

RL with Linear Features: When Does It Work & When Doesn't It Work?

Part 2: Linear BC & The Offline RL Case
CS 2284: Foundations of Reinforcement Learning

Kianté Brantley & Sham Kakade

$\int \gamma_c !$

Agenda

Announcements

- HW1 due Monday.

Recap

- Regression + D-opt design.

Today

- Finish linear BC analysis.
- Offline RL

Recap

Least-Squares Value Iteration (LSVI)

Setting: Finite Horizon H , Generative Model (Simulator).

Algorithm: Backward Induction via Regression.

① Initialize $\hat{V}_H(s) = 0$.

② For $h = H - 1, \dots, 0$:

- **Collect Data:** Generate dataset $D_h = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^N$.

- **Form Targets:** Compute regression targets using the *next* value function:

$$y_i = r_h(s_i, a_i) + \hat{V}_{h+1}(s'_i)$$

- **Regression:** Solve for parameters $\hat{\theta}_h$:

$$\hat{\theta}_h \in \operatorname{argmin}_{\theta \in \mathbb{R}^d} \sum_{i=1}^N (\theta^\top \phi(s_i, a_i) - y_i)^2$$

- **Update:** Set $\hat{Q}_h(s, a) = \hat{\theta}_h^\top \phi(s, a)$ and $\hat{V}_h(s) = \max_a \hat{Q}_h(s, a)$.

$$E[y_i | s, a] = \mathcal{R}[\theta_{h+1}]$$

$$= \operatorname{argmax}_a \hat{\theta}_{h+1}^\top \phi(s, a')$$

The Assumption Ladder

We can organize linear RL assumptions from weakest (hardest) to strongest (easiest).

(A) Agnostic Approximation

No realizability. Q^* is “close” to linear.

Status: Hard. Requires strong distribution assumptions.

(B) Linear Q^* Realizability

$$Q_h^*(s, a) = (\theta_h^*)^\top \phi(s, a).$$

Status: Insufficient. Exponential lower bounds exist.

(C) All-Policies Realizability

$$Q_h^\pi(s, a) = (\theta_h^\pi)^\top \phi(s, a) \text{ for all } \pi.$$

Status: Subtle. Fails offline even with perfect coverage.

(D) Linear Bellman Completeness

$$\mathcal{T}_h f \in \mathcal{F} \text{ for all } f \in \mathcal{F}.$$

Status: Sufficient! (This Lecture)

Fixed Design OLS (The Tool)

Consider the standard linear regression setting:

$$\|\vec{x}\|_{\Lambda}^2 = \vec{x}^T \Lambda \vec{x}$$

$$y_i = x_i^\top \theta^* + \xi_i, \quad \text{with } \mathbb{E}[\xi_i | x_i] = 0 \text{ (sub-Gaussian).}$$

The OLS estimator is $\hat{\theta} = \Lambda^{-1} \left(\frac{1}{N} \sum_{i=1}^N x_i y_i \right)$, where $\Lambda = \frac{1}{N} \sum_i x_i x_i^\top$.

Fixed-design OLS Bound

With probability at least $1 - \delta$, the prediction error is bounded in the Λ -norm:

$$\|\hat{\theta} - \theta^*\|_{\Lambda} \stackrel{z}{\lesssim} \sigma \sqrt{\frac{d \log(1/\delta)}{N}}.$$

This translates to a *pointwise* bound using leverage scores:

$$|(\hat{\theta} - \theta^*)^\top \phi(s, a)| \leq \|\hat{\theta} - \theta^*\|_{\Lambda} \sqrt{\phi(s, a)^\top \Lambda^{-1} \phi(s, a)}$$

OLS review

Consider the standard linear regression setting:

$$\vec{y} = \mathbf{X} \theta^* + \vec{\epsilon}$$

$$y_i = x_i^\top \theta^* + \xi_i, \quad \text{with } \mathbb{E}[\xi_i | x_i] = 0 \text{ (sub-Gaussian).}$$

$$\xi_i \sim \mathcal{N}(0, \sigma^2) \quad \mathbf{X} \in \mathbb{R}^{n \times d}$$

$$\vec{y}, \vec{\epsilon} \in \mathbb{R}^n$$

The OLS estimator is $\hat{\theta} = \Lambda^{-1} \left(\frac{1}{N} \sum_{i=1}^N x_i y_i \right)$, where $\Lambda = \frac{1}{N} \sum_i x_i x_i^\top$.

