

# Lecture 12: Offline RL

This lecture studies *offline reinforcement learning* (also known as *batch RL*), where the agent must learn a good policy from a fixed dataset collected a priori, without any further interaction with the environment. The central principle is *pessimism in the face of uncertainty*—the mirror image of the optimism principle studied in Lecture 11. We cover three main ideas: the suboptimality decomposition that explains why pessimism is needed, *Pessimistic Value Iteration* (PEVI) for *linear MDPs* (Jin et al., 2021), and *Bellman-consistent pessimism* (BCP) for general function classes (Xie et al., 2021). The PEVI portion of the lecture remains in the *finite-horizon episodic* setting of Lecture 11, while the BCP portion later switches to a *discounted infinite-horizon* setting following Xie et al. (2021).

**Setup.** For the PEVI part of the lecture, we continue with the *finite-horizon episodic MDP*  $M = (\mathcal{S}, \mathcal{A}, H, \{P_h^*\}, \{r_h\}, s_0)$  from Lecture 11, with state space  $\mathcal{S}$ , action space  $\mathcal{A}$ , horizon  $H$ , transition kernel  $P_h^*(\cdot | s, a)$ , deterministic reward  $r_h(s, a) \in [0, 1]$ , and fixed initial state  $s_0$ . Value functions and Q-functions are defined as before. For any Q-function  $g : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , write  $V_g(s) := \max_{a \in \mathcal{A}} g(s, a)$ . In the PEVI sections, we use the Bellman optimality operator on next-step Q-functions:

$$(\mathcal{T}_h g)(s, a) := r_h(s, a) + (P_h^* V_g)(s, a),$$

where  $(P_h^* u)(s, a) = \mathbb{E}_{s' \sim P_h^*(\cdot | s, a)}[u(s')]$  for any state-value function  $u : \mathcal{S} \rightarrow \mathbb{R}$ .

**Finite-horizon Bellman equations (used in the PEVI sections):** For any deterministic policy  $\pi$ ,

$$Q_h^\pi(s, a) = r_h(s, a) + P_h^* V_{h+1}^\pi(s, a), \quad V_h^\pi(s) = Q_h^\pi(s, \pi_h(s)), \quad V_H^\pi \equiv 0.$$

At the optimal policy  $\pi^*$ , these specialize to the *Bellman optimality equations*:

$$Q_h^*(s, a) = r_h(s, a) + P_h^* V_{h+1}^*(s, a) = (\mathcal{T}_h Q_{h+1}^*)(s, a), \quad V_h^*(s) = \max_{a \in \mathcal{A}} Q_h^*(s, a), \quad V_H^* \equiv 0.$$

The performance metric in offline RL is the *suboptimality* of a learned policy  $\pi$ :

$$\text{SubOpt}(\pi; s_0) := V_0^*(s_0) - V_0^\pi(s_0).$$

## The Offline RL Problem

### Setting and Data Collecting Process

In the offline (or batch) setting, the learner has access only to a pre-collected dataset  $\mathcal{D}$ —there is no ability to interact with the environment or collect new data. This is in stark contrast to the online exploration problem of Lecture 11, where the agent adaptively collects trajectories over  $K$  episodes.

**Definition 1** (Offline dataset). *The offline dataset consists of  $K$  trajectories:*

$$\mathcal{D} = \left\{ (s_h^\tau, a_h^\tau, r_h^\tau, s_{h+1}^\tau) \right\}_{\tau=0, h=0}^{K-1, H-1},$$

where for each trajectory  $\tau \in \{0, \dots, K-1\}$ , the experimenter takes action  $a_h^\tau$  at state  $s_h^\tau$ , observes reward  $r_h^\tau = r_h(s_h^\tau, a_h^\tau)$ , and the next state is sampled  $s_{h+1}^\tau \sim P_h^*(\cdot | s_h^\tau, a_h^\tau)$ . The actions  $a_h^\tau$  are chosen by a (possibly adaptive) experimenter and are not under the learner's control.

The key assumption on the dataset is *compliance*: the conditional distribution of  $(r_h^\tau, s_{h+1}^\tau)$  given the current and past data is determined by the true MDP.

**Assumption 1** (Compliance). *The dataset  $\mathcal{D}$  is compliant with the underlying MDP  $(\mathcal{S}, \mathcal{A}, H, P^*, r)$ : for each trajectory  $\tau$  and step  $h$ , conditioning on  $(s_h^\tau, a_h^\tau)$  and all data from earlier trajectories and earlier steps within trajectory  $\tau$ , the tuple  $(r_h^\tau, s_{h+1}^\tau)$  is generated by the reward function  $r_h$  and transition kernel  $P_h^*$  of the underlying MDP.*

**Generality of compliance.** Assumption 1 is deliberately permissive: the actions  $a_h^\tau$  can be chosen by any mechanism—a fixed behavior policy, an adaptive strategy, or even an adversary. The  $K$  trajectories need not be independent or identically distributed. All that is required is that the *transitions and rewards* follow the true MDP. This captures diverse practical scenarios: data pooled from multiple operators, data collected by an evolving policy, or data from a simulator with scripted behaviors.

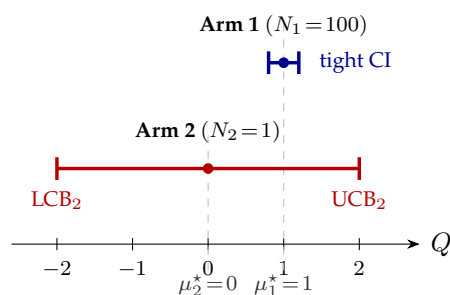
### Online vs. Offline: From Optimism to Pessimism

In Lecture 11, we saw that the online exploration problem is solved by the *optimism* principle: the agent inflates uncertain Q-values with bonuses, which drives it to explore under-visited states. The key structural feature of the analysis is:

$$\text{Online regret} \leq \text{sum of bonuses along the learned policy's trajectory.}$$

This works because the agent controls which states to visit, and optimism ensures it visits informative ones.

In the offline setting, the agent has *no control* over the data. A natural question is whether we can simply run an online optimistic algorithm (e.g., UCBVI or Lin-UCBVI) on the offline dataset. Figure 1 shows why the answer is *no*: optimism’s bonus mechanism picks in exactly the wrong direction when revisiting an arm is impossible. The remedy is to flip the sign—subtract a penalty (pessimism) instead of adding a bonus—so that poorly-covered state-action pairs are assigned *low* Q-values and the learned policy avoids them.



- |   |                                  |
|---|----------------------------------|
| <b>Greedy:</b> $\arg \max \hat{\mu}$        | • Arm 1 ✓ (but fragile to noise) |
| <b>UCB:</b> $\arg \max(\hat{\mu} + \Gamma)$ | • Arm 2 ✗ (systematically wrong) |
| <b>LCB:</b> $\arg \max(\hat{\mu} - \Gamma)$ | • Arm 1 ✓ (robust)               |

Figure 1: **Why optimism fails offline.** Bandit with  $\mu_1^* = 1$ ,  $\mu_2^* = 0$ , data  $N_1 = 100$ ,  $N_2 = 1$ , and point estimates equal to the true means ( $\hat{\mu}_1 = 1$ ,  $\hat{\mu}_2 = 0$ ; *no bad luck*). The penalty/bonus  $\Gamma(a) = c/\sqrt{N(a)}$  makes Arm 2’s confidence interval much wider than Arm 1’s. **UCB** picks the upper endpoint of the CI and thus *deterministically* selects the under-sampled Arm 2—its bonus rewards data scarcity, which is exactly the wrong signal offline. **LCB** picks the lower endpoint and robustly selects Arm 1. Online, UCB would revisit Arm 2 and the bonus would shrink; offline, no revisits are possible, and the wrong-direction correction is permanent.

