# UCB-VI for Linear MDPs

$\theta^*$  is known  
parameter  $\mu^* \in \mathbb{R}^{|S| \times d}$

• In every round:

1. Run Ridge regression for estimating the model

$\mu_h^n \in \mathbb{R}^{|S| \times d}$

$$\hat{\mu}_h^n = \operatorname{argmin}_{\mu \in \mathbb{R}^{|S| \times d}} \sum_{i=0}^{n-1} \|\mu \phi(s_h^i, a_h^i) - \delta(s_{h+1}^i)\|_2^2 + \lambda \|\mu\|_F^2.$$

Feature

$\|\cdot\|_F = \|\text{vec}(\cdot)\|_2$

label

$\phi(s_h^i, a_h^i) \in \mathbb{R}^{|S|}$

$\delta_{s_{h+1}^i} \in \mathbb{R}^{|S|}$

$$\hat{\mu}_h^n = \sum_{i=0}^{n-1} \delta(s_{h+1}^i) \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

$\Lambda_h^n = \sum_{i=0}^{n-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top$

$\mathbb{E}[\delta_{s_{h+1}^i}] = P(s_{h+1}^i)$

2. Construct the exploration bonuses

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

width of  $\phi(s, a)$

3. Run optimistic value iterations, and update greedy policy

# Optimistic value iterations

$$\begin{aligned}
 \hat{V}_H^n(s) &= 0, \forall s, \\
 \hat{Q}_h^n(s, a) &= \theta_h^* \cdot \phi(s, a) + \beta \sqrt{\phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s, a)} + \phi(s, a)^\top (\hat{\mu}_h^n)^\top \hat{V}_{h+1}^n \\
 &= \beta \sqrt{\phi(s, a)^\top (\Lambda_h^n)^{-1} \phi(s, a)} + (\theta_h^* + (\hat{\mu}_h^n)^\top \hat{V}_{h+1}^n)^\top \phi(s, a), \\
 \hat{V}_h^n(s) &= \min\{\max_a \hat{Q}_h^n(s, a), H\}, \quad \pi_h^n(s) = \operatorname{argmax}_a \hat{Q}_h^n(s, a).
 \end{aligned}$$

# Regret bound

• Choose  $\beta = Hd \left( \sqrt{\ln \frac{H}{\delta}} + \sqrt{\ln(W + H)} + \sqrt{\ln B} + \sqrt{\ln d} + \sqrt{\ln N} \right) \approx \tilde{O}(l+d)$   
 $\lambda = 1$

• Regret  $\tilde{O}(H^2 \sqrt{d^3 N})$   $N$  is the # of episodes

# Sketch of the regret analysis

Per-episode Regret:

$$V^* - V^{\pi_n} = V_0^*(s_0) - V_0^{\pi_n}(s_0) \stackrel{\text{Optimism}}{\leq} \underbrace{V_0^{\pi_n}(s_0) - V_0^{\pi_n}(s_0)}$$

Simulation Lemma  $\rightarrow \leq \sum_{h=0}^{H-1} \mathbb{E}^{\pi_n} \left[ b_h^n(S_h, A_h) + \left( P_h^n(\cdot | S_h, A_h) - P_h(\cdot | S_h, A_h) \right) \cdot \underbrace{V_{h+1}^n}_{\leq \frac{H-1}{h+2} \mathbb{E}^{\pi_n} [b_h^n(S_h, A_h)]} \right]$

Lemma (Uniform Concentration) w.p.  $1-\delta \forall s, a, h, n$ ,  
 $\sup_{f \in \mathcal{F}} \left( \sum_{h=0}^{H-1} \left( P_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot f \right) \leq \beta \| \phi(s, a) \|_{\mathcal{H}_n}^{n-1} =: b_h^n(s, a)$  if  $\mathcal{H}_n \in \mathcal{F}$

w.h.p.  $V^* - V^{\pi_n} \leq \mathbb{E}^{\pi_n} \left[ \sum_{h=0}^{H-1} b_h^n(S_h^n, A_h^n) \mid \text{hist}_n \right]$