$$\Sigma = \frac{1}{n} \mathbf{X}^\top \mathbf{X}$$

$$v = \frac{1}{n} \mathbf{X}^\top \mathbf{y}$$

$$v = \frac{1}{n} \mathbf{X}^\top (\mathbf{X} \theta^* + \vec{\epsilon}) = \Lambda \theta^* + \frac{1}{n} \mathbf{X}^\top \vec{\epsilon}$$

$$\hat{\theta} - \theta^* = \Lambda^{-1} v - \theta^* = \frac{1}{n} \Lambda^{-1} \mathbf{X}^\top \vec{\epsilon}$$

$$\sigma^2 \mathbf{I}$$

$$\begin{aligned} \mathbb{E}[\|\hat{\theta} - \theta^*\|_1^2] &= \frac{1}{n^2} \mathbb{E} \left[\vec{\epsilon}^\top \mathbf{X} \Lambda^{-1} \mathbf{X} \vec{\epsilon} \right] = \frac{1}{n^2} \text{Tr} \left(\mathbf{X}^\top \mathbb{E}[\vec{\epsilon} \vec{\epsilon}^\top] \mathbf{X} \Lambda^{-1} \right) \\ &= \frac{\sigma^2}{n} \text{Tr}(\Lambda^{-1} \Lambda) = \frac{\sigma^2}{n} \text{Tr}(\mathbf{I}) = \frac{d \sigma^2}{n} \end{aligned}$$

D-optimal design: the leverage-minimizing geometry

To guarantee uniform bounds, we must choose our training data carefully.

The feature set is:

$$\Phi := \{\phi(s, a) : (s, a) \in \mathcal{S} \times \mathcal{A}\} \subset \mathbb{R}^d.$$

D-optimal design (lemma; geometric fact)

Suppose Φ is compact. There exists a distribution ρ supported on at most $d(d+1)/2$ state-action pairs s.t. with

$$\Sigma := \mathbb{E}_{(s,a) \sim \rho} [\phi(s, a) \phi(s, a)^\top],$$

we have $\Sigma \succ 0$ and

$$\sup_{(s,a)} \phi(s, a)^\top \Sigma^{-1} \phi(s, a) = d.$$

Furthermore, no distribution ρ can achieve a lower (worst-case) leverage score.

Today: Analysis of LSVI

Roadmap: Proving LSVI Works

Recall our strategy:

Regression Error $\xrightarrow{\text{D-Opt}}$ Pointwise Error $\xrightarrow{\text{Completeness}}$ Bellman Residual $\xrightarrow{\text{Sim. Lemma}}$ Policy Loss

We will formalize this in three steps:

- 1 **Regression Analysis:** Bound the error $\|\hat{\theta}_h - \theta_h^*\|$ using OLS.
- 2 **Bellman Residuals:** D-Optimal Design lets us translate parameter err into uniform value function error Res_h .
- 3 **The Simulation Lemma:** Show how residuals propagate to the final value/policy.

Step 1: The Regression at Stage h

Consider a fixed stage h . We have the "next" value function \widehat{V}_{h+1} .

The Data:

- Dataset $D_h = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^N$ collected via D-optimal design ρ .

Step 1: The Regression at Stage h

Consider a fixed stage h . We have the "next" value function \widehat{V}_{h+1} .

The Data:

- Dataset $D_h = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^N$ collected via D-optimal design ρ .
- Design Matrix $\Lambda_h = \frac{1}{N} \sum_{i=1}^N \phi(s_i, a_i) \phi(s_i, a_i)^\top \approx \Sigma_\rho$.

Step 1: The Regression at Stage h

Consider a fixed stage h . We have the "next" value function \widehat{V}_{h+1} .

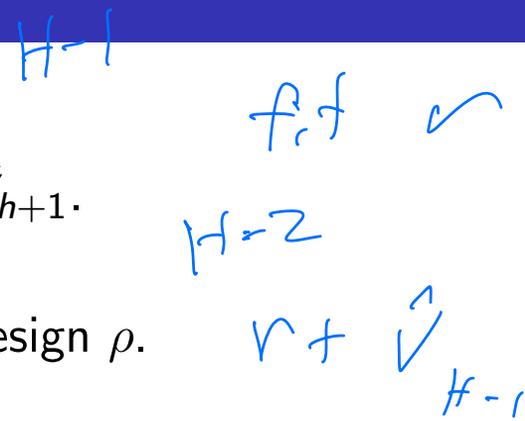
The Data:

- Dataset $D_h = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^N$ collected via D-optimal design ρ .
- Design Matrix $\Lambda_h = \frac{1}{N} \sum_{i=1}^N \phi(s_i, a_i) \phi(s_i, a_i)^\top \approx \Sigma_\rho$.