**The pessimism principle:**

$$\underbrace{\hat{Q}_h = r_h + \hat{P}_h \hat{V}_{h+1} + b_h}_{\text{optimism (online, Lecture 11)}} \longrightarrow \underbrace{\hat{Q}_h = r_h + \hat{P}_h \hat{V}_{h+1} - \Gamma_h}_{\text{pessimism (offline, this lecture)}}$$

The penalty  $\Gamma_h(s, a) \geq 0$  has the *same form* as the exploration bonus, but the sign is flipped.

## What Causes Suboptimality?

Before presenting the PEVI algorithm, we analyze *what causes suboptimality* for any algorithm that constructs estimated Q-functions and a greedy policy from the dataset. This decomposition, due to Jin et al. (2021), is the finite-horizon offline analogue of the Perfor-

mance Difference via Bellman Error lemma from Lecture 7, and also the same standard policy-loss decomposition reused in the episodic regret analysis of Lecture 11. It reveals why pessimism is the right principle in the offline setting.

## Decomposition of Suboptimality

Suppose an algorithm produces estimated Q-functions  $\widehat{Q}_h : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  from the dataset  $\mathcal{D}$ . Define

$$\widehat{Q}_H(\cdot, \cdot) := 0, \quad \widehat{V}_h(s) := \max_{a \in \mathcal{A}} \widehat{Q}_h(s, a), \quad \widehat{\pi}_h(s) \in \operatorname{argmax}_{a \in \mathcal{A}} \widehat{Q}_h(s, a),$$

and note that the *Bellman error* at step  $h$  is

$$(\mathcal{T}_h \widehat{Q}_{h+1})(s, a) - \widehat{Q}_h(s, a),$$

which measures how well  $\widehat{Q}_h$  approximates the *true* Bellman backup of  $\widehat{Q}_{h+1}$ .

**Lemma 1** (Decomposition of suboptimality). *Let  $\widehat{\pi} = \{\widehat{\pi}_h\}_{h=0}^{H-1}$  be the greedy policy induced by  $\widehat{Q}$ . Then for any initial state  $s_0$ ,*

$$\begin{aligned} \text{SubOpt}(\widehat{\pi}; s_0) &= \underbrace{\sum_{h=0}^{H-1} \mathbb{E}_{\widehat{\pi}} \left[ \widehat{Q}_h(s_h, a_h) - (\mathcal{T}_h \widehat{Q}_{h+1})(s_h, a_h) \mid s_0 \right]}_{\text{(i): Pessimism-controlled term}} \\ &+ \underbrace{\sum_{h=0}^{H-1} \mathbb{E}_{\pi^*} \left[ (\mathcal{T}_h \widehat{Q}_{h+1})(s_h, a_h) - \widehat{Q}_h(s_h, a_h) \mid s_0 \right]}_{\text{(ii): Optimism-controlled term}} \\ &+ \underbrace{\sum_{h=0}^{H-1} \mathbb{E}_{\pi^*} \left[ \widehat{Q}_h(s_h, \pi_h^*(s_h)) - \widehat{Q}_h(s_h, \widehat{\pi}_h(s_h)) \mid s_0 \right]}_{\text{(iii): Optimization Error}}. \end{aligned} \quad (1)$$

Here  $\mathbb{E}_{\widehat{\pi}}$  and  $\mathbb{E}_{\pi^*}$  denote expectations over trajectories induced by  $\widehat{\pi}$  and  $\pi^*$  in the true MDP, given the fixed functions  $\widehat{Q}_{h+1}$  and  $\widehat{Q}_h$ .

*Proof.* The proof is the same telescoping argument behind the *Performance Difference via Bellman Error* lemma from Lecture 7, applied once with  $\pi = \pi^*$  and once with  $\pi = \widehat{\pi}$  and then subtracted. Fix  $s_0$  and abbreviate  $\mathbb{E}_{\pi}[\cdot] = \mathbb{E}_{\pi}[\cdot \mid s_0]$ . For any policy  $\pi$ , telescoping the Bellman equation  $V_h^\pi(s) = Q_h^\pi(s, \pi_h(s)) = r_h(s, \pi_h(s)) + P_h^* V_{h+1}^\pi(s, \pi_h(s))$  and using

$(\mathcal{T}_h \widehat{Q}_{h+1})(s, a) = r_h(s, a) + P_h^* \widehat{V}_{h+1}(s, a)$  gives

$$V_0^\pi(s_0) - \widehat{V}_0(s_0) = \sum_{h=0}^{H-1} \mathbb{E}_\pi \left[ (\mathcal{T}_h \widehat{Q}_{h+1})(s_h, a_h) - \widehat{Q}_h(s_h, a_h) + \widehat{Q}_h(s_h, \pi_h(s_h)) - \widehat{V}_h(s_h) \right]. \quad (2)$$

Since  $\widehat{V}_h(s) = \widehat{Q}_h(s, \widehat{\pi}_h(s))$ , applying (2) to  $\pi = \pi^*$  yields

$$V_0^*(s_0) - \widehat{V}_0(s_0) = \sum_{h=0}^{H-1} \mathbb{E}_{\pi^*} \left[ (\mathcal{T}_h \widehat{Q}_{h+1})(s_h, a_h) - \widehat{Q}_h(s_h, a_h) + \widehat{Q}_h(s_h, \pi_h^*(s_h)) - \widehat{Q}_h(s_h, \widehat{\pi}_h(s_h)) \right],$$

and applying (2) to  $\pi = \widehat{\pi}$  yields

$$V_0^{\widehat{\pi}}(s_0) - \widehat{V}_0(s_0) = \sum_{h=0}^{H-1} \mathbb{E}_{\widehat{\pi}} \left[ (\mathcal{T}_h \widehat{Q}_{h+1})(s_h, a_h) - \widehat{Q}_h(s_h, a_h) \right],$$

since  $\widehat{Q}_h(s_h, \widehat{\pi}_h(s_h)) - \widehat{V}_h(s_h) = 0$ . Subtracting the second from the first and rewriting the first difference as  $\widehat{Q}_h - \mathcal{T}_h \widehat{Q}_{h+1}$  gives (1).  $\square$