Regret =  $\mathbb{E} \left[ \sum_{i=0}^{N-1} V^* - V^{\pi_n} \right] = \mathbb{E} \left[ \sum_{i=0}^{N-1} (V^* - V^{\pi_n}) \mathbb{I}(N_{\text{not fail}}) \right] + \mathbb{E} \left[ \sum_{i=0}^{N-1} (V^* - V^{\pi_n}) \mathbb{I}(F_{\text{fail}}) \right]$

$\leq \mathbb{E} \left[ \sum_{i=1}^{N-1} \sum_{h=0}^{H-1} b_h^n(S_h^i, A_h^i) \mathbb{I}(N_{\text{not fail}}) \right] + \underbrace{\delta \cdot N \cdot H}_{\uparrow \text{worst possible failure mode}}$

tabular case

# Sketch of the regret analysis

$$\sum_{h=0}^{H-1} \sum_{n=0}^{N-1} b_h(S_h^n, A_h^n) = \beta \sum_{h=0}^{H-1} \sum_{n=0}^{N-1} \sqrt{\phi(S_h^n, A_h^n)^T (\Lambda_h^n)^{-1} \phi(S_h^n, A_h^n)}$$

$$\leq \beta \sum_{h=0}^{H-1} \sqrt{N \sum_{n=0}^{N-1} \phi(S_h^n, A_h^n)^T (\Lambda_h^n)^{-1} \phi(S_h^n, A_h^n)} = d t \sqrt{N d \log U}$$

$$\sqrt{\frac{1}{N} \sum_{h=0}^{H-1} \sum_{n=0}^{N-1} \frac{1}{\sqrt{b_h^n}}$$

$$\sum_{i=1}^M \frac{1}{\sqrt{b_i}}$$

Lemma (Information gain bound)

$$\forall S_h^n, A_h^n \text{ sequence } \sum_{n=0}^{N-1} \phi(S_h^n, A_h^n)^T (\Lambda_h^n)^{-1} \phi(S_h^n, A_h^n) = \tilde{O}(d \log U)$$

$$\tilde{O}\left(\frac{1}{d^2} d^3 N\right)$$

$$\beta = \tilde{O}(t d)$$

# Sketch of the regret analysis



# It remains to prove

$$\sup_f |(\hat{P} - P) \cdot f| \leq \epsilon$$

$n \geq \frac{1}{\epsilon^2}$

- 1. Uniform convergence bound

- 2. “Optimism”

The same induction argument as in the UCB-VI for tabular MDP  
(Read Lemma 7.10 in AJKS)

- 3. “Information gain” bound

The same argument as in the Linear Bandits case.  
(Read Lemma 7.12 in AJKS)

# "Optimism" from Optimistic Value Iterations

$$\text{Claim: } \hat{V}_0^{\pi_n}(s_0) \geq V_0^*(s_0)$$

if  $\hat{V}_{h+1}^n(s) \geq V_{h+1}^*(s)$ ,

$$\hat{Q}_h^n(s_a) - Q_h^*(s_a) = \overbrace{r(s_a) + \beta \int \phi(s_a)^T (\hat{\Lambda}_h^n)^{-1} \phi(s_a)}^{-r(s_a)} + \phi(s_a)^T (\hat{\Lambda}_h^n)^T \hat{V}_{h+1}^n - \phi(s_a)^T (\Lambda_h^*)^T V_{h+1}^*$$

