

CSE 542: Statistical Reinforcement Learning

Lecture 3: Offline Reinforcement Learning

Kevin Jamieson

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

Today's Plan

- 1 Motivation and Setting
- 2 MDP Basics Recap
- 3 The Offline Setup
- 4 The Plug-in MLE Policy
- 5 The Principle of Pessimism

The Data-Driven World

Modern decision-making is awash in **historical data**:

- **Healthcare:** years of electronic health records documenting treatments and outcomes
- **Autonomous driving:** millions of miles of logged human driving behavior
- **Recommender systems:** logs of user interactions with content
- **Finance:** decades of market microstructure data
- **Robotics:** teleoperated demonstrations and prior deployments

A natural question

Can we extract a *good policy* from this data alone, **without further interaction**?

Why Not Just Roll Out More?

Online RL (upcoming lectures) learns by interacting with the environment. It works well — but interaction can be:

- **Dangerous:** an ICU dosing agent that explores could harm patients
- **Expensive:** a real robot arm exploring costs time and hardware
- **Irreversible:** bad decisions in financial markets cannot be undone
- **Impossible:** the environment from which data was collected no longer exists

Why Not Just Roll Out More?

Online RL (upcoming lectures) learns by interacting with the environment. It works well — but interaction can be:

- **Dangerous:** an ICU dosing agent that explores could harm patients
- **Expensive:** a real robot arm exploring costs time and hardware
- **Irreversible:** bad decisions in financial markets cannot be undone
- **Impossible:** the environment from which data was collected no longer exists

Offline (Batch) RL

Learn a policy from a fixed dataset \mathcal{D} of historical transitions.

No further interaction allowed.

The Core Challenge: Distribution Shift

Offline RL is fundamentally harder than the generative model (Lecture 2).

Generative Model (Lec. 2)

Query *any* (s, a) on demand.

Data coverage is **our choice**.

Sample complexity: $\frac{|S||A|H^3}{\epsilon^2}$.

Offline RL

Data fixed by a **behavior policy**.

Some (s, a) pairs may be **unobserved**.

Cannot query missing transitions.

The Core Challenge: Distribution Shift

Offline RL is fundamentally harder than the generative model (Lecture 2).

Generative Model (Lec. 2)

Query *any* (s, a) on demand.
Data coverage is **our choice**.
Sample complexity: $\frac{|S||A|H^3}{\epsilon^2}$.

Offline RL

Data fixed by a **behavior policy**.
Some (s, a) pairs may be **unobserved**.
Cannot query missing transitions.

The danger: a policy that deviates from data support enters regions where the empirical model \hat{P} is unreliable — potentially catastrophically wrong.

Main message of this lecture

Two algorithms for offline RL:

- 1 **Plug-in MLE** — natural but suboptimal; pays for uncertainty under *both* π^* and $\hat{\pi}$.
- 2 **Pessimistic Value Iteration** — optimal; pays only for uncertainty under π^* .

Today's Plan

- 1 Motivation and Setting
- 2 MDP Basics Recap**
- 3 The Offline Setup
- 4 The Plug-in MLE Policy
- 5 The Principle of Pessimism

Definition

A finite-horizon MDP is $(\mathcal{S}, \mathcal{A}, \{P_h\}_h, \{r_h\}_h, H, \mu)$ where

- \mathcal{S} : finite state space, $S = |\mathcal{S}|$; \mathcal{A} : finite action space, $A = |\mathcal{A}|$
- $P_h : \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathcal{S}}$: transition kernel at stage h
- $r_h : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$: reward function at stage h
- $H \in \mathbb{N}$: horizon; $\mu \in \Delta_{\mathcal{S}}$: initial state distribution

A policy $\pi = \{\pi_h\}_h$ prescribes actions at each stage. An episode generates a trajectory $\tau = (s_1, a_1, r_1, s_2, \dots, s_H, a_H, r_H)$ where $s_1 \sim \mu$, $a_h = \pi_h(s_h)$, $s_{h+1} \sim P_h(\cdot | s_h, a_h)$.

Objective: maximize expected cumulative reward $V_0^\pi = \mathbb{E}_\pi \left[\sum_{h=1}^H r_h(s_h, a_h) \right]$.

Value Functions

State value function

$$V_h^\pi(s) = \mathbb{E} \left[\sum_{t=h}^H r_t(s_t, a_t) \mid \pi, s_h = s \right]$$

State-action value (Q) function

$$Q_h^\pi(s, a) = \mathbb{E} \left[\sum_{t=h}^H r_t(s_t, a_t) \mid \pi, s_h = s, a_h = a \right]$$

Value Functions

State value function

$$V_h^\pi(s) = \mathbb{E} \left[\sum_{t=h}^H r_t(s_t, a_t) \mid \pi, s_h = s \right]$$

State-action value (Q) function

$$Q_h^\pi(s, a) = \mathbb{E} \left[\sum_{t=h}^H r_t(s_t, a_t) \mid \pi, s_h = s, a_h = a \right]$$

Key facts:

- $V_h^\pi(s) \in [0, H - h + 1]$ and $Q_h^\pi(s, a) \in [0, H - h + 1]$
- $V_h^\pi(s) = Q_h^\pi(s, \pi_h(s))$
- **Bellman backup:** $Q_h^\pi(s, a) = r_h(s, a) + \mathbb{E}_{s' \sim P_h(\cdot | s, a)} [V_{h+1}^\pi(s')]$

We write $V_0^\pi = \mathbb{E}_{s_1 \sim \mu} [V_1^\pi(s_1)]$ for the overall value.