#### Interpreting the three sources of suboptimality.

- **Term (i): pessimism-controlled term.** This is the Bellman error along the trajectory of  $\widehat{\pi}$ , written in the order  $\widehat{Q}_h - \mathcal{T}_h \widehat{Q}_{h+1}$ . If  $\widehat{Q}_h \leq (\mathcal{T}_h \widehat{Q}_{h+1})$  pointwise, equivalently  $\widehat{Q}_h - (\mathcal{T}_h \widehat{Q}_{h+1}) \leq 0$ , then term (i) is non-positive. Without such sign control, this term can be badly biased because  $\widehat{\pi}$  is chosen from the same data used to estimate  $\widehat{Q}$ .
- **Term (ii): optimism-controlled term.** This is the Bellman error along the trajectory of  $\pi^*$ . It is the comparator-policy term that already appeared in the telescoping / Bellman-error decompositions from earlier lectures. In the symmetric optimistic direction, if  $\widehat{Q}_h \geq (\mathcal{T}_h \widehat{Q}_{h+1})$  pointwise, then this term is controlled with the favorable sign. In offline RL, we instead upper-bound it by uncertainty along  $\pi^*$ 's trajectory.
- **Term (iii): Optimization error.** This term measures the gap  $\widehat{Q}_h(s_h, \pi_h^*(s_h)) - \widehat{Q}_h(s_h, \widehat{\pi}_h(s_h))$  under  $\pi^*$ 's trajectory. It is non-positive when  $\widehat{\pi}$  is greedy with respect to  $\widehat{Q}_h$  (since then  $\widehat{Q}_h(s, \widehat{\pi}_h(s)) \geq \widehat{Q}_h(s, a)$  for all  $a$ ), and hence can be dropped.

The key insight: for the greedy policy  $\widehat{\pi}$ , term (iii)  $\leq 0$ , and pessimism eliminates the pessimism-controlled term (i), so the remaining burden is to control the optimism-controlled term (ii). This depends on how well  $\mathcal{D}$  covers the *optimal policy's* trajectory, not on uniform coverage of the entire state-action space.

## Illustration via a Special Case: Multi-Armed Bandit

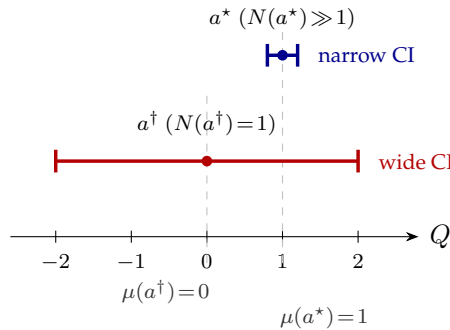
To build intuition, consider the multi-armed bandit (MAB), a special case with  $|\mathcal{S}| = 1$  and  $H = 1$ . There are  $A$  arms with mean rewards  $\mu(a)$ ,  $a \in \mathcal{A}$ . The dataset consists of  $K$  pulls  $\mathcal{D} = \{(a^\tau, r^\tau)\}_{\tau=0}^{K-1}$ , and the sample average estimator is

$$\hat{\mu}(a) = \frac{1}{N(a)} \sum_{\tau: a^\tau = a} r^\tau, \quad N(a) = \sum_{\tau=0}^{K-1} \mathbb{1}\{a^\tau = a\}.$$

**Bandit Decomposition.** Let  $\hat{Q} : \mathcal{A} \rightarrow \mathbb{R}$  be any bandit score and let  $\hat{\pi} = \operatorname{argmax}_a \hat{Q}(a)$ . Since the Bellman error is  $\mu(a) - \hat{Q}(a)$  and the optimization term is  $\hat{Q}(\pi^*) - \hat{Q}(\hat{\pi}) \leq 0$ , the decomposition (1) becomes

$$\text{SubOpt}(\hat{\pi}) \leq \underbrace{\mathbb{E}_{\hat{\pi}} \left[ \hat{Q}(\hat{\pi}) - \mu(\hat{\pi}) \right]}_{\text{(i)}} + \underbrace{\mathbb{E}_{\pi^*} \left[ \mu(\pi^*) - \hat{Q}(\pi^*) \right]}_{\text{(ii)}}.$$

For naive greedy,  $\hat{Q} = \hat{\mu}$ , so  $\hat{\pi} = \operatorname{argmax}_a \hat{\mu}(a)$  and the estimation error  $\hat{\mu}(a) - \mu(a)$  is mean zero for each arm. However, term (i) becomes  $\hat{\mu}(\hat{\pi}) - \mu(\hat{\pi})$ , which is *positively biased* because  $\hat{\pi}$  is selected to maximize  $\hat{\mu}$ .



- Term (ii)** **Oracle-controlled:** only uses the well-covered arm  $a^*$ .
- Term (i)** **Selection-biased:**  $\hat{\pi}$  may jump to the poorly covered arm  $a^\dagger$ .
- Pessimism** **Pointwise pessimism:** if  $\hat{Q}(a) \leq \mu(a)$  for all  $a$ , then  $\mathbb{E}_{\hat{\pi}}[\hat{Q}(\hat{\pi}) - \mu(\hat{\pi})] \leq 0$ .

Figure 2: **Bandit view of the offline decomposition.** Term (ii) is small when the optimal arm  $a^*$  is well covered. Term (i) is not oracle-controlled because the data-dependent policy  $\hat{\pi}$  can jump to the wide-CI arm  $a^\dagger$ . Pessimism penalizes  $a^\dagger$  and makes term (i) non-positive.

Figure 2 isolates the asymmetry in (1): term (ii) is oracle-controlled by the coverage of  $a^*$ , whereas term (i) is evaluated at the data-dependent arm  $\hat{\pi}$  and can therefore be driven by the poorly covered arm  $a^\dagger$ . Pessimism is stronger than concentration: if  $\hat{Q}(a) = \hat{\mu}(a) - \Gamma(a)$

satisfies  $\widehat{Q}(a) \leq \mu(a)$  for all  $a$ , then

$$\mathbb{E}_{\widehat{\pi}} \left[ \widehat{Q}(\widehat{\pi}) - \mu(\widehat{\pi}) \right] \leq 0,$$

so only the oracle-controlled term remains.

## PEVI for Linear MDPs

We now study Pessimistic Value Iteration (PEVI) as a *linear-MDP algorithm*. The presentation is intentionally parallel to Lecture 11's treatment of Lin-UCBVI: the Bellman backup is estimated by linear regression, uncertainty is measured by an ellipsoidal norm, and the only conceptual change is that the uncertainty term is *subtracted* rather than added.

### The Linear MDP

We reuse the linear transition model from Lecture 11 and add a linear reward model, restated here for convenience.

**Definition 2** (Linear MDP). *An episodic MDP satisfies the linear MDP structure with feature map  $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$  if there exist unknown signed measures  $\mu_h^* = (\mu_h^{*(1)}, \dots, \mu_h^{*(d)})$  over  $\mathcal{S}$  and an unknown vector  $\theta_h \in \mathbb{R}^d$  such that*

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

We assume  $\|\phi(s, a)\|_2 \leq 1$  for all  $(s, a)$  and  $\max\{\|\mu_h^*(\mathcal{S})\|_2, \|\theta_h\|_2\} \leq \sqrt{d}$ .

The key property is that the Bellman backup remains linear in  $\phi$ . By the linearization of the transition term from Lecture 11, for any state-value function  $f : \mathcal{S} \rightarrow \mathbb{R}$ ,

$$(P_h^* f)(s, a) = \langle \phi(s, a), (\mu_h^*)^\top f \rangle.$$

Therefore, for any Q-function  $g : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ ,

$$(\mathcal{T}_h g)(s, a) = r_h(s, a) + (P_h^* V_g)(s, a) = \langle \phi(s, a), w_{h,g}^* \rangle, \quad \text{where } w_{h,g}^* := \theta_h + (\mu_h^*)^\top V_g \in \mathbb{R}^d.$$

This linearization holds regardless of whether  $g$  is itself linear—the Bellman backup of any bounded next-step Q-function remains linear in  $\phi$ .

