

CSE 542: Statistical Reinforcement Learning

Lecture 13: Linear MDPs and Online RL

Kevin Jamieson

Paul G. Allen School of Computer Science & Engineering
University of Washington

Outline

- 1 Recap of Lecture 12
- 2 Online RL: Why More Structure Is Needed
- 3 Linear MDPs
- 4 LSVI-UCB
- 5 Self-Normalized Concentration
- 6 Optimism & Regret Decomposition
- 7 The Elliptical Potential Lemma
- 8 Regret Bound for LSVI-UCB
- 9 Discussion & Outlook

Where We Left Off

Linear function approximation: feature map $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$, $\|\phi\|_2 \leq 1$, class $\mathcal{F} = \{\langle \phi, w \rangle : w \in \mathbb{R}^d\}$.

Two model-free assumptions:

(R) $Q^* \in \mathcal{F}$ (linear realizability).

(C) $\mathcal{T}_h \mathcal{F} \subseteq \mathcal{F}$ (Bellman completeness): the algorithmic assumption.

What we proved:

- Generative model + (C) $\Rightarrow \tilde{O}(d^2 H^5 / \varepsilon^2)$ samples via LSVI with G -optimal design.
- (R) alone $\Rightarrow \Omega(2^{\min(d, H)})$ lower bound.
- Offline + (C) \Rightarrow vanilla LSVI under all-policy concentrability. Pessimism (PEVI) would buy single-policy concentrability but already needs more than (C) — the same obstruction we're about to face online.

Where We Left Off

Linear function approximation: feature map $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$, $\|\phi\|_2 \leq 1$, class $\mathcal{F} = \{\langle \phi, w \rangle : w \in \mathbb{R}^d\}$.

Two model-free assumptions:

(R) $Q^* \in \mathcal{F}$ (linear realizability).

(C) $\mathcal{T}_h \mathcal{F} \subseteq \mathcal{F}$ (Bellman completeness): the algorithmic assumption.

What we proved:

- Generative model + (C) $\Rightarrow \tilde{O}(d^2 H^5 / \varepsilon^2)$ samples via LSVI with G -optimal design.
- (R) alone $\Rightarrow \Omega(2^{\min(d, H)})$ lower bound.
- Offline + (C) \Rightarrow vanilla LSVI under all-policy concentrability. Pessimism (PEVI) would buy single-policy concentrability but already needs more than (C) — the same obstruction we're about to face online.

Today: online RL

No generative model. Must explore via on-policy trajectories. We'll see that (C) is *not* enough — need a stronger **model-based** assumption.

Today's Plan

- 1 Recap of Lecture 12
- 2 Online RL: Why More Structure Is Needed**
- 3 Linear MDPs
- 4 LSVI-UCB
- 5 Self-Normalized Concentration
- 6 Optimism & Regret Decomposition
- 7 The Elliptical Potential Lemma
- 8 Regret Bound for LSVI-UCB
- 9 Discussion & Outlook

What Changes Without a Generative Model

Generative model

- Pick *any* (s, a) to query.
- G -optimal design covers feature space.
- Estimation only, no exploration.

Online RL

- Roll out a policy from $s_1 \sim \nu$.
- Visited (s, a) depend on past behavior.
- Must **explore** actively.

What Changes Without a Generative Model

Generative model

- Pick *any* (s, a) to query.
- G -optimal design covers feature space.
- Estimation only, no exploration.

Online RL

- Roll out a policy from $s_1 \sim \nu$.
- Visited (s, a) depend on past behavior.
- Must **explore** actively.

Tabular UCB-VI (Lecture 9): added a count-based bonus $H\sqrt{\log(\cdot)/n_h^k(s, a)}$ to the empirical Bellman backup.

What Changes Without a Generative Model

Generative model

- Pick *any* (s, a) to query.
- G -optimal design covers feature space.
- Estimation only, no exploration.

Online RL

- Roll out a policy from $s_1 \sim \nu$.
- Visited (s, a) depend on past behavior.
- Must **explore** actively.

Tabular UCB-VI (Lecture 9): added a count-based bonus $H\sqrt{\log(\cdot)/n_h^k(s, a)}$ to the empirical Bellman backup.

Natural attempt for linear FA: compute \hat{w}_h^k by ridge regression on past trajectory data, then use

$$\hat{Q}_h^k(s, a) = \langle \phi(s, a), \hat{w}_h^k \rangle + \beta \|\phi(s, a)\|_{(\Lambda_h^k)^{-1}}.$$

What Changes Without a Generative Model

Generative model

- Pick *any* (s, a) to query.
- G -optimal design covers feature space.
- Estimation only, no exploration.

Online RL

- Roll out a policy from $s_1 \sim \nu$.
- Visited (s, a) depend on past behavior.
- Must **explore** actively.

Tabular UCB-VI (Lecture 9): added a count-based bonus $H\sqrt{\log(\cdot)/n_h^k(s, a)}$ to the empirical Bellman backup.

Natural attempt for linear FA: compute \hat{w}_h^k by ridge regression on past trajectory data, then use

$$\hat{Q}_h^k(s, a) = \langle \phi(s, a), \hat{w}_h^k \rangle + \beta \|\phi(s, a)\|_{(\Lambda_h^k)^{-1}}.$$

Question: is Bellman completeness enough to make this work?

What Optimism Demands

For the optimism argument from Lecture 9 to go through, we need: for any value function V that the algorithm might compute,

$$|((P_h - \hat{P}_h^k)V)(s, a)| \lesssim \beta \cdot \|\phi(s, a)\|_{(\Lambda_h^k)^{-1}}.$$

This is the linear analog of the tabular concentration $|((P_h - \hat{P}_h^k)V)| \leq H\sqrt{\log/n}$.

What Optimism Demands

For the optimism argument from Lecture 9 to go through, we need: for any value function V that the algorithm might compute,

$$|((P_h - \hat{P}_h^k)V)(s, a)| \lesssim \beta \cdot \|\phi(s, a)\|_{(\Lambda_h^k)^{-1}}.$$

This is the linear analog of the tabular concentration $|((P_h - \hat{P}_h^k)V)| \leq H\sqrt{\log/n}$.

Crucial requirement. The function $(s, a) \mapsto (P_h V)(s, a)$ must lie in $\text{span}(\phi)$ for the value functions V that LSVI-UCB actually produces.

What Optimism Demands

For the optimism argument from Lecture 9 to go through, we need: for any value function V that the algorithm might compute,

$$|((P_h - \widehat{P}_h^k)V)(s, a)| \lesssim \beta \cdot \|\phi(s, a)\|_{(\Lambda_h^k)^{-1}}.$$

This is the linear analog of the tabular concentration $|((P_h - \widehat{P}_h^k)V)| \leq H\sqrt{\log/n}$.

Crucial requirement. The function $(s, a) \mapsto (P_h V)(s, a)$ must lie in $\text{span}(\phi)$ for the value functions V that LSVI-UCB actually produces.

What does LSVI-UCB produce?

$$\widehat{V}_h^k(s) = \max_a \left[\langle \phi(s, a), \widehat{w}_h^k \rangle + \beta \|\phi(s, a)\|_{(\Lambda_h^k)^{-1}} \right].$$

The max over a and the bonus mean \widehat{V}_h^k is **not** an element of \mathcal{F} .

Bellman Completeness Doesn't Suffice for Online Optimism

What (C) gives

Closure of \mathcal{F} under \mathcal{T}_h : $f \in \mathcal{F} \implies \mathcal{T}_h f \in \mathcal{F}$.

Bellman Completeness Doesn't Suffice for Online Optimism

What (C) gives

Closure of \mathcal{F} under \mathcal{T}_h : $f \in \mathcal{F} \implies \mathcal{T}_h f \in \mathcal{F}$.

Problem. $\widehat{V}_h^k \notin \mathcal{F}$ (it has a max and a bonus term). So (C) tells us *nothing* about $(P_h \widehat{V}_h^k)(s, a)$.

Bellman Completeness Doesn't Suffice for Online Optimism

What (C) gives

Closure of \mathcal{F} under \mathcal{T}_h : $f \in \mathcal{F} \implies \mathcal{T}_h f \in \mathcal{F}$.

Problem. $\widehat{V}_h^k \notin \mathcal{F}$ (it has a max and a bonus term). So (C) tells us *nothing* about $(P_h \widehat{V}_h^k)(s, a)$.

What we actually need

For every bounded $V : \mathcal{S} \rightarrow [0, H]$ — not just $V \in \mathcal{F}$ — the function $(s, a) \mapsto (P_h V)(s, a)$ should be linear in $\phi(s, a)$.

Bellman completeness only controls the linear class; we need control over a **much larger** class of V 's.

Bellman Completeness Doesn't Suffice for Online Optimism

What (C) gives

Closure of \mathcal{F} under \mathcal{T}_h : $f \in \mathcal{F} \implies \mathcal{T}_h f \in \mathcal{F}$.

Problem. $\widehat{V}_h^k \notin \mathcal{F}$ (it has a max and a bonus term). So (C) tells us *nothing* about $(P_h \widehat{V}_h^k)(s, a)$.