Optimal Policy and Value Iteration

Bellman Optimality Equations

Define $Q_h^*(s, a) = \sup_{\pi} Q_h^{\pi}(s, a)$. Then $Q_h^* = Q_h^*$ for all h iff

$$Q_h(s, a) = r_h(s, a) + \mathbb{E}_{s' \sim P_h(\cdot | s, a)} [\max_{a'} Q_{h+1}(s', a')], \quad Q_{H+1} \equiv 0.$$

The greedy policy $\pi_h^*(s) = \arg \max_a Q_h^*(s, a)$ is optimal.

Value Iteration (planning algorithm)

- 1 Set $Q_H(s, a) = r_H(s, a)$ for all (s, a) .
- 2 For $h = H - 1, \dots, 1$: $Q_h(s, a) = r_h(s, a) + \mathbb{E}_{s' \sim P_h(\cdot | s, a)} [\max_{a'} Q_{h+1}(s', a')]$.

Value iteration is **exact** when P_h is known. Today: P_h is **unknown** and we only have a fixed offline dataset.

Hoeffding's Inequality

The central concentration inequality we will use throughout.

Theorem (Hoeffding's Inequality)

Let X_1, \dots, X_N be i.i.d. random variables with $\mathbb{E}[X_i] = \mu$ and $a \leq X_i \leq b$ almost surely. For any $\delta \in (0, 1)$, with probability at least $1 - \delta$:

$$\left| \frac{1}{N} \sum_{i=1}^N X_i - \mu \right| \leq (b - a) \sqrt{\frac{\log(2/\delta)}{2N}}.$$

Hoeffding's Inequality

The central concentration inequality we will use throughout.

Theorem (Hoeffding's Inequality)

Let X_1, \dots, X_N be i.i.d. random variables with $\mathbb{E}[X_i] = \mu$ and $a \leq X_i \leq b$ almost surely. For any $\delta \in (0, 1)$, with probability at least $1 - \delta$:

$$\left| \frac{1}{N} \sum_{i=1}^N X_i - \mu \right| \leq (b - a) \sqrt{\frac{\log(2/\delta)}{2N}}.$$

Key features:

- Deviation scales as $\frac{b-a}{\sqrt{N}}$ — the *range* of X_i divided by \sqrt{N} .
- No assumption on the distribution beyond bounded support.
- **Union bound over K events:** replace $\log(2/\delta)$ with $\log(2K/\delta)$.

Hoeffding Applied to Transition Estimation

Fix a function $f : \mathcal{S} \rightarrow [0, H]$. Given $n_h(s, a)$ i.i.d. transitions from $P_h(\cdot | s, a)$:

$$|(\hat{P}_h - P_h)f|(s, a) = \left| \frac{1}{n_h(s, a)} \sum_{i=1}^{n_h(s, a)} f(s'_i) - \mathbb{E}_{s' \sim P_h(\cdot | s, a)}[f(s')] \right|.$$

Hoeffding Applied to Transition Estimation

Fix a function $f : \mathcal{S} \rightarrow [0, H]$. Given $n_h(s, a)$ i.i.d. transitions from $P_h(\cdot | s, a)$:

$$|(\hat{P}_h - P_h)f|(s, a) = \left| \frac{1}{n_h(s, a)} \sum_{i=1}^{n_h(s, a)} f(s'_i) - \mathbb{E}_{s' \sim P_h(\cdot | s, a)}[f(s')] \right|.$$

Since $f(s') \in [0, H]$, Hoeffding gives with probability $\geq 1 - \delta$:

$$|(\hat{P}_h - P_h)f|(s, a) \leq H \sqrt{\frac{\log(2/\delta)}{2 n_h(s, a)}}.$$

Hoeffding Applied to Transition Estimation

Fix a function $f : \mathcal{S} \rightarrow [0, H]$. Given $n_h(s, a)$ i.i.d. transitions from $P_h(\cdot | s, a)$:

$$|(\hat{P}_h - P_h)f|(s, a) = \left| \frac{1}{n_h(s, a)} \sum_{i=1}^{n_h(s, a)} f(s'_i) - \mathbb{E}_{s' \sim P_h(\cdot | s, a)}[f(s')] \right|.$$

Since $f(s') \in [0, H]$, Hoeffding gives with probability $\geq 1 - \delta$:

$$|(\hat{P}_h - P_h)f|(s, a) \leq H \sqrt{\frac{\log(2/\delta)}{2 n_h(s, a)}}.$$

Union bounding over all (h, s, a) triples (at most SAH of them):

$$|(\hat{P}_h - P_h)f|(s, a) \leq H \sqrt{\frac{\log(2SAH/\delta)}{2 n_h(s, a)}} \quad \text{simultaneously for all } (h, s, a).$$

This is the engine behind almost every result today.

Today's Plan

- 1 Motivation and Setting
- 2 MDP Basics Recap
- 3 The Offline Setup**
- 4 The Plug-in MLE Policy
- 5 The Principle of Pessimism

The Offline Data Model

We are given a fixed dataset of historical transitions.