**Algorithm 1** PEVI (episodic, linear MDP)

**Require:** Dataset  $\mathcal{D} = \{(s_h^\tau, a_h^\tau, r_h^\tau, s_{h+1}^\tau)\}_{\tau, h=0}^{K-1, H-1}$ , regularization  $\lambda \geq 1$ , confidence  $\beta > 0$

- 1: **Initialize:**  $\widehat{Q}_H(\cdot, \cdot) \leftarrow 0, \widehat{V}_H(\cdot) \leftarrow 0$
- 2: **for**  $h = H - 1, H - 2, \dots, 0$  **do**
- 3:    $\Lambda_h \leftarrow \lambda I + \sum_{\tau=0}^{K-1} \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top$  ▷ Design matrix
- 4:    $\widehat{w}_h \leftarrow \Lambda_h^{-1} \sum_{\tau=0}^{K-1} \phi(s_h^\tau, a_h^\tau) (r_h^\tau + \widehat{V}_{h+1}(s_{h+1}^\tau))$  ▷ Ridge regression
- 5:    $\Gamma_h(\cdot, \cdot) \leftarrow \beta \cdot (\phi(\cdot, \cdot)^\top \Lambda_h^{-1} \phi(\cdot, \cdot))^{1/2}$  ▷ Uncertainty
- 6:    $\widehat{Q}_h(\cdot, \cdot) \leftarrow \min\{\langle \phi(\cdot, \cdot), \widehat{w}_h \rangle - \Gamma_h(\cdot, \cdot), H - h\}^+$  ▷ Pessimistic clipping
- 7:    $\widehat{\pi}_h(\cdot) \leftarrow \operatorname{argmax}_a \widehat{Q}_h(\cdot, a)$  ▷ Optimization
- 8:    $\widehat{V}_h(\cdot) \leftarrow \widehat{Q}_h(\cdot, \widehat{\pi}_h(\cdot))$  ▷ Evaluation
- 9: **end for**
- 10: **Output:**  $\text{Pess}(\mathcal{D}) = \{\widehat{\pi}_h\}_{h=0}^{H-1}$

**Algorithm: Linear PEVI**

**Ridge regression for the Bellman backup.** At each step  $h$ , define the regularized design matrix

$$\Lambda_h := \lambda I + \sum_{\tau=0}^{K-1} \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top, \quad \lambda \geq 1, \quad (3)$$

and the ridge regression estimate of  $w_{h, \widehat{Q}_{h+1}}^*$ :

$$\widehat{w}_h := \Lambda_h^{-1} \sum_{\tau=0}^{K-1} \phi(s_h^\tau, a_h^\tau) (r_h^\tau + \widehat{V}_{h+1}(s_{h+1}^\tau)), \quad (4)$$

and, more generally, for any Q-function  $g$  define the empirical Bellman backup

$$(\widehat{\mathcal{T}}_h g)(s, a) := \phi(s, a)^\top \Lambda_h^{-1} \sum_{\tau=0}^{K-1} \phi(s_h^\tau, a_h^\tau) (r_h^\tau + V_g(s_{h+1}^\tau)).$$

In particular,  $(\widehat{\mathcal{T}}_h \widehat{Q}_{h+1})(s, a) = \langle \phi(s, a), \widehat{w}_h \rangle$ .

**Ellipsoidal penalty.** The penalty function is the Mahalanobis norm:

$$\Gamma_h(s, a) := \beta \cdot (\phi(s, a)^\top \Lambda_h^{-1} \phi(s, a))^{1/2} = \beta \cdot \|\phi(s, a)\|_{\Lambda_h^{-1}}, \quad (5)$$

where  $\beta > 0$  is a scaling parameter set below. This is the *same form* as the Lin-UCBVI bonus from Lecture 11, but now it is *subtracted* from the estimate rather than added.

**PEVI for linear MDP:**

$$\widehat{Q}_h(s, a) = \left[ \underbrace{\langle \phi(s, a), \widehat{w}_h \rangle}_{\text{ridge regression estimate}} - \underbrace{\beta \cdot \|\phi(s, a)\|_{\Lambda_h^{-1}}}_{\text{ellipsoidal penalty}} \right]_{H-h}^+.$$

**Lin-UCBVI vs. PEVI.**

	Lin-UCBVI (Lecture 11)	PEVI (this lecture)
<b>Data</b>	Adaptive online trajectories	Fixed offline dataset
<b>Regression target</b>	$r_h + \widehat{V}_{h+1}(s')$	$r_h + \widehat{V}_{h+1}(s')$
<b>Correction</b>	$+\beta \ \phi\ _{(\Lambda_h^k)^{-1}}$	$-\beta \ \phi\ _{\Lambda_h^{-1}}$
<b>Effect</b>	Encourages exploration	Avoids uncertain regions
<b>Performance term</b>	Bonuses along $\pi^k$	Penalties along $\pi^*$

**Main Result**

**Theorem 2** (PEVI suboptimality for linear MDP). *Suppose Assumption 1 holds and the underlying MDP is a linear MDP (Definition 2). In Algorithm 1, set*

$$\lambda = 1, \quad \beta = c \cdot dH \sqrt{\zeta}, \quad \zeta = \log(2dHK/\xi).$$

Then with probability at least  $1 - \xi$ , for any  $s_0 \in \mathcal{S}$ ,

$$\text{SubOpt}(\text{Pess}(\mathcal{D}); s_0) \leq 2\beta \sum_{h=0}^{H-1} \mathbb{E}_{\pi^*} \left[ \left( \phi(s_h, a_h)^\top \Lambda_h^{-1} \phi(s_h, a_h) \right)^{1/2} \mid s_0 \right]. \quad (6)$$

**PEVI suboptimality (linear MDP):**

$$\text{SubOpt}(\text{Pess}(\mathcal{D}); s_0) \leq \widetilde{O}(dH) \cdot \sum_{h=0}^{H-1} \mathbb{E}_{\pi^*} \left[ \|\phi(s_h, a_h)\|_{\Lambda_h^{-1}} \mid s_0 \right].$$

The bound depends on the penalty evaluated along  $\pi^*$ 's trajectory, with the design matrix  $\Lambda_h$  determined by the dataset.

## Analysis of Linear PEVI

As in Lecture 11, the key step is a confidence bound for the linear Bellman backup. The only extra issue is that the target  $\widehat{V}_{h+1}$  is fitted from the same offline dataset, so the usual fixed-target ridge bound must be uniformized over the PEVI iterate class.