What we actually need

For every bounded $V : \mathcal{S} \rightarrow [0, H]$ — not just $V \in \mathcal{F}$ — the function $(s, a) \mapsto (P_h V)(s, a)$ should be linear in $\phi(s, a)$.

Bellman completeness only controls the linear class; we need control over a **much larger** class of V 's.

This is exactly what Linear MDPs guarantee.

Today's Plan

- 1 Recap of Lecture 12
- 2 Online RL: Why More Structure Is Needed
- 3 Linear MDPs**
- 4 LSVI-UCB
- 5 Self-Normalized Concentration
- 6 Optimism & Regret Decomposition
- 7 The Elliptical Potential Lemma
- 8 Regret Bound for LSVI-UCB
- 9 Discussion & Outlook

Definition: Linear MDP (Jin, Yang, Wang, Jordan 2020)

Linear MDP assumption (L)

An MDP is *linear* with feature map $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$ (with $\|\phi(s, a)\|_2 \leq 1$) if for each $h \in [H]$ there exist

- an unknown signed measure $\mu_h = (\mu_h^{(1)}, \dots, \mu_h^{(d)})$ on \mathcal{S} ,
- an unknown vector $\theta_h \in \mathbb{R}^d$,

such that

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

Normalization: $\|\mu_h(\mathcal{S})\|_2 \leq \sqrt{d}$ and $\|\theta_h\|_2 \leq \sqrt{d}$.

Definition: Linear MDP (Jin, Yang, Wang, Jordan 2020)

Linear MDP assumption (L)

An MDP is *linear* with feature map $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$ (with $\|\phi(s, a)\|_2 \leq 1$) if for each $h \in [H]$ there exist

- an unknown signed measure $\mu_h = (\mu_h^{(1)}, \dots, \mu_h^{(d)})$ on \mathcal{S} ,
- an unknown vector $\theta_h \in \mathbb{R}^d$,

such that

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

Normalization: $\|\mu_h(\mathcal{S})\|_2 \leq \sqrt{d}$ and $\|\theta_h\|_2 \leq \sqrt{d}$.

This is a model-based assumption: we commit to a specific form of P_h and r_h .

Why \sqrt{d} Is the “Natural” Scaling

The bounds $\|\mu_h(\mathcal{S})\|_2 \leq \sqrt{d}$ and $\|\theta_h\|_2 \leq \sqrt{d}$ may look arbitrary. They are calibrated so that the **tabular case sits exactly at the boundary**.

Tabular case. $\phi(s, a) = e_{(s,a)} \in \mathbb{R}^d$ with $d = SA$.

- Signed measures: $\mu_h^{(s,a)}(s') = P_h(s'|s, a) \in [0, 1]$, so $\mu_h^{(s,a)}(\mathcal{S}) = 1$ for each component.

$$\|\mu_h(\mathcal{S})\|_2 = \sqrt{\sum_{(s,a)} 1^2} = \sqrt{SA} = \sqrt{d}.$$

- Reward weights: $\theta_h^{(s,a)} = r_h(s, a) \in [0, 1]$, so

$$\|\theta_h\|_2 \leq \sqrt{SA} = \sqrt{d}.$$

Why \sqrt{d} Is the “Natural” Scaling

The bounds $\|\mu_h(\mathcal{S})\|_2 \leq \sqrt{d}$ and $\|\theta_h\|_2 \leq \sqrt{d}$ may look arbitrary. They are calibrated so that the **tabular case sits exactly at the boundary**.

Tabular case. $\phi(s, a) = e_{(s,a)} \in \mathbb{R}^d$ with $d = SA$.

- Signed measures: $\mu_h^{(s,a)}(s') = P_h(s'|s, a) \in [0, 1]$, so $\mu_h^{(s,a)}(\mathcal{S}) = 1$ for each component.

$$\|\mu_h(\mathcal{S})\|_2 = \sqrt{\sum_{(s,a)} 1^2} = \sqrt{SA} = \sqrt{d}.$$

- Reward weights: $\theta_h^{(s,a)} = r_h(s, a) \in [0, 1]$, so

$$\|\theta_h\|_2 \leq \sqrt{SA} = \sqrt{d}.$$

Interpretation

Each “feature direction” $j \in [d]$ contributes $O(1)$ to the relevant operator norms; with d such directions the ℓ_2 aggregate scales as \sqrt{d} . This is what allows the bounds $\|w_h^\pi\|_2 \leq 2H\sqrt{d}$ that drive all subsequent regret analyses.

What (L) Buys: Universal Linearity of $P_h V$

Lemma (Universal linearity)

Under (L), for every bounded function $V : \mathcal{S} \rightarrow \mathbb{R}$,

$$(P_h V)(s, a) = \left\langle \phi(s, a), \int V(s') \mu_h(s') ds' \right\rangle.$$

That is, P_h maps any V into $\text{span}(\phi)$ — not just functions in \mathcal{F} .

What (L) Buys: Universal Linearity of $P_h V$

Lemma (Universal linearity)

Under (L), for every bounded function $V : \mathcal{S} \rightarrow \mathbb{R}$,

$$(P_h V)(s, a) = \left\langle \phi(s, a), \int V(s') \mu_h(s') ds' \right\rangle.$$

That is, P_h maps any V into $\text{span}(\phi)$ — not just functions in \mathcal{F} .

Proof. $(P_h V)(s, a) = \int V(s') \langle \phi(s, a), \mu_h(s') \rangle ds' = \langle \phi(s, a), \int V \mu_h ds' \rangle.$ □

What (L) Buys: Universal Linearity of $P_h V$

Lemma (Universal linearity)

Under (L), for every bounded function $V : \mathcal{S} \rightarrow \mathbb{R}$,

$$(P_h V)(s, a) = \left\langle \phi(s, a), \int V(s') \mu_h(s') ds' \right\rangle.$$

That is, P_h maps any V into $\text{span}(\phi)$ — not just functions in \mathcal{F} .

Proof. $(P_h V)(s, a) = \int V(s') \langle \phi(s, a), \mu_h(s') \rangle ds' = \langle \phi(s, a), \int V \mu_h ds' \rangle$. □

Exactly what online optimism needed

For any \widehat{V}_h^k produced by LSVI-UCB (with bonus, max, clipping),

$$(P_h \widehat{V}_h^k)(s, a) = \langle \phi(s, a), w \rangle \text{ for some } w \in \mathbb{R}^d.$$

So $((P_h - \widehat{P}_h^k) \widehat{V}_h^k)$ becomes a self-normalized linear-regression error, controlled at rate $\|\phi\|_{(\Lambda_h^k)^{-1}}$.

Consequence: Q^π Is Linear for Every Policy

Lemma (Linearity of Q^π)

In a linear MDP, for every policy π and every $h \in [H]$, there exists $w_h^\pi \in \mathbb{R}^d$ with

$$Q_h^\pi(s, a) = \langle \phi(s, a), w_h^\pi \rangle \quad \forall (s, a).$$

Moreover $\|w_h^\pi\|_2 \leq 2H\sqrt{d}$.

Consequence: Q^π Is Linear for Every Policy

Lemma (Linearity of Q^π)

In a linear MDP, for every policy π and every $h \in [H]$, there exists $w_h^\pi \in \mathbb{R}^d$ with

$$Q_h^\pi(s, a) = \langle \phi(s, a), w_h^\pi \rangle \quad \forall (s, a).$$

Moreover $\|w_h^\pi\|_2 \leq 2H\sqrt{d}$.

Proof. Backward induction. $h = H + 1$: $Q_{H+1}^\pi \equiv 0$. \checkmark

Step. Apply universal linearity to V_{h+1}^π :

$$Q_h^\pi(s, a) = \langle \phi(s, a), \theta_h \rangle + (P_h V_{h+1}^\pi)(s, a) = \langle \phi(s, a), \theta_h + \int V_{h+1}^\pi d\mu_h \rangle.$$

So $w_h^\pi = \theta_h + \int V_{h+1}^\pi d\mu_h$, with $\|w_h^\pi\| \leq \sqrt{d} + H\sqrt{d} \leq 2H\sqrt{d}$. \square

Consequence: Q^π Is Linear for Every Policy

Lemma (Linearity of Q^π)

In a linear MDP, for every policy π and every $h \in [H]$, there exists $w_h^\pi \in \mathbb{R}^d$ with

$$Q_h^\pi(s, a) = \langle \phi(s, a), w_h^\pi \rangle \quad \forall (s, a).$$

Moreover $\|w_h^\pi\|_2 \leq 2H\sqrt{d}$.

Proof. Backward induction. $h = H + 1$: $Q_{H+1}^\pi \equiv 0$. \checkmark

Step. Apply universal linearity to V_{h+1}^π :

$$Q_h^\pi(s, a) = \langle \phi(s, a), \theta_h \rangle + (P_h V_{h+1}^\pi)(s, a) = \langle \phi(s, a), \theta_h + \int V_{h+1}^\pi d\mu_h \rangle.$$

So $w_h^\pi = \theta_h + \int V_{h+1}^\pi d\mu_h$, with $\|w_h^\pi\| \leq \sqrt{d} + H\sqrt{d} \leq 2H\sqrt{d}$. \square

Hierarchy of linear assumptions

(R) linear Q^* \subsetneq (C) Bellman completeness \subsetneq (L) Linear MDP.