Offline data assumption

For each stage $h \in \{1, \dots, H - 1\}$ and pair $(s, a) \in \mathcal{S} \times \mathcal{A}$:

- We observe $n_h(s, a) \geq 1$ i.i.d. samples $s' \sim P_h(\cdot | s, a)$.
- Sample sets are **independent** across triples (h, s, a) .

The empirical kernel is $\hat{P}_h(\cdot | s, a) = \frac{1}{n_h(s, a)} \sum_i \mathbf{1}\{s'_i = \cdot\}$.

Remarks:

- The counts $\{n_h(s, a)\}$ are fixed (determined by a behavior policy); we cannot choose them.
- This stylized model highlights the role of **data coverage**: pairs with small $n_h(s, a)$ are unreliable.
- We also assume deterministic rewards $r_h(s, a)$ are known (or can be estimated trivially).

Uncertainty Bonus

Define the **log factor** and **uncertainty bonus**:

$$\iota \triangleq \log\left(\frac{8SAH}{\delta}\right), \quad b_h(s, a) \triangleq H\sqrt{\frac{\iota}{2n_h(s, a)}}.$$

Interpretation: $b_h(s, a)$ is the Hoeffding confidence interval width for transition-weighted values at (s, a, h) , holding simultaneously over all pairs after a union bound.

Uncertainty Bonus

Define the **log factor** and **uncertainty bonus**:

$$\iota \triangleq \log\left(\frac{8SAH}{\delta}\right), \quad b_h(s, a) \triangleq H\sqrt{\frac{\iota}{2n_h(s, a)}}.$$

Interpretation: $b_h(s, a)$ is the Hoeffding confidence interval width for transition-weighted values at (s, a, h) , holding simultaneously over all pairs after a union bound.

The “Good Event” (Lemma)

With probability at least $1 - \delta$, simultaneously for every $h \in \{1, \dots, H - 1\}$ and $(s, a) \in \mathcal{S} \times \mathcal{A}$:

$$|(\hat{P}_h - P_h)V_{h+1}^{\pi^*}|(s, a) \leq b_h(s, a), \quad |(\hat{P}_h - P_h)V_{h+1}^{\hat{\pi}}|(s, a) \leq b_h(s, a),$$

$$|(\hat{P}_h - P_h)\Delta_{h+1}|(s, a) \leq b_h(s, a), \quad |(\hat{P}_h - P_h)\Gamma_{h+1}|(s, a) \leq b_h(s, a),$$

where $\Delta_h \triangleq V_h^* - \hat{V}_h$ and $\Gamma_h \triangleq \hat{V}_h - V_h^{\hat{\pi}}$.

Why Four Inequalities?

The good event controls four error quantities simultaneously:

- 1 $|(\hat{P}_h - P_h)V_{h+1}^{\pi^*}|$ — error in estimating transitions weighted by **optimal values**
- 2 $|(\hat{P}_h - P_h)V_{h+1}^{\hat{\pi}}|$ — error weighted by **learned policy's values**
- 3 $|(\hat{P}_h - P_h)\Delta_{h+1}|$ — error in the gap $V^* - \hat{V}$ (used in MLE proof)
- 4 $|(\hat{P}_h - P_h)\Gamma_{h+1}|$ — error in the gap $\hat{V} - V^{\hat{\pi}}$ (used in MLE proof)

Why Four Inequalities?

The good event controls four error quantities simultaneously:

- 1 $|(\hat{P}_h - P_h)V_{h+1}^{\pi^*}|$ — error in estimating transitions weighted by **optimal values**
- 2 $|(\hat{P}_h - P_h)V_{h+1}^{\hat{\pi}}|$ — error weighted by **learned policy's values**
- 3 $|(\hat{P}_h - P_h)\Delta_{h+1}|$ — error in the gap $V^* - \hat{V}$ (used in MLE proof)
- 4 $|(\hat{P}_h - P_h)\Gamma_{h+1}|$ — error in the gap $\hat{V} - V^{\hat{\pi}}$ (used in MLE proof)

Why can we union bound over all four?

Each function depends only on empirical kernels from *future* stages $h + 1, \dots, H - 1$. After conditioning on that future-stage data, each function is **fixed and independent** of the samples at stage h — enabling Hoeffding.

There are at most $4SAH$ such bounds; a union bound at level $\delta/(4SAH)$ each concludes the proof.

Occupancy Measures

A key quantity in the bounds: how often does policy π visit (s, a) at stage t ?

Occupancy measure

$$w_t^\pi(s, a | s', h) \triangleq \Pr_\pi[s_t = s, a_t = a | s_h = s'] .$$

$$w_t^\pi(s, a) \triangleq \Pr_\pi[s_t = s, a_t = a] \quad (\text{starting from } s_1 \sim \mu).$$

Properties:

- $\sum_{s,a} w_t^\pi(s, a) = 1$ for each t (it is a distribution over (s, a) at stage t).
- The bound we will derive depends on $\sum_{t,s,a} w_t^{\pi^*}(s, a) b_t(s, a)$.
- If π^* concentrates on well-covered pairs (large n_t), this sum is small.

Occupancy Measures

A key quantity in the bounds: how often does policy π visit (s, a) at stage t ?