$\geq 0$

applying inductive hypothesis

$$\rightarrow \geq \beta \int \phi(s_a)^T (\hat{\Lambda}_h^n)^{-1} \phi(s_a) + \phi(s_a)^T (\hat{\Lambda}_h^n - \Lambda_h^*)^T \hat{V}_{h+1}^n$$

if we choose  $\beta$  s.t. w.h.p

$$\beta \int \phi(s_a)^T (\hat{\Lambda}_h^n)^{-1} \phi(s_a) \geq \left| \phi(s_a)^T (\hat{\Lambda}_h^n - \Lambda_h^*)^T \hat{V}_{h+1}^n \right|$$

$$\sup_{f \in F} \phi(s_a)^T (\hat{\Lambda}_h^n - \Lambda_h^*)^T f \leq \delta$$

$$\hat{\Lambda}_h^n \in \bar{F}$$

where  $F$  is finite, Woeffel's + union bound

$$\delta = \frac{\delta}{\prod_{j=1}^n}$$

# Let us start with pointwise convergence for a fixed $V$

- Recall: Hoeffding's inequality + Union bound

- Recall: Ridge regression

$\{(S_h^i, A_h^i, S_{h+1}^i) \text{ for } i=1, 2, \dots, n-1\}$

$$\hat{\mu}_h^n = \underset{\mu \in \mathbb{R}^{|S_h|}}{\text{argmin}} \sum_{i=0}^{n-1} \|\mu \circ \phi(S_h^i, A_h^i) - S_{h+1}^i(\cdot)\|_2^2 + \lambda \|\mu\|_F^2$$

$$\hat{\mu}_h^n = \frac{\sum_{i=1}^{n-1} S_{h+1}^i(\cdot) \cdot \phi(S_h^i, A_h^i)^T}{|S|} \cdot (\Lambda_h^n)^{-1} \quad \left\{ \begin{array}{l} \text{for a fixed } V \in \mathbb{R}^{|S|} \\ (\hat{P}_h^n(\cdot, S_a) - P_h(\cdot, S_a)) \cdot V \in \mathbb{R} \end{array} \right.$$

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

$$P_h(\cdot, S_a) = \mu^* \cdot \phi(S, a)$$

$$\phi(S_a)^T \left( \sum_{i=1}^{n-1} (\Lambda_h^n)^{-1} \phi(S_h^i, A_h^i) \langle S_{h+1}^i, V \rangle \right)$$

$$\langle S_{h+1}^i, V \rangle$$

19

# Error of ridge regression estimate

- Lemma 7.3 AJKS

$$\hat{\mu}_h^n - \mu_h^* = \underbrace{-\lambda \mu_h^* (\Lambda_h^n)^{-1}}_{\text{bias}} + \underbrace{\sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^T (\Lambda_h^n)^{-1}}_{\text{Variance}}.$$

$\epsilon_h^i = \underbrace{y_{h+1}^i - P(-|S_h^i, A_h^i)}_{\mu_h^* \cdot \phi(S_h^i, a_h^i)}$   
 $\epsilon_{\text{RISKE}}$

- The quantity of interest is a inner product with this:

$$\left[ (\hat{\mu}_h^n - \mu_h^*) \cdot \phi(s, a) \right]^T V = \phi(s, a)^T \underbrace{(\hat{\mu}_h^n - \mu_h^*)^T}_{\text{Variance } \sigma^2 \cdot V} \cdot V$$

$$\text{bias} = -\lambda \phi(s, a)^T (\Lambda_h^n)^{-1} \mu_h^{*T} V$$

$$= -\lambda \phi(s, a)^T \mu_h^* (\Lambda_h^n)^{-\frac{1}{2}} (\Lambda_h^n)^{-\frac{1}{2}} \mu_h^{*T} V$$

$$\leq \lambda \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \|\mu_h^{*T} V\|_{\Lambda_h^n} \leq \lambda \|\phi(s, a)\|_{(\Lambda_h^n)^{-1}} \sqrt{d} H$$

# Recap: Self-normalized Martingale concentration bound.

## Lemma (Self-Normalized Bound for Vector-Valued Martingales)

(Abassi et. al '11) Suppose  $\{\varepsilon_i\}_{i=1}^{\infty}$  are mean zero random variables (can be generalized to martingales), and  $\varepsilon_i$  is bounded by  $\sigma$ . Let  $\{X_i\}_{i=1}^{\infty}$  be a stochastic process. Define  $\Sigma_t = \Sigma_0 + \sum_{i=1}^t X_i X_i^T$ . With probability at least  $1 - \delta$ , we have for all  $t \geq 1$ :

$$\left\| \sum_{i=1}^t X_i \varepsilon_i \right\|_{\Sigma_t^{-1}}^2 \leq \sigma^2 \log \left( \frac{\det(\Sigma_t) \det(\Sigma_0)^{-1}}{\delta^2} \right).$$

# Apply the above concentration

- How?