(L) gives linearity of Q^π for *all* π — the strongest of the three.

Today's Plan

- 1 Recap of Lecture 12
- 2 Online RL: Why More Structure Is Needed
- 3 Linear MDPs
- 4 LSVI-UCB**
- 5 Self-Normalized Concentration
- 6 Optimism & Regret Decomposition
- 7 The Elliptical Potential Lemma
- 8 Regret Bound for LSVI-UCB
- 9 Discussion & Outlook

From UCB-VI to LSVI-UCB

Previous section: linear FA with a **generative model** — queries chosen via **G-optimal** design.

Now: no generative model; we see only on-policy trajectories and must **explore** to discover relevant (s, a) pairs.

Recall UCB-VI in the tabular case:

$$\widehat{Q}_h^k(s, a) := r_h(s, a) + (\widehat{P}_h^k \widehat{V}_{h+1}^k)(s, a) + \underbrace{H \sqrt{\frac{\log(\cdot)}{n_h^k(s, a)}}}_{\text{count-based bonus}}.$$

Linear MDP analogue.

- Replace the empirical Bellman backup with a **ridge regression** of $r_h + \widehat{V}_{h+1}^k$ onto ϕ .
- Replace counts $n_h^k(s, a)$ by the **elliptical norm** $\|\phi(s, a)\|_{(\Lambda_h^k)^{-1}}^2$.

From UCB-VI to LSVI-UCB

Previous section: linear FA with a **generative model** — queries chosen via **G-optimal** design.

Now: no generative model; we see only on-policy trajectories and must **explore** to discover relevant (s, a) pairs.

Recall UCB-VI in the tabular case:

$$\widehat{Q}_h^k(s, a) := r_h(s, a) + (\widehat{P}_h^k \widehat{V}_{h+1}^k)(s, a) + \underbrace{H \sqrt{\frac{\log(\cdot)}{n_h^k(s, a)}}}_{\text{count-based bonus}}.$$

Linear MDP analogue.

- Replace the empirical Bellman backup with a **ridge regression** of $r_h + \widehat{V}_{h+1}^k$ onto ϕ .
- Replace counts $n_h^k(s, a)$ by the **elliptical norm** $\|\phi(s, a)\|_{(\Lambda_h^k)^{-1}}$.

Slogan

$$\text{Counts} \longleftrightarrow \text{Elliptical norm} \quad \sqrt{\frac{1}{n_h^k(s, a)}} \longleftrightarrow \|\phi(s, a)\|_{(\Lambda_h^k)^{-1}}$$

Ridge Regression Recap

Given data $\{(\phi_i, y_i)\}_{i=1}^n$ with $\phi_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$, ridge regression with regularizer $\lambda > 0$:

$$\hat{w} = \arg \min_{w \in \mathbb{R}^d} \sum_{i=1}^n (\langle \phi_i, w \rangle - y_i)^2 + \lambda \|w\|_2^2.$$

Closed form:

$$\hat{w} = \Lambda^{-1} \sum_{i=1}^n \phi_i y_i, \quad \Lambda := \lambda I + \sum_{i=1}^n \phi_i \phi_i^\top.$$

Ridge Regression Recap

Given data $\{(\phi_i, y_i)\}_{i=1}^n$ with $\phi_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$, ridge regression with regularizer $\lambda > 0$:

$$\hat{w} = \arg \min_{w \in \mathbb{R}^d} \sum_{i=1}^n (\langle \phi_i, w \rangle - y_i)^2 + \lambda \|w\|_2^2.$$

Closed form:

$$\hat{w} = \Lambda^{-1} \sum_{i=1}^n \phi_i y_i, \quad \Lambda := \lambda I + \sum_{i=1}^n \phi_i \phi_i^\top.$$

Elliptical norm

$$\|x\|_\Lambda^2 := x^\top \Lambda x, \quad \|x\|_{\Lambda^{-1}}^2 := x^\top \Lambda^{-1} x.$$

Interpretation. $\|\phi(s, a)\|_{\Lambda^{-1}}^2$ measures how “new” the direction $\phi(s, a)$ is, given past features:

- Tabular case ($\phi(s, a) = e_{(s,a)}$): $\Lambda = \text{diag}(\lambda + n(s, a))$, so $\|\phi(s, a)\|_{\Lambda^{-1}}^2 = 1/(\lambda + n(s, a))$.
- Hence $\|\phi(s, a)\|_{\Lambda^{-1}}^2$ is the linear-algebraic analog of the inverse count.

Algorithm

Input: feature map ϕ , regularizer $\lambda > 0$, bonus scale $\beta > 0$.

For episode $k = 1, 2, \dots, K$:

- 1 Set $\widehat{V}_{H+1}^k(\cdot) \equiv 0$. For $h = H, H-1, \dots, 1$:

$$\Lambda_h^k = \lambda I + \sum_{\tau=1}^{k-1} \phi(\mathbf{s}_h^\tau, \mathbf{a}_h^\tau) \phi(\mathbf{s}_h^\tau, \mathbf{a}_h^\tau)^\top,$$

$$\widehat{\mathbf{w}}_h^k = (\Lambda_h^k)^{-1} \sum_{\tau=1}^{k-1} \phi(\mathbf{s}_h^\tau, \mathbf{a}_h^\tau) [r_h^\tau + \widehat{V}_{h+1}^k(\mathbf{s}_{h+1}^\tau)],$$

$$b_h^k(\mathbf{s}, \mathbf{a}) = \beta \|\phi(\mathbf{s}, \mathbf{a})\|_{(\Lambda_h^k)^{-1}},$$

$$\widehat{Q}_h^k(\mathbf{s}, \mathbf{a}) = \min\{H, \langle \phi(\mathbf{s}, \mathbf{a}), \widehat{\mathbf{w}}_h^k \rangle + b_h^k(\mathbf{s}, \mathbf{a})\},$$

$$\widehat{V}_h^k(\mathbf{s}) = \max_a \widehat{Q}_h^k(\mathbf{s}, \mathbf{a}), \quad \pi_h^k(\mathbf{s}) = \arg \max_a \widehat{Q}_h^k(\mathbf{s}, \mathbf{a}).$$

- 2 Execute π^k , collect trajectory $(\mathbf{s}_1^k, \mathbf{a}_1^k, r_1^k, \dots, \mathbf{s}_H^k, \mathbf{a}_H^k, r_H^k)$.

Why the Closed Form Looks Right

If we knew V_{h+1}^* exactly, the linearity lemma says

$$r_h(s, a) + \mathbb{E}_{s'} V_{h+1}^*(s') = \langle \phi(s, a), w_h^* \rangle,$$

with $w_h^* = \theta_h + \int V_{h+1}^* d\mu_h$.

Why the Closed Form Looks Right

If we knew V_{h+1}^* exactly, the linearity lemma says

$$r_h(s, a) + \mathbb{E}_{s'} V_{h+1}^*(s') = \langle \phi(s, a), w_h^* \rangle,$$

with $w_h^* = \theta_h + \int V_{h+1}^* d\mu_h$.

With \widehat{V}_{h+1}^k as a stand-in for V_{h+1}^* , the regression target

$$y_h^\tau = r_h^\tau + \widehat{V}_{h+1}^k(s_{h+1}^\tau)$$

satisfies

$$\mathbb{E}[y_h^\tau \mid s_h^\tau, a_h^\tau] = \langle \phi(s_h^\tau, a_h^\tau), w_h^k \rangle, \quad w_h^k \triangleq \theta_h + \int \widehat{V}_{h+1}^k d\mu_h.$$

Ridge regression $\Rightarrow \widehat{w}_h^k$ approximates w_h^k .

Why the Closed Form Looks Right

If we knew V_{h+1}^* exactly, the linearity lemma says

$$r_h(s, a) + \mathbb{E}_{s'} V_{h+1}^*(s') = \langle \phi(s, a), w_h^* \rangle,$$

with $w_h^* = \theta_h + \int V_{h+1}^* d\mu_h$.

With \widehat{V}_{h+1}^k as a stand-in for V_{h+1}^* , the regression target

$$y_h^\tau = r_h^\tau + \widehat{V}_{h+1}^k(s_{h+1}^\tau)$$

satisfies

$$\mathbb{E}[y_h^\tau \mid s_h^\tau, a_h^\tau] = \langle \phi(s_h^\tau, a_h^\tau), w_h^k \rangle, \quad w_h^k \triangleq \theta_h + \int \widehat{V}_{h+1}^k d\mu_h.$$

Ridge regression $\Rightarrow \widehat{w}_h^k$ approximates w_h^k .

The key issue

\widehat{V}_{h+1}^k is itself a function of past data, so the regression targets are *not* i.i.d. around a fixed w_h^k .

We need a **uniform** concentration bound over a function class containing all candidate \widehat{V}_{h+1}^k .