Occupancy measure

$$w_t^\pi(s, a | s', h) \triangleq \Pr_\pi[s_t = s, a_t = a | s_h = s'] .$$

$$w_t^\pi(s, a) \triangleq \Pr_\pi[s_t = s, a_t = a] \quad (\text{starting from } s_1 \sim \mu).$$

Properties:

- $\sum_{s,a} w_t^\pi(s, a) = 1$ for each t (it is a distribution over (s, a) at stage t).
- The bound we will derive depends on $\sum_{t,s,a} w_t^{\pi^*}(s, a) b_t(s, a)$.
- If π^* concentrates on well-covered pairs (large n_t), this sum is small.

Key intuition

We can only do as well as our data *covers the optimal policy's trajectory*. Pairs visited by π^* but unseen in data are the fundamental bottleneck.

Today's Plan

- 1 Motivation and Setting
- 2 MDP Basics Recap
- 3 The Offline Setup
- 4 The Plug-in MLE Policy**
- 5 The Principle of Pessimism

Natural Algorithm: Plug-in MLE

The most natural approach: **estimate** P_h , **then plan**.

Plug-in MLE algorithm

- 1 Build empirical transitions: $\hat{P}_h(s'|s, a) = \frac{\#\{s' \text{ observed at } (s, a, h)\}}{n_h(s, a)}$.
- 2 Run value iteration on the **empirical MDP** $\hat{\mathcal{M}} = (\mathcal{S}, \mathcal{A}, \{\hat{P}_h\}, \{r_h\}, H, \mu)$.
- 3 Output: $\hat{\pi} = \{\hat{\pi}_h\}$ where $\hat{\pi}_h(s) = \arg \max_a \hat{Q}_h(s, a)$.

Natural Algorithm: Plug-in MLE

The most natural approach: **estimate** P_h , **then plan**.

Plug-in MLE algorithm

- 1 Build empirical transitions: $\hat{P}_h(s'|s, a) = \frac{\#\{s' \text{ observed at } (s, a, h)\}}{n_h(s, a)}$.
- 2 Run value iteration on the **empirical MDP** $\hat{\mathcal{M}} = (\mathcal{S}, \mathcal{A}, \{\hat{P}_h\}, \{r_h\}, H, \mu)$.
- 3 Output: $\hat{\pi} = \{\hat{\pi}_h\}$ where $\hat{\pi}_h(s) = \arg \max_a \hat{Q}_h(s, a)$.

Notation: \hat{V}_h^π is the value of policy π in $\hat{\mathcal{M}}$ (using \hat{P}). We write $\hat{V}_h = \hat{V}_h^{\hat{\pi}}$ for the value of the MLE policy in its own MDP.

Question: how suboptimal is $\hat{\pi}$ in the *true* MDP?

$$V_0^* - V_0^{\hat{\pi}} = ?$$

MLE Bound: Main Result

Theorem (MLE bound)

On the good event, for every stage h and state s :

$$\mathbb{E}_{s \sim \mu} [V_h^*(s) - V_h^{\hat{\pi}}(s)] \leq 2 \sum_{t=1}^{H-1} \sum_{s', a'} (w_t^{\pi^*}(s', a') + w_t^{\hat{\pi}}(s', a')) b_t(s', a')$$

MLE Bound: Main Result

Theorem (MLE bound)

On the good event, for every stage h and state s :

$$\mathbb{E}_{s \sim \mu} [V_h^*(s) - V_h^{\hat{\pi}}(s)] \leq 2 \sum_{t=1}^{H-1} \sum_{s', a'} (w_t^{\pi^*}(s', a') + w_t^{\hat{\pi}}(s', a')) b_t(s', a')$$

Two terms, two occupancies:

- First sum: uncertainty under the **optimal policy** π^* trajectory.
- Second sum: uncertainty under the **learned policy** $\hat{\pi}$ trajectory.

MLE Bound: Main Result

Theorem (MLE bound)

On the good event, for every stage h and state s :

$$\mathbb{E}_{s \sim \mu} [V_h^*(s) - V_h^{\hat{\pi}}(s)] \leq 2 \sum_{t=1}^{H-1} \sum_{s', a'} (w_t^{\pi^*}(s', a') + w_t^{\hat{\pi}}(s', a')) b_t(s', a')$$

Two terms, two occupancies:

- First sum: uncertainty under the **optimal policy** π^* trajectory.
- Second sum: uncertainty under the **learned policy** $\hat{\pi}$ trajectory.

The problem with MLE

The bound depends on where $\hat{\pi}$ goes — but $\hat{\pi}$ could stray into poorly-covered regions where $b_t(s, a)$ is large. The second term is hard to control without strong coverage assumptions on *all* policies.

Decomposing the Sub-optimality

Write $\widehat{V}_h = \widehat{V}_h^{\widehat{\pi}}$ and decompose:

$$V_h^* - V_h^{\widehat{\pi}} = \underbrace{(V_h^* - \widehat{V}_h)}_{\Delta_h} + \underbrace{(\widehat{V}_h - V_h^{\widehat{\pi}})}_{\Gamma_h}.$$

Decomposing the Sub-optimality

Write $\widehat{V}_h = \widehat{V}_h^{\widehat{\pi}}$ and decompose:

$$V_h^* - V_h^{\widehat{\pi}} = \underbrace{(V_h^* - \widehat{V}_h)}_{\Delta_h} + \underbrace{(\widehat{V}_h - V_h^{\widehat{\pi}})}_{\Gamma_h}.$$

$$\Delta_h = V_h^* - \widehat{V}_h$$

How much does the empirical MDP *underestimate* the optimal value? (Optimism of \widehat{V} would make this negative.)

$$\Gamma_h = \widehat{V}_h - V_h^{\widehat{\pi}}$$

How much does $\widehat{\pi}$'s empirical value *overestimate* its true value? (Evaluation error.)