The Target:

$$y_i = r_h(s_i, a_i) + \widehat{V}_{h+1}(s'_i)$$

Step 1: The Regression at Stage h



Consider a fixed stage h . We have the "next" value function \widehat{V}_{h+1} .

The Data:

- Dataset $D_h = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^N$ collected via D-optimal design ρ .
- Design Matrix $\Lambda_h = \frac{1}{N} \sum_{i=1}^N \phi(s_i, a_i) \phi(s_i, a_i)^\top \approx \Sigma_\rho$.

The Target:

$$y_i = r_h(s_i, a_i) + \widehat{V}_{h+1}(s'_i)$$

By **Completeness**, the true expected target is linear:

$$\mathbb{E}[y_i | s_i, a_i] = (\mathcal{T}_h \widehat{Q}_{h+1})(s_i, a_i) = (\tilde{\theta}_h^\star)^\top \phi(s_i, a_i)$$

Thus, $y_i = (\tilde{\theta}_h^\star)^\top \phi(s_i, a_i) + \xi_i$, where ξ_i is sub-Gaussian noise.

$\tilde{\theta}$ is better notation
(best not to confuse $\tilde{\theta}$
with opt θ^\star)
 $\tilde{\theta}$ is from BC.
[fixed in pdf slides]
& chapter

Step 1: The Fixed Design Bound

We apply the standard Fixed-Design OLS bound.

OLS Generalization Bound

With probability $1 - \delta$, for all h

$$\|\hat{\theta}_h - \theta_h^*\|_{\Lambda_h} \leq c \cdot H \sqrt{\frac{d \log(H/\delta)}{N}}$$

(Here the noise scale is H because values are bounded by H).

$$\sup_{\xi, a} \left| \hat{\theta}_n \cdot \phi(\xi, a) - \theta_n \cdot \phi(\xi, a) \right| \leq ??$$

Step 2: From Parameters to Pointwise Error

We want to bound the prediction error at *any* (s, a) .

\sqrt{d}

$$|(\hat{\theta}_h - \theta_h^*)^\top \phi(s, a)| \leq \|\hat{\theta}_h - \theta_h^*\|_{\Lambda_h} \sqrt{\phi(s, a)^\top \Lambda_h^{-1} \phi(s, a)}$$

\top
 $H \sqrt{d}$

Step 2: From Parameters to Pointwise Error

We want to bound the prediction error at *any* (s, a) .

$$|(\hat{\theta}_h - \theta_h^*)^\top \phi(s, a)| \leq \|\hat{\theta}_h - \theta_h^*\|_{\Lambda_h} \sqrt{\phi(s, a)^\top \Lambda_h^{-1} \phi(s, a)}$$

Plug in our bounds (using $\Lambda_h \approx \Sigma_\rho$):

- Parameter Error (OLS): $\approx H \sqrt{\frac{d}{N}}$
- Leverage Score (D-Opt): $\sqrt{\phi^\top \Lambda_h^{-1} \phi} \approx \sqrt{\phi^\top \Sigma_\rho^{-1} \phi} \leq \sqrt{d}$

Step 2: From Parameters to Pointwise Error

We want to bound the prediction error at *any* (s, a) .

$$|(\hat{\theta}_h - \theta_h^*)^\top \phi(s, a)| \leq \|\hat{\theta}_h - \theta_h^*\|_{\Lambda_h} \sqrt{\phi(s, a)^\top \Lambda_h^{-1} \phi(s, a)}$$

Plug in our bounds (using $\Lambda_h \approx \Sigma_\rho$):

- Parameter Error (OLS): $\approx H\sqrt{\frac{d}{N}}$
- Leverage Score (D-Opt): $\sqrt{\phi^\top \Lambda_h^{-1} \phi} \approx \sqrt{\phi^\top \Sigma_\rho^{-1} \phi} \leq \sqrt{d}$