**Proposition 3** (Fixed-target model confidence for linear PEVI). *Fix  $h \in \{0, \dots, H-1\}$  and a deterministic stage- $(h+1)$  Q-function  $Q$  with  $V_Q(\cdot) \in [0, H-h-1]$ . Under the conditions of Theorem 2, with probability at least  $1 - \xi$ , simultaneously for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,*

$$|(\widehat{\mathcal{T}}_h Q)(s, a) - (\mathcal{T}_h Q)(s, a)| \leq \beta \cdot \|\phi(s, a)\|_{\Lambda_h^{-1}}.$$

*Proof sketch.* Fix  $h$  and first fix a deterministic stage- $(h+1)$  Q-function  $Q$ . Write  $V := V_Q$  and  $\phi_\tau := \phi(s_h^\tau, a_h^\tau)$ . By Bellman closure, there exists  $w_{h,Q}^* \in \mathbb{R}^d$  such that

$$r_h(s, a) + (P_h^* V_Q)(s, a) = \langle \phi(s, a), w_{h,Q}^* \rangle, \quad \|w_{h,Q}^*\|_2 \leq H\sqrt{d}.$$

Also define

$$\eta_\tau(Q) := V_Q(s_{h+1}^\tau) - (P_h^* V_Q)(s_h^\tau, a_h^\tau).$$

Then

$$r_h^\tau + V_Q(s_{h+1}^\tau) = r_h(s_h^\tau, a_h^\tau) + (P_h^* V_Q)(s_h^\tau, a_h^\tau) + \eta_\tau(Q),$$

where  $\{\phi_\tau \eta_\tau(Q)\}_{\tau=0}^{K-1}$  is a martingale difference sequence and  $|\eta_\tau(Q)| \leq H$ . Since  $\Lambda_h = \lambda I + \sum_{\tau=0}^{K-1} \phi_\tau \phi_\tau^\top$ ,

$$\begin{aligned} & (\widehat{\mathcal{T}}_h Q)(s, a) - (\mathcal{T}_h Q)(s, a) \\ &= \phi(s, a)^\top \Lambda_h^{-1} \sum_{\tau=0}^{K-1} \phi_\tau \eta_\tau(Q) - \lambda \phi(s, a)^\top \Lambda_h^{-1} w_{h,Q}^*. \end{aligned}$$

Hence, for every fixed  $Q$ ,

$$\begin{aligned} & |(\widehat{\mathcal{T}}_h Q)(s, a) - (\mathcal{T}_h Q)(s, a)| \\ & \leq \|\phi(s, a)\|_{\Lambda_h^{-1}} \left\| \sum_{\tau=0}^{K-1} \phi_\tau \eta_\tau(Q) \right\|_{\Lambda_h^{-1}} + \lambda |\phi(s, a)^\top \Lambda_h^{-1} w_{h,Q}^*|. \end{aligned}$$

By a self-normalized bound ([Abbasi-Yadkori et al., 2011](#)),

$$\left\| \sum_{\tau=0}^{K-1} \phi_\tau \eta_\tau(Q) \right\|_{\Lambda_h^{-1}} = O(H\sqrt{d\zeta}),$$

and since  $\Lambda_h \succeq \lambda I$ ,

$$\lambda |\phi(s, a)^\top \Lambda_h^{-1} w_{h,Q}^*| \leq \sqrt{\lambda} \|w_{h,Q}^*\|_2 \cdot \|\phi(s, a)\|_{\Lambda_h^{-1}} \leq H\sqrt{\lambda d} \|\phi(s, a)\|_{\Lambda_h^{-1}}.$$

Therefore,

$$\begin{aligned} & |(\widehat{\mathcal{T}}_h Q)(s, a) - (\mathcal{T}_h Q)(s, a)| \\ & \leq O(H\sqrt{d\zeta} + H\sqrt{\lambda d}) \cdot \|\phi(s, a)\|_{\Lambda_h^{-1}}. \end{aligned}$$

This proves the fixed-target bound.  $\square$

**Corollary 4** (Confidence for the PEVI iterate). *Under the conditions of Theorem 2, with probability at least  $1 - \xi$ , simultaneously for all  $h \in \{0, \dots, H - 1\}$  and all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,*

$$|(\widehat{\mathcal{T}}_h \widehat{Q}_{h+1})(s, a) - (\mathcal{T}_h \widehat{Q}_{h+1})(s, a)| \leq \beta \|\phi(s, a)\|_{\Lambda_h^{-1}}.$$

*External ingredient.* Proposition 3 gives the bound for every fixed deterministic  $Q$ . To instantiate it at the data-dependent PEVI iterate  $\widehat{Q}_{h+1}$ , Jin et al. (2021) uniformize over an  $\ell_\infty$ -cover of the PEVI iterate class. The log-covering number is  $\widetilde{O}(d^2)$ , which is absorbed by the choice  $\beta = c d H \sqrt{\zeta}$ .  $\square$

**Lemma 5** (Pointwise pessimism of  $\widehat{Q}$ ). *On the event of Corollary 4, simultaneously for all  $h \in \{0, \dots, H - 1\}$  and all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,*

$$(\mathcal{T}_h \widehat{Q}_{h+1})(s, a) - 2\Gamma_h(s, a) \leq \widehat{Q}_h(s, a) \leq (\mathcal{T}_h \widehat{Q}_{h+1})(s, a).$$

*Equivalently,*

$$0 \leq (\mathcal{T}_h \widehat{Q}_{h+1})(s, a) - \widehat{Q}_h(s, a) \leq 2\Gamma_h(s, a) = 2\beta \|\phi(s, a)\|_{\Lambda_h^{-1}}.$$

*Proof.*

$$\widehat{Q}_h(s, a) = \min \left\{ (\widehat{\mathcal{T}}_h \widehat{Q}_{h+1})(s, a) - \Gamma_h(s, a), H - h \right\}^+.$$

By construction,  $\widehat{V}_{h+1}(\cdot) \in [0, H - h - 1]$ , so

$$0 \leq (\mathcal{T}_h \widehat{Q}_{h+1})(s, a) = r_h(s, a) + P_h^* \widehat{V}_{h+1}(s, a) \leq H - h.$$

On the event of Corollary 4,

$$(\mathcal{T}_h \widehat{Q}_{h+1})(s, a) - 2\Gamma_h(s, a) \leq (\widehat{\mathcal{T}}_h \widehat{Q}_{h+1})(s, a) - \Gamma_h(s, a) \leq (\mathcal{T}_h \widehat{Q}_{h+1})(s, a) \leq H - h.$$

Therefore the upper truncation is inactive, and

$$\widehat{Q}_h(s, a) = \max\left\{(\widehat{\mathcal{T}}_h \widehat{Q}_{h+1})(s, a) - \Gamma_h(s, a), 0\right\}.$$

Using the display above together with  $0 \leq (\mathcal{T}_h \widehat{Q}_{h+1})(s, a)$ , we obtain

$$\widehat{Q}_h(s, a) \leq (\mathcal{T}_h \widehat{Q}_{h+1})(s, a),$$

because  $\max\{u, 0\} \leq x$  whenever  $u \leq x$  and  $0 \leq x$ . Likewise,

$$\widehat{Q}_h(s, a) \geq (\widehat{\mathcal{T}}_h \widehat{Q}_{h+1})(s, a) - \Gamma_h(s, a) \geq (\mathcal{T}_h \widehat{Q}_{h+1})(s, a) - 2\Gamma_h(s, a),$$

since  $\max\{u, 0\} \geq u$ . This proves the sandwich bound, and the equivalent Bellman-error bound follows by rearranging.  $\square$