We bound each term separately by unrolling Bellman recursions along the respective policy trajectories.

Controlling Δ_h : Proof Sketch

Let $a^* = \pi_h^*(s)$. Since $\widehat{V}_h(s) = \max_a \{r_h(s, a) + (\widehat{P}_h \widehat{V}_{h+1})(s, a)\}$:

$$\Delta_h(s) \leq r_h(s, a^*) + (P_h V_{h+1}^*)(s, a^*) - r_h(s, a^*) - (\widehat{P}_h \widehat{V}_{h+1})(s, a^*).$$

Controlling Δ_h : Proof Sketch

Let $a^* = \pi_h^*(s)$. Since $\widehat{V}_h(s) = \max_a \{r_h(s, a) + (\widehat{P}_h \widehat{V}_{h+1})(s, a)\}$:

$$\Delta_h(s) \leq r_h(s, a^*) + (P_h V_{h+1}^*)(s, a^*) - r_h(s, a^*) - (\widehat{P}_h \widehat{V}_{h+1})(s, a^*).$$

Add and subtract $(\widehat{P}_h V_{h+1}^*)(s, a^*)$:

$$\Delta_h(s) \leq \underbrace{(P_h - \widehat{P}_h)V_{h+1}^*(s, a^*)}_{|\cdot| \leq b_h(s, a^*)} + (\widehat{P}_h \Delta_{h+1})(s, a^*).$$

Add and subtract $(P_h \Delta_{h+1})(s, a^*)$:

$$\Delta_h(s) \leq (P_h \Delta_{h+1})(s, \pi_h^*(s)) + 2b_h(s, \pi_h^*(s)).$$

Controlling Δ_h : Proof Sketch

Let $a^* = \pi_h^*(s)$. Since $\widehat{V}_h(s) = \max_a \{r_h(s, a) + (\widehat{P}_h \widehat{V}_{h+1})(s, a)\}$:

$$\Delta_h(s) \leq r_h(s, a^*) + (P_h V_{h+1}^*)(s, a^*) - r_h(s, a^*) - (\widehat{P}_h \widehat{V}_{h+1})(s, a^*).$$

Add and subtract $(\widehat{P}_h V_{h+1}^*)(s, a^*)$:

$$\Delta_h(s) \leq \underbrace{(P_h - \widehat{P}_h) V_{h+1}^*}_{|\cdot| \leq b_h(s, a^*)}(s, a^*) + (\widehat{P}_h \Delta_{h+1})(s, a^*).$$

Add and subtract $(P_h \Delta_{h+1})(s, a^*)$:

$$\Delta_h(s) \leq (P_h \Delta_{h+1})(s, \pi_h^*(s)) + 2b_h(s, \pi_h^*(s)).$$

Unroll along π^* 's trajectory (using $\Delta_{H+1} \equiv 0$):

$$\Delta_h(s) \leq 2 \sum_{t=h}^{H-1} \sum_{s', a'} w_t^{\pi^*}(s', a' | s, h) b_t(s', a').$$

Controlling Γ_h : Proof Sketch

Since $\hat{\pi}$ is greedy for $\hat{\mathcal{M}}$ and $V_h^{\hat{\pi}}$ satisfies the true Bellman equation:

$$\begin{aligned}\hat{V}_h(s) &= r_h(s, \hat{\pi}_h(s)) + (\hat{P}_h \hat{V}_{h+1})(s, \hat{\pi}_h(s)), \\ V_h^{\hat{\pi}}(s) &= r_h(s, \hat{\pi}_h(s)) + (P_h V_{h+1}^{\hat{\pi}})(s, \hat{\pi}_h(s)).\end{aligned}$$

Subtracting:

$$\Gamma_h(s) = \underbrace{(\hat{P}_h - P_h)V_{h+1}^{\hat{\pi}}(s, \hat{\pi}_h(s))}_{|\cdot| \leq b_h(s, \hat{\pi}_h(s))} + (\hat{P}_h \Gamma_{h+1})(s, \hat{\pi}_h(s)).$$

Controlling Γ_h : Proof Sketch

Since $\hat{\pi}$ is greedy for $\hat{\mathcal{M}}$ and $V_h^{\hat{\pi}}$ satisfies the true Bellman equation:

$$\begin{aligned}\hat{V}_h(s) &= r_h(s, \hat{\pi}_h(s)) + (\hat{P}_h \hat{V}_{h+1})(s, \hat{\pi}_h(s)), \\ V_h^{\hat{\pi}}(s) &= r_h(s, \hat{\pi}_h(s)) + (P_h V_{h+1}^{\hat{\pi}})(s, \hat{\pi}_h(s)).\end{aligned}$$