Result (Uniform Error): Multiply them together:

$$\sup_{(s,a)} |\hat{Q}_h(s, a) - (\mathcal{T}_h \hat{Q}_{h+1})(s, a)| \lesssim H\sqrt{\frac{d}{N}} \cdot \sqrt{d} = \frac{Hd}{\sqrt{N}}$$

Step 3: The Simulation Lemma

Definitions:

- Let $\hat{V}_h(s) := \max_a \hat{Q}_h(s, a)$ and $\hat{\pi}_h(s) := \operatorname{argmax}_a \hat{Q}_h(s, a)$.
- Define the stage- h **Bellman Residual** as the prediction error:

$$\operatorname{Res}_h := \left\| \hat{Q}_h - \mathcal{T}_h \hat{Q}_{h+1} \right\|_{\infty}$$

Bellman residual
of g w.r.t. f
 $\|g - \mathcal{T}f\|_{\infty}$

Step 3: The Simulation Lemma

Definitions:

- Let $\widehat{V}_h(s) := \max_a \widehat{Q}_h(s, a)$ and $\widehat{\pi}_h(s) := \operatorname{argmax}_a \widehat{Q}_h(s, a)$.
- Define the stage- h **Bellman Residual** as the prediction error:

$$\operatorname{Res}_h := \left\| \widehat{Q}_h - \mathcal{T}_h \widehat{Q}_{h+1} \right\|_{\infty}$$

$$\eta \approx \frac{H\alpha}{\sqrt{N}}$$

Lemma 3.3 (Small residual \Rightarrow good policy)

Assume $\widehat{Q}_H \equiv 0$ and $\operatorname{Res}_h \leq \eta$ for all h . Then:

- 1 **Value Accuracy:** For all $h \in [H]$,

$$\left\| \widehat{Q}_h - Q_h^* \right\|_{\infty} \leq (H - h)\eta.$$

Step 3: The Simulation Lemma

Definitions:

- Let $\widehat{V}_h(s) := \max_a \widehat{Q}_h(s, a)$ and $\widehat{\pi}_h(s) := \operatorname{argmax}_a \widehat{Q}_h(s, a)$.
- Define the stage- h **Bellman Residual** as the prediction error:

$$\operatorname{Res}_h := \left\| \widehat{Q}_h - \mathcal{T}_h \widehat{Q}_{h+1} \right\|_\infty$$

Lemma 3.3 (Small residual \Rightarrow good policy)

Assume $\widehat{Q}_H \equiv 0$ and $\operatorname{Res}_h \leq \eta$ for all h . Then:

- 1 **Value Accuracy:** For all $h \in [H]$,

$$\left\| \widehat{Q}_h - Q_h^* \right\|_\infty \leq (H - h)\eta.$$

- 2 **Policy Loss:** For the greedy policy $\widehat{\pi}$,

$$V_0^*(s) - V_0^{\widehat{\pi}}(s) \leq 2H^2\eta.$$

Proof of Claim 1: Value Accuracy

We prove $\|\hat{Q}_h - Q_h^*\|_\infty \leq (H - h)\eta$ by backward induction.

Base Case ($h = H$): $Q_H^* = \hat{Q}_H = 0$, so error is 0.

$$\pm \gamma_n \hat{Q}_{n+1}$$

Inductive Step: Consider the error at stage h :

$$\begin{aligned} \|\hat{Q}_h - Q_h^*\|_\infty &= \|\hat{Q}_h - \gamma_n Q_{h+1}^*\|_\infty \\ &\leq \|\hat{Q}_h - \gamma_n \hat{Q}_{h+1}\|_\infty + \|\gamma_n \hat{Q}_{h+1} - \gamma_n Q_{h+1}^*\|_\infty \end{aligned}$$

Proof of Claim 1: Value Accuracy

We prove $\|\widehat{Q}_h - Q_h^*\|_\infty \leq (H - h)\eta$ by backward induction.

Base Case ($h = H$): $Q_H^* = \widehat{Q}_H = 0$, so error is 0.

Inductive Step: Consider the error at stage h :

$$\begin{aligned}\|\widehat{Q}_h - Q_h^*\|_\infty &= \|\widehat{Q}_h - \mathcal{T}_h Q_{h+1}^*\|_\infty \\ &\leq \underbrace{\|\widehat{Q}_h - \mathcal{T}_h \widehat{Q}_{h+1}\|_\infty}_{\text{Residual } \leq \eta} + \underbrace{\|\mathcal{T}_h \widehat{Q}_{h+1} - \mathcal{T}_h Q_{h+1}^*\|_\infty}_{\text{Recursive Error}}\end{aligned}$$

Proof of Claim 1: Value Accuracy

We prove $\|\widehat{Q}_h - Q_h^*\|_\infty \leq (H - h)\eta$ by backward induction.