*Proof of Theorem 2.* Let  $\widehat{\pi} = \text{Pess}(\mathcal{D})$ , and let  $E$  denote the event of Corollary 4. On  $E$ , Lemma 5 yields pointwise control of the Bellman error:

$$0 \leq (\mathcal{T}_h \widehat{Q}_{h+1})(s, a) - \widehat{Q}_h(s, a) \leq 2\beta \|\phi(s, a)\|_{\Lambda_h^{-1}}.$$

Now fix the dataset on the event  $E$ , so that  $\widehat{Q}$  and  $\widehat{\pi}$  are deterministic functions and Lemma 1 applies directly. Lemma 1 gives

$$\begin{aligned} \text{SubOpt}(\widehat{\pi}; s_0) &= \sum_{h=0}^{H-1} \mathbb{E}_{\widehat{\pi}} \left[ \widehat{Q}_h(s_h, a_h) - (\mathcal{T}_h \widehat{Q}_{h+1})(s_h, a_h) \mid s_0 \right] \\ &\quad + \sum_{h=0}^{H-1} \mathbb{E}_{\pi^*} \left[ (\mathcal{T}_h \widehat{Q}_{h+1})(s_h, a_h) - \widehat{Q}_h(s_h, a_h) \mid s_0 \right] \\ &\quad + \sum_{h=0}^{H-1} \mathbb{E}_{\pi^*} \left[ \widehat{Q}_h(s_h, \pi_h^*(s_h)) - \widehat{Q}_h(s_h, \widehat{\pi}_h(s_h)) \mid s_0 \right]. \end{aligned}$$

We now bound the three terms on the right-hand side:

- Since  $\widehat{\pi}_h(s) \in \arg\max_a \widehat{Q}_h(s, a)$ , we have

$$\widehat{Q}_h(s, \pi_h^*(s)) - \widehat{Q}_h(s, \widehat{\pi}_h(s)) \leq 0 \quad \text{for every } (h, s),$$

so the optimization error term (iii) is non-positive.

- Since  $\widehat{Q}_h(s, a) - (\mathcal{T}_h \widehat{Q}_{h+1})(s, a) \leq 0$  pointwise on  $E$ , the pessimism-controlled term

$$\sum_{h=0}^{H-1} \mathbb{E}_{\widehat{\pi}} \left[ \widehat{Q}_h(s_h, a_h) - (\mathcal{T}_h \widehat{Q}_{h+1})(s_h, a_h) \mid s_0 \right]$$

is also non-positive.

- Finally, using the upper bound

$$(\mathcal{T}_h \widehat{Q}_{h+1})(s, a) - \widehat{Q}_h(s, a) \leq 2\beta \|\phi(s, a)\|_{\Lambda_h^{-1}}$$

pointwise on  $E$  and taking expectation along  $\pi^*$  yields

$$\begin{aligned} \sum_{h=0}^{H-1} \mathbb{E}_{\pi^*} \left[ (\mathcal{T}_h \widehat{Q}_{h+1})(s_h, a_h) - \widehat{Q}_h(s_h, a_h) \mid s_0 \right] &\leq \sum_{h=0}^{H-1} \mathbb{E}_{\pi^*} \left[ 2\beta \|\phi(s_h, a_h)\|_{\Lambda_h^{-1}} \mid s_0 \right] \\ &= 2\beta \sum_{h=0}^{H-1} \mathbb{E}_{\pi^*} \left[ \|\phi(s_h, a_h)\|_{\Lambda_h^{-1}} \mid s_0 \right]. \end{aligned}$$

Combining the three displays above, we obtain

$$\text{SubOpt}(\widehat{\pi}; s_0) \leq 2\beta \sum_{h=0}^{H-1} \mathbb{E}_{\pi^*} \left[ \|\phi(s_h, a_h)\|_{\Lambda_h^{-1}} \mid s_0 \right].$$

Since  $\Pr(E) \geq 1 - \xi$  by Corollary 4, the theorem follows.  $\square$

#### How linear PEVI works.

1. The linear MDP structure turns the Bellman backup into a linear regression problem in the same feature map  $\phi$  used by Lin-UCBVI.
2. The ellipsoidal penalty  $\beta \|\phi\|_{\Lambda_h^{-1}}$  makes the estimated Q-function a *lower* bound on the Bellman backup, so the pessimism-controlled term in Lemma 1 has the favorable sign.
3. The remaining suboptimality is paid only along  $\pi^*$ 's trajectory. This is the offline analog of “paying for the bonus” in Lecture 11, with the comparator trajectory switching from the learned policy to the optimal policy.

## The Oracle Property

A remarkable feature of the bound (6) is that the suboptimality only depends on how well the dataset  $\mathcal{D}$  covers the trajectory of  $\pi^*$ . We call this the *oracle property*: the algorithm automatically adapts to the optimal policy, even though  $\pi^*$  is unknown.

**Illustration via tabular MDP.** Consider the tabular MDP as a special case of the linear MDP with  $\phi(s, a) = e_{(s,a)} \in \mathbb{R}^{SA}$  (the one-hot encoding). Then

$$\Lambda_h = \lambda I + \sum_{\tau} \phi_{\tau} \phi_{\tau}^{\top} = \text{diag}(\lambda + N_h(s, a) : (s, a) \in \mathcal{S} \times \mathcal{A}),$$

where  $N_h(s, a) = \sum_{\tau} \mathbf{1}\{(s_h^{\tau}, a_h^{\tau}) = (s, a)\}$  is the number of visits to  $(s, a)$  at step  $h$ . The penalty becomes

$$\Gamma_h(s, a) = \frac{\beta}{\sqrt{\lambda + N_h(s, a)}},$$

and the bound (6) specializes to

$$\text{SubOpt} \leq 2\beta \sum_{h=0}^{H-1} \mathbb{E}_{\pi^*} \left[ \frac{1}{\sqrt{\lambda + N_h(s_h, a_h)}} \right].$$

The suboptimality only depends on the counts  $N_h(s_h^*, a_h^*)$  along the optimal trajectory. State-action pairs  $(s, a)$  that  $\pi^*$  never visits contribute zero to the bound, regardless of how often (or rarely) they appear in the dataset.

**Oracle property.** Suppose the transition is deterministic, so  $\pi^*$  follows a unique trajectory  $\{(s_h^*, a_h^*)\}_{h=0}^{H-1}$ . Then

$$\text{SubOpt} \leq \sum_{h=0}^{H-1} \frac{2\beta}{\sqrt{\lambda + N_h(s_h^*, a_h^*)}}.$$

This depends only on how many times the dataset visits the states on  $\pi^*$ 's path. A behavior policy  $\pi^b$  that visits many “wrong” states does not hurt; and a behavior policy that focuses on the optimal trajectory leads to low suboptimality. Meanwhile, a state-action pair  $(s^{\circ}, a^{\circ})$  off  $\pi^*$ 's path—even one visited millions of times—does not affect the bound at all.

This is analogous to *local* sample complexity: the difficulty of learning  $\pi^*$  depends on the *local* coverage near  $\pi^*$ , not on global coverage.

## Sufficient Coverage and Sample Complexity

To translate the data-dependent bound (6) into a sample complexity guarantee, we need to quantify how well the dataset covers  $\pi^*$ 's trajectory.