Subtracting:


$$\Gamma_h(s) = \underbrace{(\hat{P}_h - P_h)V_{h+1}^{\hat{\pi}}(s, \hat{\pi}_h(s))}_{|\cdot| \leq b_h(s, \hat{\pi}_h(s))} + (\hat{P}_h \Gamma_{h+1})(s, \hat{\pi}_h(s)).$$

Insert and subtract $(P_h \Gamma_{h+1})(s, \hat{\pi}_h(s))$:

$$\Gamma_h(s) \leq (P_h \Gamma_{h+1})(s, \hat{\pi}_h(s)) + 2b_h(s, \hat{\pi}_h(s)).$$

Unroll along $\hat{\pi}$'s trajectory:

$$\Gamma_h(s) \leq 2 \sum_{t=h}^{H-1} \sum_{s', a'} w_t^{\hat{\pi}}(s', a' | s, h) b_t(s', a').$$

Adding the bounds for Δ_h and Γ_h proves the theorem. 

Today's Plan

- 1 Motivation and Setting
- 2 MDP Basics Recap
- 3 The Offline Setup
- 4 The Plug-in MLE Policy
- 5 The Principle of Pessimism**

Why MLE Fails: The Exploitation Trap

The MLE policy maximizes value in the **empirical MDP**, ignoring estimation uncertainty.

Illustrative failure

Suppose state s_{rare} is visited very rarely in the data (small $n_h(s_{\text{rare}}, a)$ for all a).

- $\hat{P}_h(\cdot | s_{\text{rare}}, a)$ is noisy — could erroneously suggest high reward.
- The MLE policy may route through s_{rare} to exploit this “phantom” high reward.
- In truth, the policy gets poor reward and the bound’s second term (under $\hat{\pi}$) is large.

Why MLE Fails: The Exploitation Trap

The MLE policy maximizes value in the **empirical MDP**, ignoring estimation uncertainty.

Illustrative failure

Suppose state s_{rare} is visited very rarely in the data (small $n_h(s_{\text{rare}}, a)$ for all a).

- $\hat{P}_h(\cdot | s_{\text{rare}}, a)$ is noisy — could erroneously suggest high reward.
- The MLE policy may route through s_{rare} to exploit this “phantom” high reward.
- In truth, the policy gets poor reward and the bound’s second term (under $\hat{\pi}$) is large.

The fix: penalize transitions into poorly-covered regions.

An agent should be **pessimistic** about poorly-covered state-action pairs — it should avoid them, not exploit imaginary rewards there.

This mirrors the **optimism** principle from online RL, but in the opposite direction.

Pessimistic Value Iteration (PEVI)

Algorithm: Pessimistic Value Iteration

Initialize $\underline{V}_{H+1} \equiv 0$. Backward for $h = H, H - 1, \dots, 1$:

$$\underline{Q}_h(s, a) \triangleq r_h(s, a) + (\hat{P}_h \underline{V}_{h+1})(s, a) - b_h(s, a),$$

$$\underline{V}_h(s) \triangleq \max_{a \in \mathcal{A}} \underline{Q}_h(s, a).$$

Output: $\tilde{\pi}_h(s) = \arg \max_a \underline{Q}_h(s, a)$ (greedy with respect to \underline{V}).

Pessimistic Value Iteration (PEVI)

Algorithm: Pessimistic Value Iteration

Initialize $\underline{V}_{H+1} \equiv 0$. Backward for $h = H, H-1, \dots, 1$:

$$\underline{Q}_h(s, a) \triangleq r_h(s, a) + (\hat{P}_h \underline{V}_{h+1})(s, a) - b_h(s, a),$$

$$\underline{V}_h(s) \triangleq \max_{a \in \mathcal{A}} \underline{Q}_h(s, a).$$

Output: $\tilde{\pi}_h(s) = \arg \max_a \underline{Q}_h(s, a)$ (greedy with respect to \underline{V}).

Key difference from MLE: subtract the bonus $b_h(s, a)$ at each backup.

- High-uncertainty pairs (small $n_h(s, a)$, large b_h) are **penalized**.
- The policy is discouraged from visiting regions not well-covered by data.
- Think of $\underline{Q}_h(s, a)$ as a **lower confidence bound** on the true $Q_h^*(s, a)$.

Key Lemma: Sandwiching

Lemma (Pessimistic sandwiching)

On the good event, for all h and s :

$$\underline{V}_h(s) \leq V_h^{\tilde{\pi}}(s) \leq V_h^*(s).$$

Key Lemma: Sandwiching

Lemma (Pessimistic sandwiching)

On the good event, for all h and s :

$$\underline{V}_h(s) \leq V_h^{\tilde{\pi}}(s) \leq V_h^*(s).$$

Proof: The upper bound ($V_h^{\tilde{\pi}} \leq V_h^*$) is immediate (no policy beats optimal).