Base Case ($h = H$): $Q_H^* = \widehat{Q}_H = 0$, so error is 0.

Inductive Step: Consider the error at stage h :

$$\begin{aligned}\|\widehat{Q}_h - Q_h^*\|_\infty &= \|\widehat{Q}_h - \mathcal{T}_h Q_{h+1}^*\|_\infty \\ &\leq \underbrace{\|\widehat{Q}_h - \mathcal{T}_h \widehat{Q}_{h+1}\|_\infty}_{\text{Residual} \leq \eta} + \underbrace{\|\mathcal{T}_h \widehat{Q}_{h+1} - \mathcal{T}_h Q_{h+1}^*\|_\infty}_{\text{Recursive Error}}\end{aligned}$$

Using the contraction property (non-expansion) of \mathcal{T}_h :

$$\|\mathcal{T}_h \widehat{Q}_{h+1} - \mathcal{T}_h Q_{h+1}^*\|_\infty \leq \|\widehat{Q}_{h+1} - Q_{h+1}^*\|_\infty$$

Thus, $\text{Error}_h \leq \eta + \text{Error}_{h+1}$. Unrolling gives $\sum \eta = (H - h)\eta$. □

Proof of Claim 2: Policy Loss

This follows the standard "Performance Difference" logic.

The Argument:

- The loss of a greedy policy is bounded by the sum of suboptimality gaps at each step.
- Since $\hat{\pi}$ is greedy with respect to \hat{Q} , the single-step loss is bounded by $2 \times$ estimation error:

$$Q_h^*(s, \pi^*) - Q_h^*(s, \hat{\pi}) \leq 2 \|\hat{Q}_h - Q_h^*\|_\infty$$

Total Loss: Summing over H steps:

$$V_0^* - V_0^{\hat{\pi}} \leq \sum_{h=0}^{H-1} 2 \underbrace{\|\hat{Q}_h - Q_h^*\|_\infty}_{\leq (H-h)\eta} \leq \sum_{h=0}^{H-1} 2(H-h)\eta \approx H^2\eta$$

Final Result: LSVI Sample Complexity

$$\|\hat{Q}_n - Q_n^*\| \leq H^2 \mathcal{N}$$

Putting it all together:

- 1 We want final error $V^* - V^{\hat{\pi}} \leq \epsilon$.
- 2 By Simulation Lemma, we need $\text{Res}_h \leq \epsilon/(2H^2)$. $= \mathcal{N}$
- 3 By Regression Analysis, we need $Hd/\sqrt{N} \approx \epsilon/H^2$.

\ll
 \mathcal{N}

Final Result: LSVI Sample Complexity

Putting it all together:

- 1 We want final error $V^* - V^{\hat{\pi}} \leq \epsilon$.
- 2 By Simulation Lemma, we need $\text{Res}_h \leq \epsilon/(2H^2)$.
- 3 By Regression Analysis, we need $Hd/\sqrt{N} \approx \epsilon/H^2$.

Solve for N :

$$\frac{Hd}{\sqrt{N}} \approx \frac{\epsilon}{H^2} \implies \sqrt{N} \approx \frac{H^3 d}{\epsilon} \implies N \approx \frac{H^6 d^2}{\epsilon^2}$$

Theorem (LSVI Generative)

LSVI with D-optimal design yields an ϵ -optimal policy with $\tilde{O}(H^6 d^2 / \epsilon^2)$ samples.

Offline RL ~~Recap~~

Switching to Offline RL

The Setting:

- We can no longer query the simulator.
- We are given static datasets D_0, \dots, D_{H-1} .
- **Key Question:** When does LSVI still work?

The Challenge:

- In Generative Mode, we used *D-Optimal Design* to ensure:

$$\Lambda_h \approx \Sigma_{\rho^*} \implies \text{Good Coverage Everywhere}$$

- In Offline RL, we are stuck with the behavior policy's distribution.

The Coverage Assumption

To guarantee success, the offline data must "cover" the feature space at least as well as the optimal design (up to a constant).

The Coverage Assumption

To guarantee success, the offline data must "cover" the feature space at least as well as the optimal design (up to a constant).