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

$\downarrow \downarrow$   
 $\sum_{s_{t-1}} P(\cdot | s_{t-1})$

- This is a martingale difference sequence.

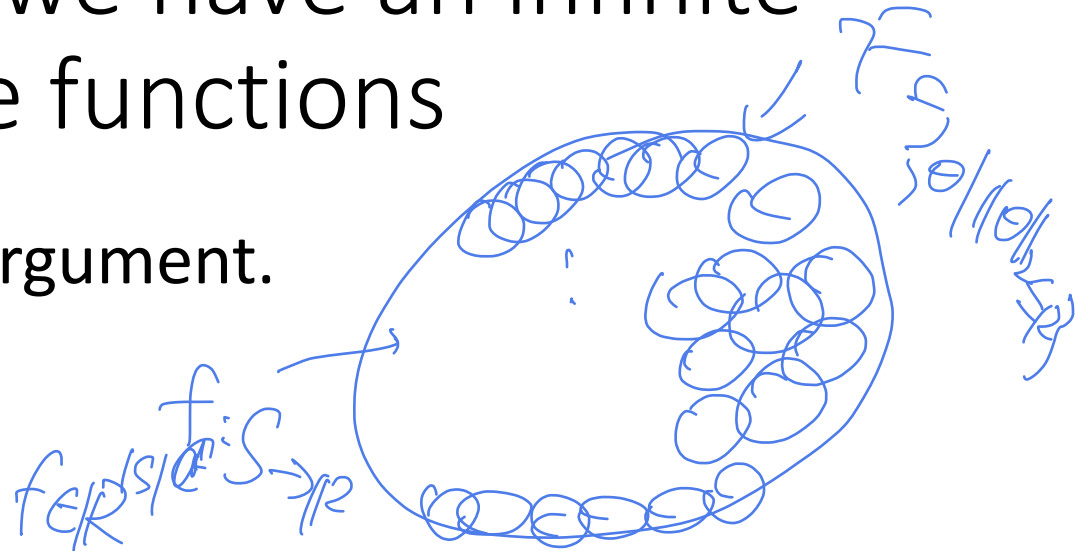
- Thus by the “Self-Normalized bound”:

$$\left\| \sum_{i=0}^{n-1} \phi(s_h^i, a_h^i) (V^\top \epsilon_h^i) \right\|_{(\Lambda_h^n)^{-1}} \leq 3H \sqrt{\ln \frac{H \det(\Lambda_h^n)^{1/2} \det(\lambda I)^{-1/2}}{\delta}}.$$

$\swarrow$  Information Path

Challenge: we cannot use union bound because we have an infinite number of value functions

- A covering number argument.



- Covering number: the number of balls with radius  $\epsilon$  that is needed to cover all points in a set.

$$\begin{aligned}
 & (\hat{P} - P)^T \hat{V} \\
 &= \underbrace{(\hat{P} - P)^T (\hat{V} - \tilde{V})}_{\leq \epsilon} + \underbrace{(\hat{P} - P)^T \tilde{V}}_{\text{pointwise union bound}}
 \end{aligned}$$

$$\forall f \in F.$$

$$\exists \tilde{f} \in \tilde{F} \text{ such that } \boxed{d(f, \tilde{f}) \leq \epsilon}$$

$$\tilde{F} \subset F, |\tilde{F}| \leq N$$

# Family of value functions we consider

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

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

What is a finite set to cover this class such that for every  $f$  in this set, there is a function in the finite set, such that they are  $\varepsilon$ -close in sup-norm?



# Covering number calculations

# From covering number to a uniform convergence bound