For the lower bound, argue **backward induction**. Trivially $\underline{V}_{H+1} \leq V_{H+1}^{\tilde{\pi}}$. Suppose $\underline{V}_{h+1} \leq V_{h+1}^{\tilde{\pi}}$. By greediness of $\tilde{\pi}$:

$$\underline{V}_h(s) = r_h(s, \tilde{\pi}_h(s)) + (\hat{P}_h \underline{V}_{h+1})(s, \tilde{\pi}_h(s)) - b_h(s, \tilde{\pi}_h(s)).$$

On the good event:

$$(\hat{P}_h \underline{V}_{h+1})(s, \tilde{\pi}_h(s)) - b_h(s, \tilde{\pi}_h(s)) \leq (P_h \underline{V}_{h+1})(s, \tilde{\pi}_h(s)) \leq (P_h V_{h+1}^{\tilde{\pi}})(s, \tilde{\pi}_h(s)).$$

Therefore $\underline{V}_h(s) \leq V_h^{\tilde{\pi}}(s)$. □

Pessimism Bound: Main Result

Theorem (Pessimism bound)

On the good event,

$$\mathbb{E}_{s \sim \mu} [V_1^*(s) - V_1^{\tilde{\pi}}(s)] \leq 2 \sum_{t=1}^{H-1} \sum_{s', a'} w_t^{\pi^*}(s', a') b_t(s', a').$$

Consequently, the same bound holds with probability at least $1 - \delta$.

Pessimism Bound: Main Result

Theorem (Pessimism bound)

On the good event,

$$\mathbb{E}_{s \sim \mu} [V_1^*(s) - V_1^{\tilde{\pi}}(s)] \leq 2 \sum_{t=1}^{H-1} \sum_{s', a'} w_t^{\pi^*}(s', a') b_t(s', a').$$

Consequently, the same bound holds with probability at least $1 - \delta$.

Crucial improvement over MLE

Only one occupancy term — under π^* , not $\tilde{\pi}$.

The second term (evaluation error under the learned policy) has **vanished**.

Pessimism Bound: Proof

By the sandwiching lemma: $V_h^* - V_h^{\tilde{\pi}} \leq V_h^* - \underline{V}_h =: \Delta_h$.

Let $a^* = \pi_h^*(s)$. Since $\underline{V}_h(s) \geq \underline{Q}_h(s, a^*)$:

$$\Delta_h(s) \leq r_h(s, a^*) + (P_h V_{h+1}^*)(s, a^*) - r_h(s, a^*) - (\widehat{P}_h \underline{V}_{h+1})(s, a^*) + b_h(s, a^*).$$

Pessimism Bound: Proof

By the sandwiching lemma: $V_h^* - V_h^{\tilde{\pi}} \leq V_h^* - \underline{V}_h =: \Delta_h$.

Let $a^* = \pi_h^*(s)$. Since $\underline{V}_h(s) \geq \underline{Q}_h(s, a^*)$:

$$\Delta_h(s) \leq r_h(s, a^*) + (P_h V_{h+1}^*)(s, a^*) - r_h(s, a^*) - (\widehat{P}_h \underline{V}_{h+1})(s, a^*) + b_h(s, a^*).$$

Add and subtract $(P_h \underline{V}_{h+1})(s, a^*)$:

$$\Delta_h(s) \leq (P_h \Delta_{h+1})(s, a^*) + \underbrace{(P_h - \widehat{P}_h) \underline{V}_{h+1}(s, a^*)}_{|\cdot| \leq b_h(s, a^*)} + b_h(s, a^*).$$

On the good event, the middle term is at most $b_h(s, a^*)$, hence:

$$\Delta_h(s) \leq (P_h \Delta_{h+1})(s, \pi_h^*(s)) + 2b_h(s, \pi_h^*(s)).$$

Pessimism Bound: Proof

By the sandwiching lemma: $V_h^* - V_h^{\tilde{\pi}} \leq V_h^* - \underline{V}_h =: \Delta_h$.

Let $a^* = \pi_h^*(s)$. Since $\underline{V}_h(s) \geq \underline{Q}_h(s, a^*)$:

$$\Delta_h(s) \leq r_h(s, a^*) + (P_h V_{h+1}^*)(s, a^*) - r_h(s, a^*) - (\widehat{P}_h \underline{V}_{h+1})(s, a^*) + b_h(s, a^*).$$

Add and subtract $(P_h \underline{V}_{h+1})(s, a^*)$:

$$\Delta_h(s) \leq (P_h \Delta_{h+1})(s, a^*) + \underbrace{(P_h - \widehat{P}_h) \underline{V}_{h+1}(s, a^*)}_{|\cdot| \leq b_h(s, a^*)} + b_h(s, a^*).$$

On the good event, the middle term is at most $b_h(s, a^*)$, hence:

$$\Delta_h(s) \leq (P_h \Delta_{h+1})(s, \pi_h^*(s)) + 2b_h(s, \pi_h^*(s)).$$

Unroll along π^* 's trajectory (using $\Delta_{H+1} \equiv 0$):

$$\Delta_h(s) \leq 2 \sum_{t=h}^{H-1} \sum_{s', a'} w_t^{\pi^*}(s', a' | s, h) b_t(s', a'). \quad \square$$

MLE bound

$$V_h^* - V_h^{\hat{\pi}} \leq 2 \sum_{t,s',a'} w_t^{\pi^*}(s', a') b_t(s', a') + 2 \sum_{t,s',a'} w_t^{\hat{\pi}}(s', a') b_t(s', a').$$

Side-by-Side Comparison

MLE bound

$$V_h^* - V_h^{\hat{\pi}} \leq 2 \sum_{t,s',a'} w_t^{\pi^*}(s', a') b_t(s', a') + 2 \sum_{t,s',a'} w_t^{\hat{\pi}}(s', a') b_t(s', a').$$