**Corollary 6** (Well-explored dataset). *Suppose there exists  $c^{\dagger} > 0$  such that the event*

$$\mathcal{E}^{\dagger} = \{ \Lambda_h \succeq I + c^{\dagger} \cdot K \cdot \mathbb{E}_{\pi^*} [\phi(s_h, a_h) \phi(s_h, a_h)^{\top} \mid s_0] \text{ for all } h \}$$

*holds with probability at least  $1 - \xi/2$ . Set  $\lambda = 1$ ,  $\beta = c \cdot dH\sqrt{\zeta}$ ,  $\zeta = \log(4dHK/\xi)$ . Then with*

probability at least  $1 - \xi$ ,

$$\text{SubOpt}(\text{Pess}(\mathcal{D}); s_0) \leq c' \cdot d^{3/2} H^2 K^{-1/2} \sqrt{\zeta},$$

where  $c'$  depends only on  $c^\dagger$  and  $c$ .

Under stronger conditions on the behavior policy (e.g.,  $\lambda_{\min}(\Sigma_h) \geq \underline{c}/d$  for all  $h$ , where  $\Sigma_h = \mathbb{E}_{\pi^b}[\phi(s_h, a_h)\phi(s_h, a_h)^\top]$ ), Corollary 4.6 of Jin et al. (2021) shows that  $K = \tilde{O}(d^3 H^4 / \varepsilon^2)$  i.i.d. trajectories suffice for  $\varepsilon$ -suboptimality.

## Information-Theoretic Lower Bound

PEVI's suboptimality bound is *minimax optimal* in the linear MDP setting, as established by the following information-theoretic lower bound. Recall that  $P_{\mathcal{D}}$  denotes the joint distribution of the data collecting process.

**Theorem 7** (Information-theoretic lower bound). *There exists an absolute constant  $c > 0$  such that for the output  $\text{Algo}(\mathcal{D})$  of any algorithm, there exist a linear MDP  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, H, P^*, r)$ , an initial state  $s_0$ , and a compliant dataset  $\mathcal{D}$  such that*

$$\mathbb{E}_{\mathcal{D} \sim P_{\mathcal{D}}} \left[ \frac{\text{SubOpt}(\text{Algo}(\mathcal{D}); s_0)}{\sum_{h=0}^{H-1} \mathbb{E}_{\pi^*} [\|\phi(s_h, a_h)\|_{\Lambda_h^{-1}} \mid s_0]} \right] \geq c.$$

**Minimax optimality of PEVI:** Up to the scaling parameter  $\beta = \tilde{O}(dH)$ , the suboptimality of PEVI (Theorem 2) matches the lower bound (Theorem 7). In other words, the data-dependent quantity  $\sum_h \mathbb{E}_{\pi^*} [\|\phi\|_{\Lambda_h^{-1}}]$  is the *fundamental statistical difficulty* of the offline RL problem in linear MDPs.

## Beyond Pointwise Pessimism: Bellman-Consistent Pessimism

The linear PEVI analysis above relies on a *pointwise* penalty  $\Gamma_h(s, a)$  such that  $\hat{Q}_h \leq \mathcal{T}_h \hat{Q}_{h+1}$  at every state-action pair. In structured settings such as tabular and linear MDPs, these penalties can be written explicitly. For general function approximation, however, no comparable closed-form pointwise penalty is available.

This mirrors online RL (Lecture 11): pointwise bonuses work for tabular (UCBVI) and linear (Lin-UCBVI) MDPs, but GOLF uses *global* optimism for general function classes. We now develop the offline analog: *Bellman-consistent pessimism* (BCP) (Xie et al., 2021).

**Pointwise vs. global pessimism (cf. Lecture 11).**

	Pointwise pessimism (PEVI)	Global pessimism (BCP)
<b>Goal</b>	$\widehat{Q}_h \leq (\mathcal{T}_h \widehat{Q}_{h+1})$ for all $(s, a, h)$	Pessimistic evaluation at $s_0$
<b>Mechanism</b>	Subtract penalty $\Gamma_h(s, a)$	Version space + min
<b>Scope</b>	Requires closed-form penalty	General $\mathcal{F}$
<b>Online analog</b>	UCBVI / Lin-UCBVI	GOLF

**Setup**

We now switch notation and follow the discounted infinite-horizon framework of [Xie et al. \(2021\)](#): an MDP  $(\mathcal{S}, \mathcal{A}, P, R, \gamma, s_0)$  with discount factor  $\gamma \in [0, 1)$ , transition  $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ , reward  $R : \mathcal{S} \times \mathcal{A} \rightarrow [0, R_{\max}]$ , and value bound  $V_{\max} = R_{\max}/(1 - \gamma)$ . We write  $J(\pi) = V^\pi(s_0)$  for the performance of policy  $\pi$ , and  $f(s_0, \pi) := \sum_a f(s_0, a) \pi(a | s_0)$  for a function  $f$  evaluated under a (possibly stochastic) policy. The policy-specific Bellman operator and transition operator are

$$(\mathcal{T}^\pi f)(s, a) = R(s, a) + \gamma (\mathcal{P}^\pi f)(s, a), \quad (\mathcal{P}^\pi f)(s, a) = \mathbb{E}_{s' \sim P(\cdot | s, a)} [f(s', \pi(s'))].$$

**Discounted Bellman equation (used in the BCP section):**

$$Q^\pi = \mathcal{T}^\pi Q^\pi, \quad \text{equivalently} \quad Q^\pi(s, a) = R(s, a) + \gamma (\mathcal{P}^\pi Q^\pi)(s, a).$$

The *discounted state-action occupancy* is  $d_\pi(s, a) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \Pr(s_t = s, a_t = a | s_0, \pi)$ .

The offline dataset  $\mathcal{D} = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^n$  consists of  $n$  i.i.d. tuples:  $(s, a) \sim \mu$ ,  $r = R(s, a)$ ,  $s' \sim P(\cdot | s, a)$ , for a *data distribution*  $\mu$  over  $\mathcal{S} \times \mathcal{A}$ . We write  $\|g\|_{2, \nu}^2 = \mathbb{E}_\nu [g(s, a)^2]$  for the  $\nu$ -weighted  $L^2$  norm.

We are given a finite value-function class  $\mathcal{F} \subset (\mathcal{S} \times \mathcal{A} \rightarrow [0, V_{\max}])$  and policy class  $\Pi$ . For clarity, we present the theory under exact realizability and completeness; the full result in [Xie et al. \(2021\)](#) handles approximate versions of both assumptions (with errors  $\varepsilon_{\mathcal{F}}, \varepsilon_{\mathcal{F}, \mathcal{F}} \geq 0$ ) and decomposes the bound into on-support and off-support errors that automatically adapt to the best bias-variance tradeoff.

**Assumption 2 (Realizability).**  $Q^\pi \in \mathcal{F}$  for all  $\pi \in \Pi$ .

**Assumption 3 (Completeness).**  $\mathcal{T}^\pi f \in \mathcal{F}$  for all  $\pi \in \Pi$  and  $f \in \mathcal{F}$ .

**Definition 3** (Distribution shift coefficient).

$$\mathcal{C}(\nu; \mu, \mathcal{F}, \pi) := \max_{f \in \mathcal{F}} \frac{(\mathbb{E}_\nu[|f - \mathcal{T}^\pi f|])^2}{\mathbb{E}_\mu[(f - \mathcal{T}^\pi f)^2]}.$$

This is the *performance-relevant* shift coefficient for our proof: it measures how well squared Bellman error under  $\mu$  (data) controls *average* Bellman error under  $\nu$  (target). When  $\mathcal{C}(d_\pi; \mu, \mathcal{F}, \pi) \leq C_2$ , any function with small Bellman error on the data also has small average Bellman error along  $\pi$ 's trajectory.