Today's Plan

- 1 Recap of Lecture 12
- 2 Online RL: Why More Structure Is Needed
- 3 Linear MDPs
- 4 LSVI-UCB
- 5 Self-Normalized Concentration**
- 6 Optimism & Regret Decomposition
- 7 The Elliptical Potential Lemma
- 8 Regret Bound for LSVI-UCB
- 9 Discussion & Outlook

Self-Normalized Concentration: Warm-up

Setup: features $\phi_1, \dots, \phi_n \in \mathbb{R}^d$ (possibly chosen adaptively) and noises η_1, \dots, η_n with $\mathbb{E}[\eta_i \mid \mathcal{F}_{i-1}] = 0$ and $|\eta_i| \leq R$.

Let $\Lambda_n = \lambda I + \sum_{i \leq n} \phi_i \phi_i^\top$ and $S_n = \sum_{i \leq n} \phi_i \eta_i$.

Self-Normalized Concentration: Warm-up

Setup: features $\phi_1, \dots, \phi_n \in \mathbb{R}^d$ (possibly chosen adaptively) and noises η_1, \dots, η_n with $\mathbb{E}[\eta_i \mid \mathcal{F}_{i-1}] = 0$ and $|\eta_i| \leq R$.

Let $\Lambda_n = \lambda I + \sum_{i \leq n} \phi_i \phi_i^\top$ and $S_n = \sum_{i \leq n} \phi_i \eta_i$.

Theorem (Abbasi-Yadkori, Pál, Szepesvári 2011)

With probability at least $1 - \delta$, for all $n \geq 1$ simultaneously:

$$\|S_n\|_{\Lambda_n^{-1}} \leq R \sqrt{2 \log \left(\frac{\det(\Lambda_n)^{1/2}}{\delta \det(\lambda I)^{1/2}} \right)} \leq R \sqrt{d \log(1 + n/(\lambda d)) + 2 \log(1/\delta)}.$$

Self-Normalized Concentration: Warm-up

Setup: features $\phi_1, \dots, \phi_n \in \mathbb{R}^d$ (possibly chosen adaptively) and noises η_1, \dots, η_n with $\mathbb{E}[\eta_i \mid \mathcal{F}_{i-1}] = 0$ and $|\eta_i| \leq R$.

Let $\Lambda_n = \lambda I + \sum_{i \leq n} \phi_i \phi_i^\top$ and $S_n = \sum_{i \leq n} \phi_i \eta_i$.

Theorem (Abbasi-Yadkori, Pál, Szepesvári 2011)

With probability at least $1 - \delta$, for all $n \geq 1$ simultaneously:

$$\|S_n\|_{\Lambda_n^{-1}} \leq R \sqrt{2 \log \left(\frac{\det(\Lambda_n)^{1/2}}{\delta \det(\lambda I)^{1/2}} \right)} \leq R \sqrt{d \log(1 + n/(\lambda d)) + 2 \log(1/\delta)}.$$

Interpretation. This is the linear analogue of Hoeffding's bound: $\frac{1}{\sqrt{n}}$ becomes $\|\phi\|_{\Lambda_n^{-1}}$, and $\log(1/\delta)$ becomes $\log(\det \Lambda_n / \det(\lambda I)) \leq d \log(1 + n/(\lambda d))$.

Confidence Set for w_h^k

Lemma (Confidence ellipsoid)

Suppose the regression target satisfies $y_h^\tau = \langle \phi(s_h^\tau, a_h^\tau), w_h^k \rangle + \eta_h^\tau$ with $\mathbb{E}[\eta_h^\tau \mid \mathcal{F}_{\tau-1}] = 0$ and $|\eta_h^\tau| \leq H$. Then with probability $\geq 1 - \delta$, for all k, h :

$$\|\widehat{w}_h^k - w_h^k\|_{\Lambda_h^k} \leq \beta_*, \quad \beta_* := H\sqrt{d \log(1 + K/(\lambda d))} + \sqrt{\lambda} \|w_h^k\|_2.$$

Confidence Set for w_h^k

Lemma (Confidence ellipsoid)

Suppose the regression target satisfies $y_h^\tau = \langle \phi(s_h^\tau, a_h^\tau), w_h^k \rangle + \eta_h^\tau$ with $\mathbb{E}[\eta_h^\tau \mid \mathcal{F}_{\tau-1}] = 0$ and $|\eta_h^\tau| \leq H$. Then with probability $\geq 1 - \delta$, for all k, h :

$$\|\widehat{w}_h^k - w_h^k\|_{\Lambda_h^k} \leq \beta_*, \quad \beta_* := H\sqrt{d \log(1 + K/(\lambda d))} + \sqrt{\lambda} \|w_h^k\|_2.$$

Consequence. By Cauchy–Schwarz:

$$|\langle \phi(s, a), \widehat{w}_h^k - w_h^k \rangle| \leq \|\widehat{w}_h^k - w_h^k\|_{\Lambda_h^k} \cdot \|\phi(s, a)\|_{(\Lambda_h^k)^{-1}} \leq \beta_* \|\phi(s, a)\|_{(\Lambda_h^k)^{-1}}.$$

Confidence Set for w_h^k

Lemma (Confidence ellipsoid)

Suppose the regression target satisfies $y_h^\tau = \langle \phi(s_h^\tau, a_h^\tau), w_h^k \rangle + \eta_h^\tau$ with $\mathbb{E}[\eta_h^\tau \mid \mathcal{F}_{\tau-1}] = 0$ and $|\eta_h^\tau| \leq H$. Then with probability $\geq 1 - \delta$, for all k, h :

$$\|\widehat{w}_h^k - w_h^k\|_{\Lambda_h^k} \leq \beta_*, \quad \beta_* := H\sqrt{d \log(1 + K/(\lambda d))} + \sqrt{\lambda} \|w_h^k\|_2.$$

Consequence. By Cauchy–Schwarz:

$$|\langle \phi(s, a), \widehat{w}_h^k - w_h^k \rangle| \leq \|\widehat{w}_h^k - w_h^k\|_{\Lambda_h^k} \cdot \|\phi(s, a)\|_{(\Lambda_h^k)^{-1}} \leq \beta_* \|\phi(s, a)\|_{(\Lambda_h^k)^{-1}}.$$

The bonus has the right form

Choosing $\beta = \widetilde{O}(H\sqrt{d})$, we get $|\langle \phi(s, a), \widehat{w}_h^k - w_h^k \rangle| \leq b_h^k(s, a)$ on the good event.

The Problem: Adaptive Targets

The confidence-ellipsoid lemma assumed a *fixed* target w_h^k . But $w_h^k = \theta_h + \int \widehat{V}_{h+1}^k d\mu_h$ depends on \widehat{V}_{h+1}^k , which depends on the data through $\widehat{w}_{h+1}^k, \widehat{w}_{h+2}^k, \dots$

The Problem: Adaptive Targets

The confidence-ellipsoid lemma assumed a *fixed* target w_h^k . But $w_h^k = \theta_h + \int \widehat{V}_{h+1}^k d\mu_h$ depends on \widehat{V}_{h+1}^k , which depends on the data through $\widehat{w}_{h+1}^k, \widehat{w}_{h+2}^k, \dots$

Naive worry. If we apply the self-normalized bound *conditional* on \widehat{V}_{h+1}^k , the conditioning destroys the martingale structure — \widehat{V}_{h+1}^k was built from the very samples η_h^τ that enter the noise.

The Problem: Adaptive Targets

The confidence-ellipsoid lemma assumed a *fixed* target w_h^k . But $w_h^k = \theta_h + \int \widehat{V}_{h+1}^k d\mu_h$ depends on \widehat{V}_{h+1}^k , which depends on the data through $\widehat{w}_{h+1}^k, \widehat{w}_{h+2}^k, \dots$

Naive worry. If we apply the self-normalized bound *conditional* on \widehat{V}_{h+1}^k , the conditioning destroys the martingale structure — \widehat{V}_{h+1}^k was built from the very samples η_h^T that enter the noise.

Fix: union-bound over a cover of all possible V 's. Define

$$\mathcal{V} := \left\{ V(s) = \min \left\{ H, \max_a \left[\langle \phi(s, a), w \rangle + \beta \|\phi(s, a)\|_{\Lambda^{-1}} \right] \right\} : \|w\| \leq L, \Lambda \succeq \lambda I \right\}.$$

By construction, $\widehat{V}_{h+1}^k \in \mathcal{V}$ for every k .

Covering number bound (Jin et al. 2020)

\mathcal{V} admits an ε -cover \mathcal{N}_ε (in $\|\cdot\|_\infty$) of size $\log |\mathcal{N}_\varepsilon| \lesssim d \log(1 + L/\varepsilon) + d^2 \log(1 + \beta^2/(\lambda\varepsilon^2))$.