Assumption: Uniform Coverage

There exists a constant $\kappa \geq 1$ such that for all h , the empirical covariance Λ_h satisfies:

$$\Lambda_h \succeq \frac{1}{\kappa} \Sigma_{\rho^*}$$

$$\Lambda_h^{-1} \preceq \kappa \Sigma_{\rho^*}^{-1}$$

where Σ_{ρ^*} is the covariance of the D-optimal design.

The Coverage Assumption

To guarantee success, the offline data must "cover" the feature space at least as well as the optimal design (up to a constant).

Assumption: Uniform Coverage

There exists a constant $\kappa \geq 1$ such that for all h , the empirical covariance Λ_h satisfies:

$$\Lambda_h \succeq \frac{1}{\kappa} \Sigma_{\rho^*}$$

where Σ_{ρ^*} is the covariance of the D-optimal design.

Interpretation:

- κ is the "relative condition number."
- If $\kappa = 1$, our data is perfect (D-optimal).
- If κ is huge, we have missing directions (poor coverage).

Analysis of Offline LSVI

The analysis remains almost identical! We just swap the leverage bound.

Analysis of Offline LSVI

The analysis remains almost identical! We just swap the leverage bound.

1. Generative (D-Optimal):

$$\phi^\top \Lambda^{-1} \phi \approx \phi^\top \Sigma_{\rho^*}^{-1} \phi \leq d$$

Analysis of Offline LSVI

The analysis remains almost identical! We just swap the leverage bound.

1. Generative (D-Optimal):

$$\phi^\top \Lambda^{-1} \phi \approx \phi^\top \Sigma_{\rho^*}^{-1} \phi \leq d$$

2. Offline (Coverage κ):

$$\begin{aligned} \Lambda \succeq \frac{1}{\kappa} \Sigma_{\rho^*} &\implies \Lambda^{-1} \preceq \kappa \Sigma_{\rho^*}^{-1} \\ \implies \phi^\top \Lambda^{-1} \phi &\leq \kappa \left(\phi^\top \Sigma_{\rho^*}^{-1} \phi \right) \leq \kappa d \end{aligned}$$

Analysis of Offline LSVI

The analysis remains almost identical! We just swap the leverage bound.

1. Generative (D-Optimal):

$$\phi^\top \Lambda^{-1} \phi \approx \phi^\top \Sigma_{\rho^*}^{-1} \phi \leq d$$

2. Offline (Coverage κ):

$$\begin{aligned} \Lambda \succeq \frac{1}{\kappa} \Sigma_{\rho^*} &\implies \Lambda^{-1} \preceq \kappa \Sigma_{\rho^*}^{-1} \\ \implies \phi^\top \Lambda^{-1} \phi &\leq \kappa \left(\phi^\top \Sigma_{\rho^*}^{-1} \phi \right) \leq \kappa d \end{aligned}$$

Result: The error scales by $\sqrt{\kappa}$. Sample complexity scales linearly with κ :

$$\text{Error} \approx \frac{Hd\sqrt{\kappa}}{\sqrt{N}} \implies N \approx \frac{\kappa H^6 d^2}{\epsilon^2}$$

Task 2: Policy Evaluation

Sometimes we don't want to find π^* , but just evaluate a fixed policy π .

Task 2: Policy Evaluation

Sometimes we don't want to find π^* , but just evaluate a fixed policy π .

Algorithm: Least-Squares Policy Evaluation (LSPE)

- Same backward regression structure.
- **Target Change:** Instead of $\max_{a'} \hat{Q}_{h+1}(s', a')$, we use the value of our specific policy:

$$y_i = r_i + \hat{Q}_{h+1}(s'_i, \pi(s'_i))$$

Task 2: Policy Evaluation

Sometimes we don't want to find π^* , but just evaluate a fixed policy π .

Algorithm: Least-Squares Policy Evaluation (LSPE)

- Same backward regression structure.
- **Target Change:** Instead of $\max_{a'} \hat{Q}_{h+1}(s', a')$, we use the value of our specific policy:

$$y_i = r_i + \hat{Q}_{h+1}(s'_i, \pi(s'_i))$$

Weaker Assumption:

- We don't need closure under \mathcal{T} (Optimality).
- We only need closure under \mathcal{T}^π (Policy Operator).

LSPE Analysis: Why it's better

Evaluating a fixed policy is easier than finding the optimal one.