Pessimism bound

$$V_h^* - V_h^{\hat{\pi}} \leq 2 \sum_{t,s',a'} w_t^{\pi^*}(s', a') b_t(s', a').$$

Side-by-Side Comparison

MLE bound

$$V_h^* - V_h^{\hat{\pi}} \leq 2 \sum_{t,s',a'} w_t^{\pi^*}(s', a') b_t(s', a') + 2 \sum_{t,s',a'} w_t^{\hat{\pi}}(s', a') b_t(s', a').$$

Pessimism bound

$$V_h^* - V_h^{\tilde{\pi}} \leq 2 \sum_{t,s',a'} w_t^{\pi^*}(s', a') b_t(s', a').$$

Both bounds share the first term. Pessimism **eliminates the second term**.

- The second MLE term can be *arbitrarily larger* than the first if $\hat{\pi}$ explores poorly-covered regions.
- Pessimism automatically steers $\tilde{\pi}$ away from such regions, so no second term arises.

Why Does Pessimism Win?

The real benefit of pessimism is not that concentration becomes easier.

Both algorithms face the same statistical difficulty in estimating transitions. The difference is **structural**:

Key insight

The lower bound $\underline{V}_h \leq V_h^{\tilde{\pi}}$ (sandwiching lemma) prevents the evaluation-error term $(\hat{V}_h - V_h^{\tilde{\pi}})$ from appearing in the analysis.

We only need to track how V^* decomposes along π^* 's trajectory — a trajectory we can analyze using concentration on V^* , which is **fixed** (not data-dependent).

Why Does Pessimism Win?

The real benefit of pessimism is not that concentration becomes easier.

Both algorithms face the same statistical difficulty in estimating transitions. The difference is **structural**:

Key insight

The lower bound $\underline{V}_h \leq V_h^{\tilde{\pi}}$ (sandwiching lemma) prevents the evaluation-error term $(\hat{V}_h - V_h^{\tilde{\pi}})$ from appearing in the analysis.

We only need to track how V^* decomposes along π^* 's trajectory — a trajectory we can analyze using concentration on V^* , which is **fixed** (not data-dependent).

Contrast with MLE: to bound $(\hat{V}_h - V_h^{\hat{\pi}})$, we must unroll along $\hat{\pi}$'s trajectory, which itself depends on the data — creating the problematic second occupancy term.

Explicit Sample Complexity

Substituting $b_t(s, a) = H\sqrt{\ell/(2n_t(s, a))}$ into the pessimism bound and assuming **uniform coverage** $n_t(s, a) = n$ for all (t, s, a) visited by π^* :

$$V_0^* - V_0^{\tilde{\pi}} \leq 2 \sum_{t=1}^{H-1} \sum_{s,a} w_t^{\pi^*}(s, a) H \sqrt{\frac{\ell}{2n}} \leq 2H^2 \sqrt{\frac{\ell}{2n}}.$$

Setting this $\leq \varepsilon$ and solving for n :

$$n \gtrsim \frac{H^4 \log(SAH/\delta)}{\varepsilon^2}.$$

Explicit Sample Complexity

Substituting $b_t(s, a) = H\sqrt{\iota/(2n_t(s, a))}$ into the pessimism bound and assuming **uniform coverage** $n_t(s, a) = n$ for all (t, s, a) visited by π^* :

$$V_0^* - V_0^{\tilde{\pi}} \leq 2 \sum_{t=1}^{H-1} \sum_{s,a} w_t^{\pi^*}(s, a) H \sqrt{\frac{\iota}{2n}} \leq 2H^2 \sqrt{\frac{\iota}{2n}}.$$

Setting this $\leq \varepsilon$ and solving for n :

$$n \gtrsim \frac{H^4 \log(SAH/\delta)}{\varepsilon^2}.$$

Sample complexity (offline RL, pessimistic policy)

Under uniform coverage with n samples per (s, a, h) in π^* 's support:

$$\text{total samples} \lesssim \frac{H^4 \cdot |\text{supp}(\pi^*)| \cdot \log(SAH/\delta)}{\varepsilon^2},$$

where $|\text{supp}(\pi^*)|$ counts the state-action pairs reachable under π^* .

Next lecture: Online RL and Exploration

In offline RL, coverage is *given* — we either have it or we don't.

In online RL, we **choose** our data collection policy, adapting based on what we see.

- The agent must balance **exploration** (reducing uncertainty) and **exploitation** (acting well now).
- The right lens: **regret** — total sub-optimality over K episodes.
- Key algorithm: **UCB-VI** — optimistic (not pessimistic!) value iteration.

Looking Ahead

Next lecture: Online RL and Exploration

In offline RL, coverage is *given* — we either have it or we don't.

In online RL, we **choose** our data collection policy, adapting based on what we see.

- The agent must balance **exploration** (reducing uncertainty) and **exploitation** (acting well now).
- The right lens: **regret** — total sub-optimality over K episodes.
- Key algorithm: **UCB-VI** — optimistic (not pessimistic!) value iteration.

Optimism vs. Pessimism

- **Offline RL** \Rightarrow **Pessimism**: penalize uncertain actions to stay in-distribution.
- **Online RL** \Rightarrow **Optimism**: bonus uncertain actions to incentivize exploration.