Cover Argument, in Detail

Step 1: pointwise bound for fixed V . For each *fixed* $V \in \mathcal{V}$, let $w(V) := \theta_h + \int V d\mu_h$. Targets $r_h^\tau + V(s_{h+1}^\tau)$ are i.i.d. around $\langle \phi_h^\tau, w(V) \rangle$; self-normalized lemma at confidence δ' :

$$\|\widehat{w}_h^k(V) - w(V)\|_{\Lambda_h^k} \leq H \sqrt{2 \log(1/\delta') + d \log(1 + K/(\lambda d))}.$$

Cover Argument, in Detail

Step 1: pointwise bound for fixed V . For each *fixed* $V \in \mathcal{V}$, let $w(V) := \theta_h + \int V d\mu_h$. Targets $r_h^\tau + V(s_{h+1}^\tau)$ are i.i.d. around $\langle \phi_h^\tau, w(V) \rangle$; self-normalized lemma at confidence δ' :

$$\|\widehat{w}_h^k(V) - w(V)\|_{\Lambda_h^k} \leq H \sqrt{2 \log(1/\delta')} + d \log(1 + K/(\lambda d)).$$

Step 2: union bound. Apply Step 1 to all $V \in \mathcal{N}_\varepsilon$ at $\delta' = \delta/|\mathcal{N}_\varepsilon|$:

$$\forall V \in \mathcal{N}_\varepsilon : \|\widehat{w}_h^k(V) - w(V)\|_{\Lambda_h^k} \leq H \sqrt{2 \log(|\mathcal{N}_\varepsilon|/\delta)} + d \log(1 + K/(\lambda d)).$$

Cover Argument, in Detail

Step 1: pointwise bound for fixed V . For each *fixed* $V \in \mathcal{V}$, let $w(V) := \theta_h + \int V d\mu_h$. Targets $r_h^\tau + V(s_{h+1}^\tau)$ are i.i.d. around $\langle \phi_h^\tau, w(V) \rangle$; self-normalized lemma at confidence δ' :

$$\|\widehat{w}_h^k(V) - w(V)\|_{\Lambda_h^k} \leq H\sqrt{2\log(1/\delta')} + d\log(1 + K/(\lambda d)).$$

Step 2: union bound. Apply Step 1 to all $V \in \mathcal{N}_\varepsilon$ at $\delta' = \delta/|\mathcal{N}_\varepsilon|$:

$$\forall V \in \mathcal{N}_\varepsilon : \|\widehat{w}_h^k(V) - w(V)\|_{\Lambda_h^k} \leq H\sqrt{2\log(|\mathcal{N}_\varepsilon|/\delta)} + d\log(1 + K/(\lambda d)).$$

Step 3: discretization. For $\widehat{V}_{h+1}^k \in \mathcal{V}$, pick $\tilde{V} \in \mathcal{N}_\varepsilon$ with $\|\widehat{V}_{h+1}^k - \tilde{V}\|_\infty \leq \varepsilon$. Triangle inequality gives $\|\widehat{w}_h^k(\widehat{V}_{h+1}^k) - w(\widehat{V}_{h+1}^k)\|_{\Lambda_h^k} \leq \|\widehat{w}_h^k(\tilde{V}) - w(\tilde{V})\|_{\Lambda_h^k} + O(\varepsilon\sqrt{Kd})$.

Cover Argument, in Detail

Step 1: pointwise bound for fixed V . For each fixed $V \in \mathcal{V}$, let $w(V) := \theta_h + \int V d\mu_h$. Targets $r_h^\tau + V(s_{h+1}^\tau)$ are i.i.d. around $\langle \phi_h^\tau, w(V) \rangle$; self-normalized lemma at confidence δ' :

$$\|\widehat{w}_h^k(V) - w(V)\|_{\Lambda_h^k} \leq H\sqrt{2\log(1/\delta')} + d\log(1 + K/(\lambda d)).$$

Step 2: union bound. Apply Step 1 to all $V \in \mathcal{N}_\varepsilon$ at $\delta' = \delta/|\mathcal{N}_\varepsilon|$:

$$\forall V \in \mathcal{N}_\varepsilon : \|\widehat{w}_h^k(V) - w(V)\|_{\Lambda_h^k} \leq H\sqrt{2\log(|\mathcal{N}_\varepsilon|/\delta)} + d\log(1 + K/(\lambda d)).$$

Step 3: discretization. For $\widehat{V}_{h+1}^k \in \mathcal{V}$, pick $\tilde{V} \in \mathcal{N}_\varepsilon$ with $\|\widehat{V}_{h+1}^k - \tilde{V}\|_\infty \leq \varepsilon$. Triangle inequality gives $\|\widehat{w}_h^k(\widehat{V}_{h+1}^k) - w(\widehat{V}_{h+1}^k)\|_{\Lambda_h^k} \leq \|\widehat{w}_h^k(\tilde{V}) - w(\tilde{V})\|_{\Lambda_h^k} + O(\varepsilon\sqrt{Kd})$.

Choose $\varepsilon = 1/K$. Then $\log|\mathcal{N}_\varepsilon| = \tilde{O}(d^2)$, discretization = $\tilde{O}(1)$. Net: replace $\log(1/\delta)$ in the self-normalized bound by $\tilde{O}(d^2 + \log(1/\delta))$.

Cover Argument, in Detail

Step 1: pointwise bound for fixed V . For each fixed $V \in \mathcal{V}$, let $w(V) := \theta_h + \int V d\mu_h$. Targets $r_h^\tau + V(s_{h+1}^\tau)$ are i.i.d. around $\langle \phi_h^\tau, w(V) \rangle$; self-normalized lemma at confidence δ' :

$$\|\widehat{w}_h^k(V) - w(V)\|_{\Lambda_h^k} \leq H\sqrt{2\log(1/\delta')} + d\log(1 + K/(\lambda d)).$$

Step 2: union bound. Apply Step 1 to all $V \in \mathcal{N}_\varepsilon$ at $\delta' = \delta/|\mathcal{N}_\varepsilon|$:

$$\forall V \in \mathcal{N}_\varepsilon : \|\widehat{w}_h^k(V) - w(V)\|_{\Lambda_h^k} \leq H\sqrt{2\log(|\mathcal{N}_\varepsilon|/\delta)} + d\log(1 + K/(\lambda d)).$$

Step 3: discretization. For $\widehat{V}_{h+1}^k \in \mathcal{V}$, pick $\tilde{V} \in \mathcal{N}_\varepsilon$ with $\|\widehat{V}_{h+1}^k - \tilde{V}\|_\infty \leq \varepsilon$. Triangle inequality gives $\|\widehat{w}_h^k(\widehat{V}_{h+1}^k) - w(\widehat{V}_{h+1}^k)\|_{\Lambda_h^k} \leq \|\widehat{w}_h^k(\tilde{V}) - w(\tilde{V})\|_{\Lambda_h^k} + O(\varepsilon\sqrt{Kd})$.

Choose $\varepsilon = 1/K$. Then $\log|\mathcal{N}_\varepsilon| = \tilde{O}(d^2)$, discretization = $\tilde{O}(1)$. Net: replace $\log(1/\delta)$ in the self-normalized bound by $\tilde{O}(d^2 + \log(1/\delta))$.

End-to-end consequence

Choosing $\beta = cHd\sqrt{\log(dKH/\delta)}$, on the good event \mathcal{E} :

$$|\langle \phi(s, a), \widehat{w}_h^k - w_h^k \rangle| \leq \beta \|\phi(s, a)\|_{(\Lambda_h^k)^{-1}} = b_h^k(s, a) \text{ holds simultaneously for all } k, h, s, a.$$

Today's Plan

- 1 Recap of Lecture 12
- 2 Online RL: Why More Structure Is Needed
- 3 Linear MDPs
- 4 LSVI-UCB
- 5 Self-Normalized Concentration
- 6 Optimism & Regret Decomposition**
- 7 The Elliptical Potential Lemma
- 8 Regret Bound for LSVI-UCB
- 9 Discussion & Outlook

Lemma: Optimism

Lemma (Optimism)

On the good event \mathcal{E} (the uniform self-normalized bound holds), for all k, h, s, a :

$$\widehat{Q}_h^k(s, a) \geq Q_h^*(s, a), \quad \widehat{V}_h^k(s) \geq V_h^*(s).$$

Lemma: Optimism

Lemma (Optimism)

On the good event \mathcal{E} (the uniform self-normalized bound holds), for all k, h, s, a :

$$\widehat{Q}_h^k(s, a) \geq Q_h^*(s, a), \quad \widehat{V}_h^k(s) \geq V_h^*(s).$$

Proof (backward induction on h):

Base ($h = H + 1$): both sides are zero. \checkmark

Step. Suppose $\widehat{V}_{h+1}^k \geq V_{h+1}^*$. Recall $w_h^k = \theta_h + \int \widehat{V}_{h+1}^k d\mu_h$, so

$$\langle \phi(s, a), w_h^k \rangle = r_h(s, a) + (P_h \widehat{V}_{h+1}^k)(s, a) \geq r_h(s, a) + (P_h V_{h+1}^*)(s, a) = Q_h^*(s, a).$$

On \mathcal{E} , $|\langle \phi(s, a), \widehat{w}_h^k - w_h^k \rangle| \leq b_h^k(s, a)$. Therefore

$$\widehat{Q}_h^k(s, a) \geq \langle \phi(s, a), \widehat{w}_h^k \rangle + b_h^k(s, a) \geq \langle \phi(s, a), w_h^k \rangle \geq Q_h^*(s, a).$$

(If the min with H is active, $\widehat{Q}_h^k = H \geq Q_h^*$ trivially.) Take max over a .