1. No "Greedy" Error Propagation

- In LSVI, we incur error from the Bellman residual AND the greedy step.
- In LSPE, we only track the estimation error of a fixed operator.

2. Tighter Simulation Lemma

- **Control (LSVI):** $V^* - V^{\hat{\pi}} \leq 2H^2\eta.$
- **Evaluation (LSPE):** $|V^\pi - \hat{V}^\pi| \leq H\eta.$

LSPE Analysis: Why it's better

Evaluating a fixed policy is easier than finding the optimal one.

1. No "Greedy" Error Propagation

- In LSVI, we incur error from the Bellman residual AND the greedy step.
- In LSPE, we only track the estimation error of a fixed operator.

2. Tighter Simulation Lemma

- **Control (LSVI):** $V^* - V^{\hat{\pi}} \leq 2H^2\eta.$
- **Evaluation (LSPE):** $|V^\pi - \hat{V}^\pi| \leq H\eta.$

Result: Better dependence on horizon H (typically H^4 instead of H^6).

1. The Critical Decomposition

To analyze error propagation, we introduce the **Infinite-Sample Target**.

Definition (f_h^*): The function LSVI would learn with infinite data.

$$f_h^* := \operatorname{argmin}_{f \in \mathcal{F}} \mathbb{E}_\rho \left[(f(s, a) - \mathcal{T}_h \hat{Q}_{h+1}(s, a))^2 \right] = \Pi_{\mathcal{F}, \rho}(\mathcal{T}_h \hat{Q}_{h+1})$$

(Ideally, we want f_h^* to track the optimal value Q_h^*).

$$= \mathcal{T} \hat{Q} \quad \text{if B.C.}$$

1. The Critical Decomposition

To analyze error propagation, we introduce the **Infinite-Sample Target**.

Definition (f_h^*): The function LSVI would learn with infinite data.

$$f_h^* := \operatorname{argmin}_{f \in \mathcal{F}} \mathbb{E}_\rho \left[(f(s, a) - \mathcal{T}_h \widehat{Q}_{h+1}(s, a))^2 \right] = \Pi_{\mathcal{F}, \rho}(\mathcal{T}_h \widehat{Q}_{h+1})$$

(Ideally, we want f_h^* to track the optimal value Q_h^*).

The Error Triangle: We split the total error into Statistical Variance and Recursive Stability.

$$\|\widehat{Q}_h - Q_h^*\|_\infty \leq \underbrace{\|\widehat{Q}_h - f_h^*\|_\infty}_{\text{Statistical Error}} + \underbrace{\|f_h^* - Q_h^*\|_\infty}_{\text{Recursive Stability}}$$

1. The Critical Decomposition

To analyze error propagation, we introduce the **Infinite-Sample Target**.

Definition (f_h^*): The function LSVI would learn with infinite data.

$$f_h^* := \operatorname{argmin}_{f \in \mathcal{F}} \mathbb{E}_\rho \left[(f(s, a) - \mathcal{T}_h \widehat{Q}_{h+1}(s, a))^2 \right] = \Pi_{\mathcal{F}, \rho}(\mathcal{T}_h \widehat{Q}_{h+1})$$

(Ideally, we want f_h^* to track the optimal value Q_h^*).

The Error Triangle: We split the total error into Statistical Variance and Recursive Stability.

$$\|\widehat{Q}_h - Q_h^*\|_\infty \leq \underbrace{\|\widehat{Q}_h - f_h^*\|_\infty}_{\text{Statistical Error}} + \underbrace{\|f_h^* - Q_h^*\|_\infty}_{\text{Recursive Stability}}$$

- **Statistical Error:** Controlled by N and D-Optimal Design ($\approx \sqrt{d/N}$).
- **Recursive Stability:** This is where the universe splits.

2. The Split Universe

How does the Recursive Error $\|f_h^* - Q_h^*\|_\infty$ behave?

Recall $f_h^* = \Pi_{\mathcal{F},\rho}(\mathcal{T}_h \hat{Q}_{h+1})$ and $Q_h^* = \Pi_{\mathcal{F},\rho}(\mathcal{T}_h Q_{h+1}^*)$ (by realizability).

2. The Split Universe

How does the Recursive Error $\|f_h^* - Q_h^*\|_\infty$ behave?