Regret-to-Surplus Certificate (Linear Version)

Lemma (Certificate)

Define the surplus

$$\varepsilon_h^k(s, a) := \widehat{Q}_h^k(s, a) - r_h(s, a) - (P_h \widehat{V}_{h+1}^k)(s, a).$$

On the good event, for every episode k :

$$V_0^* - V_0^{\pi^k} \leq \sum_{h=1}^H \mathbb{E}[\varepsilon_h^k(s_h^k, a_h^k) \mid \pi^k].$$

Regret-to-Surplus Certificate (Linear Version)

Lemma (Certificate)

Define the surplus

$$\varepsilon_h^k(s, a) := \widehat{Q}_h^k(s, a) - r_h(s, a) - (P_h \widehat{V}_{h+1}^k)(s, a).$$

On the good event, for every episode k :

$$V_0^* - V_0^{\pi^k} \leq \sum_{h=1}^H \mathbb{E}[\varepsilon_h^k(s_h^k, a_h^k) \mid \pi^k].$$

Proof. Identical to the tabular UCB-VI case (Lecture 9): define $\Gamma_h^k := \widehat{V}_h^k - V_h^{\pi^k}$ and unroll

$$\Gamma_h^k(s_h^k) = \varepsilon_h^k(s_h^k, a_h^k) + (P_h \Gamma_{h+1}^k)(s_h^k, a_h^k).$$

Optimism gives $V_0^* \leq \widehat{V}_0^k$, so $V_0^* - V_0^{\pi^k} \leq \mathbb{E}[\Gamma_1^k(s_1^k)] = \sum_h \mathbb{E}[\varepsilon_h^k]$. □

Regret-to-Surplus Certificate (Linear Version)

Lemma (Certificate)

Define the surplus

$$\varepsilon_h^k(s, a) := \widehat{Q}_h^k(s, a) - r_h(s, a) - (P_h \widehat{V}_{h+1}^k)(s, a).$$

On the good event, for every episode k :

$$V_0^* - V_0^{\pi^k} \leq \sum_{h=1}^H \mathbb{E}[\varepsilon_h^k(s_h^k, a_h^k) \mid \pi^k].$$

Proof. Identical to the tabular UCB-VI case (Lecture 9): define $\Gamma_h^k := \widehat{V}_h^k - V_h^{\pi^k}$ and unroll

$$\Gamma_h^k(s_h^k) = \varepsilon_h^k(s_h^k, a_h^k) + (P_h \Gamma_{h+1}^k)(s_h^k, a_h^k).$$

Optimism gives $V_0^* \leq \widehat{V}_0^k$, so $V_0^* - V_0^{\pi^k} \leq \mathbb{E}[\Gamma_1^k(s_1^k)] = \sum_h \mathbb{E}[\varepsilon_h^k]$. □

Same structure as tabular UCB-VI. All the work goes into bounding the surplus.

Bounding the Surplus

Recall the good event:

$$\mathcal{E} := \{ \forall k, h, s, a : |\langle \phi(s, a), \widehat{w}_h^k - w_h^k \rangle| \leq \beta \|\phi(s, a)\|_{(\Lambda_h^k)^{-1}} \}, \mathbb{P}(\mathcal{E}) \geq 1 - \delta.$$

Unpack the surplus. When $\min\{H, \cdot\}$ is inactive, $\widehat{Q}_h^k(s, a) = \langle \phi(s, a), \widehat{w}_h^k \rangle + b_h^k(s, a)$:

$$\begin{aligned} \varepsilon_h^k(s, a) &= \widehat{Q}_h^k(s, a) - r_h(s, a) - (P_h \widehat{V}_{h+1}^k)(s, a) \\ &= \langle \phi(s, a), \widehat{w}_h^k \rangle + b_h^k(s, a) - \underbrace{[r_h(s, a) + (P_h \widehat{V}_{h+1}^k)(s, a)]}_{= \langle \phi(s, a), w_h^k \rangle \text{ by (L)}} \\ &= \underbrace{\langle \phi(s, a), \widehat{w}_h^k - w_h^k \rangle}_{\text{regression error}} + \underbrace{b_h^k(s, a)}_{\text{bonus}}. \end{aligned}$$

(The reward r_h is part of $w_h^k = \theta_h + \int \widehat{V}_{h+1}^k d\mu_h$ and cancels.)

Bounding the Surplus

Recall the good event:

$$\mathcal{E} := \{ \forall k, h, s, a : | \langle \phi(s, a), \widehat{w}_h^k - w_h^k \rangle | \leq \beta \| \phi(s, a) \|_{(\Lambda_h^k)^{-1}} \}, \mathbb{P}(\mathcal{E}) \geq 1 - \delta.$$

Unpack the surplus. When $\min\{H, \cdot\}$ is inactive, $\widehat{Q}_h^k(s, a) = \langle \phi(s, a), \widehat{w}_h^k \rangle + b_h^k(s, a)$:

$$\begin{aligned} \varepsilon_h^k(s, a) &= \widehat{Q}_h^k(s, a) - r_h(s, a) - (P_h \widehat{V}_{h+1}^k)(s, a) \\ &= \langle \phi(s, a), \widehat{w}_h^k \rangle + b_h^k(s, a) - \underbrace{[r_h(s, a) + (P_h \widehat{V}_{h+1}^k)(s, a)]}_{= \langle \phi(s, a), w_h^k \rangle \text{ by (L)}} \\ &= \underbrace{\langle \phi(s, a), \widehat{w}_h^k - w_h^k \rangle}_{\text{regression error}} + \underbrace{b_h^k(s, a)}_{\text{bonus}}. \end{aligned}$$

(The reward r_h is part of $w_h^k = \theta_h + \int \widehat{V}_{h+1}^k d\mu_h$ and cancels.)

Apply \mathcal{E} . $| \langle \phi(s, a), \widehat{w}_h^k - w_h^k \rangle | \leq b_h^k(s, a)$, so $0 \leq \varepsilon_h^k(s, a) \leq 2 b_h^k(s, a) = 2\beta \| \phi(s, a) \|_{(\Lambda_h^k)^{-1}}$.

Today's Plan

- 1 Recap of Lecture 12
- 2 Online RL: Why More Structure Is Needed
- 3 Linear MDPs
- 4 LSVI-UCB
- 5 Self-Normalized Concentration
- 6 Optimism & Regret Decomposition
- 7 The Elliptical Potential Lemma**
- 8 Regret Bound for LSVI-UCB
- 9 Discussion & Outlook

Putting It Together

Step 1: regret to sum of surpluses. By the certificate (Lemma), on \mathcal{E} ,

$$R(K) = \sum_{k=1}^K (V_0^* - V_0^{\pi^k}) \leq \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\varepsilon_h^k(\mathbf{s}_h^k, \mathbf{a}_h^k) \mid \pi^k \right].$$

Step 2: surplus to elliptical norm. Apply the surplus bound $\varepsilon_h^k \leq 2\beta \|\phi\|_{(\Lambda_h^k)^{-1}}$:

$$R(K) \leq 2\beta \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left\| \phi(\mathbf{s}_h^k, \mathbf{a}_h^k) \right\|_{(\Lambda_h^k)^{-1}} \right].$$

Putting It Together

Step 1: regret to sum of surpluses. By the certificate (Lemma), on \mathcal{E} ,

$$R(K) = \sum_{k=1}^K (V_0^* - V_0^{\pi^k}) \leq \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\varepsilon_h^k(s_h^k, a_h^k) \mid \pi^k \right].$$

Step 2: surplus to elliptical norm. Apply the surplus bound $\varepsilon_h^k \leq 2\beta \|\phi\|_{(\Lambda_h^k)^{-1}}$:

$$R(K) \leq 2\beta \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left\| \phi(s_h^k, a_h^k) \right\|_{(\Lambda_h^k)^{-1}} \right].$$

Step 3: bound the sum of elliptical norms. For each fixed h , we need

$$\sum_{k=1}^K \left\| \phi(s_h^k, a_h^k) \right\|_{(\Lambda_h^k)^{-1}}.$$

Putting It Together

Step 1: regret to sum of surpluses. By the certificate (Lemma), on \mathcal{E} ,

$$R(K) = \sum_{k=1}^K (V_0^* - V_0^{\pi^k}) \leq \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\varepsilon_h^k(s_h^k, a_h^k) \mid \pi^k \right].$$

Step 2: surplus to elliptical norm. Apply the surplus bound $\varepsilon_h^k \leq 2\beta \|\phi\|_{(\Lambda_h^k)^{-1}}$:

$$R(K) \leq 2\beta \sum_{k=1}^K \sum_{h=1}^H \mathbb{E} \left[\left\| \phi(s_h^k, a_h^k) \right\|_{(\Lambda_h^k)^{-1}} \right].$$

Step 3: bound the sum of elliptical norms. For each fixed h , we need

$$\sum_{k=1}^K \left\| \phi(s_h^k, a_h^k) \right\|_{(\Lambda_h^k)^{-1}}.$$

Linear analogue of the counting argument

Tabular: $\sum_k 1/\sqrt{n_h^k(s_h^k, a_h^k)} \leq 2\sqrt{SAK}$.

Linear: Elliptical Potential Lemma, replacing SA with d .

Elliptical Potential Lemma

Lemma (Elliptical Potential, Abbasi-Yadkori et al. 2011)

Let $\phi_1, \phi_2, \dots, \phi_K \in \mathbb{R}^d$ with $\|\phi_k\| \leq 1$, and $\Lambda_k = \lambda I + \sum_{\tau < k} \phi_\tau \phi_\tau^\top$ with $\lambda \geq 1$. Then

$$\sum_{k=1}^K \min\left\{1, \|\phi_k\|_{\Lambda_k^{-1}}^2\right\} \leq 2 \log \frac{\det(\Lambda_{K+1})}{\det(\lambda I)} \leq 2d \log\left(1 + \frac{K}{\lambda d}\right).$$

By Cauchy–Schwarz:

$$\sum_{k=1}^K \min\left\{1, \|\phi_k\|_{\Lambda_k^{-1}}^2\right\} \leq \sqrt{2Kd \log(1 + K/(\lambda d))} = \tilde{O}(\sqrt{dK}).$$

Elliptical Potential Lemma

Lemma (Elliptical Potential, Abbasi-Yadkori et al. 2011)

Let $\phi_1, \phi_2, \dots, \phi_K \in \mathbb{R}^d$ with $\|\phi_k\| \leq 1$, and $\Lambda_k = \lambda I + \sum_{\tau < k} \phi_\tau \phi_\tau^\top$ with $\lambda \geq 1$. Then

$$\sum_{k=1}^K \min\left\{1, \|\phi_k\|_{\Lambda_k}^2\right\} \leq 2 \log \frac{\det(\Lambda_{K+1})}{\det(\lambda I)} \leq 2d \log\left(1 + \frac{K}{\lambda d}\right).$$

By Cauchy–Schwarz:

$$\sum_{k=1}^K \min\left\{1, \|\phi_k\|_{\Lambda_k}^2\right\} \leq \sqrt{2Kd \log(1 + K/(\lambda d))} = \tilde{O}(\sqrt{dK}).$$

Intuition. \mathbb{R}^d contains only $\tilde{O}(d)$ “effectively new” directions to explore. Each new direction increases $\log \det \Lambda$, but $\log \det$ is bounded.

Tabular check. $\phi_k = e_{(s_k, a_k)}$, $d = SA$: gives $\sum_k (n_k(s_k, a_k))^{-1/2} \leq \tilde{O}(\sqrt{SAK})$. ✓

Proof of the Elliptical Potential Lemma

Let $u_k = \min\{1, \|\phi_k\|_{\Lambda_k^{-1}}^2\}$. Use $u \leq 2 \log(1 + u)$ for $u \in [0, 1]$:

$$\sum_k u_k \leq 2 \sum_k \log(1 + \|\phi_k\|_{\Lambda_k^{-1}}^2).$$

Proof of the Elliptical Potential Lemma

Let $u_k = \min\{1, \|\phi_k\|_{\Lambda_k^{-1}}^2\}$. Use $u \leq 2 \log(1 + u)$ for $u \in [0, 1]$:

$$\sum_k u_k \leq 2 \sum_k \log(1 + \|\phi_k\|_{\Lambda_k^{-1}}^2).$$

Matrix determinant lemma:

$$\det(\Lambda_{k+1}) = \det(\Lambda_k + \phi_k \phi_k^\top) = \det(\Lambda_k) (1 + \|\phi_k\|_{\Lambda_k^{-1}}^2).$$

Proof of the Elliptical Potential Lemma

Let $u_k = \min\{1, \|\phi_k\|_{\Lambda_k^{-1}}^2\}$. Use $u \leq 2 \log(1 + u)$ for $u \in [0, 1]$:

$$\sum_k u_k \leq 2 \sum_k \log(1 + \|\phi_k\|_{\Lambda_k^{-1}}^2).$$

Matrix determinant lemma:

$$\det(\Lambda_{k+1}) = \det(\Lambda_k + \phi_k \phi_k^\top) = \det(\Lambda_k) (1 + \|\phi_k\|_{\Lambda_k^{-1}}^2).$$

Telescoping:

$$\sum_{k=1}^K \log(1 + \|\phi_k\|_{\Lambda_k^{-1}}^2) = \log \frac{\det(\Lambda_{K+1})}{\det(\lambda I)}.$$

Proof of the Elliptical Potential Lemma

Let $u_k = \min\{1, \|\phi_k\|_{\Lambda_k^{-1}}^2\}$. Use $u \leq 2 \log(1 + u)$ for $u \in [0, 1]$:

$$\sum_k u_k \leq 2 \sum_k \log(1 + \|\phi_k\|_{\Lambda_k^{-1}}^2).$$

Matrix determinant lemma:

$$\det(\Lambda_{k+1}) = \det(\Lambda_k + \phi_k \phi_k^\top) = \det(\Lambda_k) (1 + \|\phi_k\|_{\Lambda_k^{-1}}^2).$$

Telescoping:

$$\sum_{k=1}^K \log(1 + \|\phi_k\|_{\Lambda_k^{-1}}^2) = \log \frac{\det(\Lambda_{K+1})}{\det(\lambda I)}.$$

Bound on $\det(\Lambda_{K+1})$. Eigenvalues of Λ_{K+1} sum to

$$\text{tr}(\Lambda_{K+1}) = \lambda d + \sum_k \|\phi_k\|^2 \leq \lambda d + K.$$

By AM-GM:

$$\det(\Lambda_{K+1}) \leq \left(\frac{\text{tr}(\Lambda_{K+1})}{d} \right)^d \leq (\lambda + K/d)^d.$$

Therefore $\log \det(\Lambda_{K+1}) / \det(\lambda I) \leq d \log(1 + K/(\lambda d))$. \square

Today's Plan

- 1 Recap of Lecture 12
- 2 Online RL: Why More Structure Is Needed
- 3 Linear MDPs
- 4 LSVI-UCB
- 5 Self-Normalized Concentration
- 6 Optimism & Regret Decomposition
- 7 The Elliptical Potential Lemma
- 8 Regret Bound for LSVI-UCB**
- 9 Discussion & Outlook

Main Theorem

Theorem (Jin, Yang, Wang, Jordan 2020)

Set $\lambda = 1$ and $\beta = c H d \sqrt{\log(dKH/\delta)}$ for an appropriate constant c . With probability at least $1 - \delta$, the regret of LSVI-UCB satisfies

$$R(K) = \sum_{k=1}^K (V_0^* - V_0^{\pi^k}) \leq \tilde{O}(\sqrt{d^3 H^4 K}).$$

The bound has no dependence on $|\mathcal{S}|$ or $|\mathcal{A}|$.

Main Theorem

Theorem (Jin, Yang, Wang, Jordan 2020)

Set $\lambda = 1$ and $\beta = c Hd \sqrt{\log(dKH/\delta)}$ for an appropriate constant c . With probability at least $1 - \delta$, the regret of LSVI-UCB satisfies

$$R(K) = \sum_{k=1}^K (V_0^* - V_0^{\pi^k}) \leq \tilde{O}(\sqrt{d^3 H^4 K}).$$

The bound has no dependence on $|\mathcal{S}|$ or $|\mathcal{A}|$.

Proof. From the previous slide, $R(K) \leq 2\beta \sum_{h,k} \mathbb{E}[\|\phi\|_{(\Lambda_h^k)} - 1]$.

Why $\min\{1, \cdot\}$? Values are in $[0, H]$, so $\varepsilon_h^k(s, a) \leq H$. Hence $\varepsilon_h^k \leq \min\{H, 2\beta \|\phi\|_{(\Lambda_h^k)}\} = 2\beta \min\{H/(2\beta), \|\phi\|_{(\Lambda_h^k)}\}$. Since $\beta \geq H/2$: $\varepsilon_h^k \leq 2\beta \min\{1, \|\phi\|_{(\Lambda_h^k)}\}$.

Elliptical potential (with min inside): $\sum_{k=1}^K \min\{1, \|\phi_h^k\|_{(\Lambda_h^k)}\} \leq \sqrt{2Kd \log(1 + K/d)}$.

Main Theorem

Theorem (Jin, Yang, Wang, Jordan 2020)

Set $\lambda = 1$ and $\beta = c Hd \sqrt{\log(dKH/\delta)}$ for an appropriate constant c . With probability at least $1 - \delta$, the regret of LSVI-UCB satisfies

$$R(K) = \sum_{k=1}^K (V_0^* - V_0^{\pi^k}) \leq \tilde{O}(\sqrt{d^3 H^4 K}).$$

The bound has no dependence on $|\mathcal{S}|$ or $|\mathcal{A}|$.