Recall $f_h^* = \Pi_{\mathcal{F},\rho}(\mathcal{T}_h \widehat{Q}_{h+1})$ and $Q_h^* = \Pi_{\mathcal{F},\rho}(\mathcal{T}_h Q_{h+1}^*)$ (by realizability).

Universe A: Completeness

Assumption: \mathcal{T}_h preserves linearity.

Since $\widehat{Q}_{h+1} \in \mathcal{F}$, we have:

$f_h^* = \Pi_{\mathcal{F},\rho}(\mathcal{T}_h \widehat{Q}_{h+1}) = \mathcal{T}_h \widehat{Q}_{h+1}$, and thus:

$$\begin{aligned}\|f_h^* - Q_h^*\|_\infty &= \|\mathcal{T}_h \widehat{Q}_{h+1} - \mathcal{T}_h Q_{h+1}^*\|_\infty \\ &\leq \|\widehat{Q}_{h+1} - Q_{h+1}^*\|_\infty\end{aligned}$$

Result: Error is stable (contraction).

2. The Split Universe

How does the Recursive Error $\|f_h^* - Q_h^*\|_\infty$ behave?

Recall $f_h^* = \Pi_{\mathcal{F},\rho}(\mathcal{T}_h \widehat{Q}_{h+1})$ and $Q_h^* = \Pi_{\mathcal{F},\rho}(\mathcal{T}_h Q_{h+1}^*)$ (by realizability).

Universe A: Completeness

Assumption: \mathcal{T}_h preserves linearity.

Since $\widehat{Q}_{h+1} \in \mathcal{F}$, we have:

$f_h^* = \Pi_{\mathcal{F},\rho}(\mathcal{T}_h \widehat{Q}_{h+1}) = \mathcal{T}_h \widehat{Q}_{h+1}$, and thus:

$$\begin{aligned} \|f_h^* - Q_h^*\|_\infty &= \|\mathcal{T}_h \widehat{Q}_{h+1} - \mathcal{T}_h Q_{h+1}^*\|_\infty \\ &\leq \|\widehat{Q}_{h+1} - Q_{h+1}^*\|_\infty \end{aligned}$$

Result: Error is stable (contraction).

Universe B: Realizability Only

Assumption: Only $Q^* \in \mathcal{F}$.

The Bellman backup $\mathcal{T}_h \widehat{Q}_{h+1}$ may be **non-linear** (off-manifold).

We must project it back to \mathcal{F} :

$$f_h^* = \Pi_{\mathcal{F},\rho}(\text{Non-Linear Target})$$

Result: We pay for the stability of the projection operator $\Pi_{\mathcal{F},\rho}$.

3. The Amplification Mechanism (Universe B)

Without Completeness, we must bound the stability of the projection $\Pi_{\mathcal{F},\rho}$. Let $\Delta = \mathcal{T}_h \widehat{Q}_{h+1} - \mathcal{T}_h Q_{h+1}^*$.

The Chain of Inequalities:

$$\begin{aligned} \|f_h^* - Q_h^*\|_\infty &= \|\Pi_{\mathcal{F},\rho} \Delta\|_\infty \\ &\leq \sqrt{d} \cdot \|\Pi_{\mathcal{F},\rho} \Delta\|_{L_2(\rho)} \quad (\text{Step 1: Norm Equivalence}) \\ &\leq \sqrt{d} \cdot \|\Delta\|_{L_2(\rho)} \quad (\text{Step 2: } L_2 \text{ Stability of LS}) \\ &\leq \sqrt{d} \cdot \|\Delta\|_\infty \quad (\text{Step 3: Norm Monotonicity}) \\ &\leq \sqrt{d} \cdot \|\widehat{Q}_{h+1} - Q_{h+1}^*\|_\infty \end{aligned}$$

The Verdict: The "price" of converting the L_2 guarantee (regression) to the L_∞ guarantee (DP) is exactly \sqrt{d} .

Total Amplification over H steps $\approx (\sqrt{d})^H$.

Summary of Lecture 2:

- **Generative Model:** LSVI + D-Optimal Design achieves $\text{poly}(d, H)$ sample complexity.
- **Offline RL:** The same algorithm works if we assume coverage (κ).
- **The Key Ingredient: Linear Bellman Completeness** ensures the regression targets remain realizable, preventing bias propagation.

Next Lecture (The Negative Results):

- What if we *don't* have Completeness?
- What if we only know Q^* is linear?
- We will show that without Completeness, sample complexity can become **Exponential in H or d** .