Proof. From the previous slide, $R(K) \leq 2\beta \sum_{h,k} \mathbb{E}[\|\phi\|_{(\Lambda_h^k)^{-1}}]$.

Why $\min\{1, \cdot\}$? Values are in $[0, H]$, so $\varepsilon_h^k(s, a) \leq H$. Hence $\varepsilon_h^k \leq \min\{H, 2\beta \|\phi\|_{(\Lambda_h^k)^{-1}}\} = 2\beta \min\{H/(2\beta), \|\phi\|_{(\Lambda_h^k)^{-1}}\}$. Since $\beta \geq H/2$: $\varepsilon_h^k \leq 2\beta \min\{1, \|\phi\|_{(\Lambda_h^k)^{-1}}\}$.

Elliptical potential (with min inside): $\sum_{k=1}^K \min\{1, \|\phi_h^k\|_{(\Lambda_h^k)^{-1}}\} \leq \sqrt{2Kd \log(1 + K/d)}$.

Assemble. $R(K) \leq 2\beta \cdot H \cdot \sqrt{2Kd \log(1 + K/d)} = \tilde{O}(H \cdot Hd \cdot \sqrt{Kd}) = \tilde{O}(\sqrt{d^3 H^4 K})$. \square

Comparing the Bounds

Setting	Algorithm	Regret
Tabular MDP	UCB-VI (basic)	$\tilde{O}(H^2 S \sqrt{SAK})$
Tabular MDP	UCB-VI (Bernstein)	$\tilde{O}(H^{3/2} \sqrt{SAK})$
Linear MDP	LSVI-UCB	$\tilde{O}(\sqrt{d^3 H^4 K})$
Linear MDP	lower bound	$\Omega(d \sqrt{H^3 K})$

Comparing the Bounds

Setting	Algorithm	Regret
Tabular MDP	UCB-VI (basic)	$\tilde{O}(H^2 S \sqrt{SAK})$
Tabular MDP	UCB-VI (Bernstein)	$\tilde{O}(H^{3/2} \sqrt{SAK})$
Linear MDP	LSVI-UCB	$\tilde{O}(\sqrt{d^3 H^4 K})$
Linear MDP	lower bound	$\Omega(d \sqrt{H^3 K})$

Take-aways.

- Replacing S, A with d is a **huge win** when features are well-chosen.
- Polynomial gap to the lower bound: closing it has been an active research direction.
- Computation: each Bellman backup requires only inverting a $d \times d$ matrix and computing one ridge-regression solution per stage.

Where the $d^3 H^4$ Comes From

Decomposition of $d^3 H^4 K$:

Source	Contribution
Bonus scale $\beta = \tilde{O}(H\sqrt{d})$	$H^2 d$
Elliptical potential (\sqrt{dK})	dK
Sum over $h \in [H]$ stages	H^2
Covering number for \mathcal{V}	extra $\tilde{O}(d)$ in β
Product	$\tilde{O}(\sqrt{d^3 H^4 K})$

Where Can We Tighten?

- **d -dependence ($d^3 \rightarrow d^2$):** improved uniform concentration / randomized LSVI.
Zanette, Lazaric, Kochenderfer, Brunskill, *Frequentist regret bounds for randomized least-squares value iteration*, AISTATS 2020.
- **H -dependence ($H^4 \rightarrow H^3$):** variance-aware (Bernstein) recursion, analogous to Lecture 10.
Hu, Yang, Wang, *Nearly minimax optimal RL with linear function approximation*, ICML 2022;
Agarwal, Jin, Zhang, *VOQL: Towards optimal regret in model-free RL with nonlinear function approximation*, COLT 2023.
- **Closing the gap to $\Omega(d\sqrt{H^3K})$:**
He, Zhou, Gu, *Nearly minimax optimal RL for linear Markov decision processes*, ICML 2023.
- **Lower bound:**
Zhou, Gu, Szepesvári, *Nearly minimax optimal RL for linear mixture Markov decision processes*, COLT 2021.

Today's Plan

- 1 Recap of Lecture 12
- 2 Online RL: Why More Structure Is Needed
- 3 Linear MDPs
- 4 LSVI-UCB
- 5 Self-Normalized Concentration
- 6 Optimism & Regret Decomposition
- 7 The Elliptical Potential Lemma
- 8 Regret Bound for LSVI-UCB
- 9 Discussion & Outlook

Linear Realizability is Not Enough

Weaker assumption: only $Q_h^*(s, a) = \langle \phi(s, a), w_h^* \rangle$ for some $w_h^* \in \mathbb{R}^d$.

Linear Realizability is Not Enough

Weaker assumption: only $Q_h^*(s, a) = \langle \phi(s, a), w_h^* \rangle$ for some $w_h^* \in \mathbb{R}^d$.

Theorem (Du, Kakade, Lee, Lovett, Mahajan, Sun, Wang 2020)

There exist MDPs with linear Q^* (in d -dimensional features) for which any algorithm requires $\Omega(\min(2^d, 2^H))$ samples to find an ε -optimal policy.

Linear Realizability is Not Enough

Weaker assumption: only $Q_h^*(s, a) = \langle \phi(s, a), w_h^* \rangle$ for some $w_h^* \in \mathbb{R}^d$.

Theorem (Du, Kakade, Lee, Lovett, Mahajan, Sun, Wang 2020)

There exist MDPs with linear Q^* (in d -dimensional features) for which any algorithm requires $\Omega(\min(2^d, 2^H))$ samples to find an ε -optimal policy.

Why is the linear MDP assumption stronger?

Linear MDP $\Rightarrow Q^\pi$ is linear for every policy π .

Linear $Q^* \Rightarrow$ only Q^* (and not other Q^π) is linear.

Linear Realizability is Not Enough

Weaker assumption: only $Q_h^*(s, a) = \langle \phi(s, a), w_h^* \rangle$ for some $w_h^* \in \mathbb{R}^d$.

Theorem (Du, Kakade, Lee, Lovett, Mahajan, Sun, Wang 2020)

There exist MDPs with linear Q^* (in d -dimensional features) for which any algorithm requires $\Omega(\min(2^d, 2^H))$ samples to find an ε -optimal policy.

Why is the linear MDP assumption stronger?

Linear MDP $\Rightarrow Q^\pi$ is linear for every policy π .

Linear $Q^* \Rightarrow$ only Q^* (and not other Q^π) is linear.

Moral

Sample-efficient RL with function approximation requires *closure under Bellman operators*, not just realizability of the optimal value function.

What we did today

- Three nested assumptions: **(R)** $Q^* \in \mathcal{F}$, **(C)** $\mathcal{T}_h \mathcal{F} \subseteq \mathcal{F}$, **(L)** linear MDP.
- **Generative model + (C)**: LSVI with G -optimal design gives $\tilde{O}(d^2 H^4 / \epsilon^2)$ samples.
- **(R) alone is not enough**: $2^{\min(d, H)}$ lower bound.
- **Online + (L)**: LSVI-UCB achieves $R(K) = \tilde{O}(\sqrt{d^3 H^4 K})$, no $|\mathcal{S}|, |\mathcal{A}|$ dependence.
- Key tools: self-normalized concentration, elliptical potential lemma.

What we did today

- Three nested assumptions: **(R)** $Q^* \in \mathcal{F}$, **(C)** $\mathcal{T}_h \mathcal{F} \subseteq \mathcal{F}$, **(L)** linear MDP.
- **Generative model + (C)**: LSVI with G -optimal design gives $\tilde{O}(d^2 H^4 / \epsilon^2)$ samples.
- **(R) alone is not enough**: $2^{\min(d, H)}$ lower bound.
- **Online + (L)**: LSVI-UCB achieves $R(K) = \tilde{O}(\sqrt{d^3 H^4 K})$, no $|\mathcal{S}|, |\mathcal{A}|$ dependence.
- Key tools: self-normalized concentration, elliptical potential lemma.

Mental dictionary:

Tabular	Linear
$n_h^k(s, a)$	$\Lambda_h^k = \lambda I + \sum_{\tau} \phi_h^\tau \phi_h^{\tau T}$
$1 / \sqrt{n_h^k(s, a)}$	$\ \phi(s, a)\ _{(\Lambda_h^k)^{-1}}$
counting bound (\sqrt{SAK})	elliptical potential (\sqrt{dK})
$S \cdot A$	d

Reading & References

Primary reference:

AJKS Agarwal, Jiang, Kakade, Sun, *Reinforcement Learning: Theory and Algorithms*, Chapter 3.

Foundational:

- Jin, Yang, Wang, Jordan, COLT 2020 (LSVI-UCB).
- Abbasi-Yadkori, Pál, Szepesvári, NeurIPS 2011 (self-normalized concentration).
- Weisz, Amortila, Szepesvári, ALT 2021 (realizability lower bound).

Improvements:

- Zanette, Lazaric, Kochenderfer, Brunskill, AISTATS 2020.
- Hu, Yang, Wang, ICML 2022.
- He, Zhou, Gu, ICML 2023.

Questions?