**Remark 1** (Relation to density ratio and covariance shift). *This is the weakest coefficient needed here, because Lemma 11 reduces performance to the average Bellman error  $\mathbb{E}_{d_\pi}[|f - \mathcal{T}^\pi f|]$ . If the Bellman-error class  $\mathcal{G}_{\mathcal{F}, \pi} := \{f - \mathcal{T}^\pi f : f \in \mathcal{F}\}$  is linear in some feature map  $\psi$ , say  $g = \langle \psi, w \rangle$ , then*

$$\mathcal{C}(\nu; \mu, \mathcal{F}, \pi) \leq \sup_{w \neq 0} \frac{(\mathbb{E}_\nu[|\langle \psi, w \rangle|])^2}{w^\top \Sigma_\mu w}, \quad \Sigma_\mu = \mathbb{E}_\mu[\psi \psi^\top].$$

Indeed, each  $f \in \mathcal{F}$  induces some  $g = f - \mathcal{T}^\pi f \in \mathcal{G}_{\mathcal{F}, \pi}$ , and by linearity  $g = \langle \psi, w \rangle$  for some  $w$ . Moreover,

$$\mathbb{E}_\mu[g^2] = \mathbb{E}_\mu[\langle \psi, w \rangle^2] = w^\top \Sigma_\mu w.$$

Substituting this into Definition 3 gives the displayed upper bound. Moreover,

$$\begin{aligned} \sup_{w \neq 0} \frac{(\mathbb{E}_\nu[|\langle \psi, w \rangle|])^2}{w^\top \Sigma_\mu w} &= \sup_{w \neq 0} \left( \mathbb{E}_\nu \left[ \frac{|\langle \psi, w \rangle|}{\sqrt{w^\top \Sigma_\mu w}} \right] \right)^2 \\ &= \sup_{w \neq 0} \left( \mathbb{E}_\nu \left[ \left| \left\langle \Sigma_\mu^{-1/2} \psi, \frac{\Sigma_\mu^{1/2} w}{\sqrt{w^\top \Sigma_\mu w}} \right\rangle \right| \right] \right)^2 \\ &\leq \sup_{w \neq 0} \left( \mathbb{E}_\nu \left[ \|\psi\|_{\Sigma_\mu^{-1}} \cdot \left\| \frac{\Sigma_\mu^{1/2} w}{\sqrt{w^\top \Sigma_\mu w}} \right\|_2 \right] \right)^2 \\ &= \sup_{w \neq 0} \left( \mathbb{E}_\nu [\|\psi\|_{\Sigma_\mu^{-1}}] \right)^2 \\ &= \left( \mathbb{E}_\nu [\|\psi\|_{\Sigma_\mu^{-1}}] \right)^2, \end{aligned}$$

Here the inequality is Cauchy–Schwarz, and the penultimate equality uses that the second factor has unit norm. Thus this current form is weaker than the PEVI-style average local uncertainty. A simpler sufficient condition is

$$\Sigma_\nu \preceq C \Sigma_\mu, \quad \Sigma_\nu = \mathbb{E}_\nu[\psi \psi^\top], \quad \Sigma_\mu = \mathbb{E}_\mu[\psi \psi^\top].$$

Indeed, for any  $w$ ,

$$(\mathbb{E}_\nu[|\langle \psi, w \rangle|])^2 \leq \mathbb{E}_\nu[\langle \psi, w \rangle^2] = w^\top \Sigma_\nu w \leq C w^\top \Sigma_\mu w.$$

Finally, the raw density ratio bound  $\|\nu/\mu\|_\infty \leq C$  is an even stronger pointwise sufficient condition.

## Algorithm: Pessimistic Policy Selection

**Bellman consistency score.** For  $f, f' \in \mathcal{F}$  and  $\pi \in \Pi$ , define

$$\mathcal{L}(f', f, \pi; \mathcal{D}) = \frac{1}{n} \sum_{i=1}^n (f'(s_i, a_i) - r_i - \gamma f(s'_i, \pi(s'_i)))^2,$$

and the *Bellman consistency score*:

$$\mathcal{E}(f, \pi; \mathcal{D}) := \mathcal{L}(f, f, \pi; \mathcal{D}) - \min_{f' \in \mathcal{F}} \mathcal{L}(f', f, \pi; \mathcal{D}). \quad (7)$$

**Version space and pessimistic selection.** The version space is  $\mathcal{F}_{\pi, \varepsilon} := \{f \in \mathcal{F} : \mathcal{E}(f, \pi; \mathcal{D}) \leq \varepsilon\}$ , and BCP selects

$$\hat{\pi} = \operatorname{argmax}_{\pi \in \Pi} \min_{f \in \mathcal{F}_{\pi, \varepsilon}} f(s_0, \pi). \quad (8)$$

**BCP:** For each  $\pi$ , pessimistically evaluate it as  $\min_{f \in \mathcal{F}_{\pi, \varepsilon}} f(s_0, \pi)$ ; then select the  $\pi$  with the highest pessimistic value.

**BCP vs. GOLF.** BCP is the offline mirror of GOLF (Lecture 11). GOLF selects the most *optimistic*  $f$  in a confidence set; BCP selects the most *pessimistic* evaluation of each policy. Both use Bellman-error-based version sets, replacing pointwise bonuses/penalties with a global selection principle.

## Analysis of BCP

The analysis uses four ingredients: (1) show that the version space contains  $Q^\pi$  (so the pessimistic estimate lower-bounds  $J(\pi)$ ), (2) show that every function in the version space has small Bellman error under  $\mu$ , (3) use a simulation lemma to express  $J(\pi) - f(s_0, \pi)$  as an average Bellman error along  $d_\pi$ , and (4) transfer that Bellman error from  $d_\pi$  back to  $\mu$  via  $\mathcal{C}$ . Throughout, let  $L := \log(|\mathcal{F}||\Pi|/\delta)$ .

### Step 1: Version space containment.

**Lemma 8** (Version space containment). *Under Assumptions 2–3, setting  $\varepsilon = \varepsilon_r := c V_{\max}^2 L/n$  for a sufficiently large constant  $c$ , we have with probability at least  $1 - \delta$ :*

$$Q^\pi \in \mathcal{F}_{\pi, \varepsilon_r} \quad \text{for all } \pi \in \Pi.$$

*Previously established fast-rate ingredient.* This is the same squared-loss Bernstein fast-rate argument already used in Lecture 8 (“Fast Rate via Bernstein’s Inequality”), which was presented there as the policy-evaluation analogue of the Lecture 7 fitted Q-iteration argument. For a fixed policy  $\pi$ , the score  $\mathcal{E}(f, \pi; \mathcal{D})$  is an empirical *excess squared loss*, and  $Q^\pi$  is the population minimizer because  $Q^\pi = \mathcal{T}^\pi Q^\pi$  by realizability. The same  $Y$ -variable trick, self-bounding variance bound, and Bernstein inequality from Lecture 8 therefore yield the fast rate

$$\mathcal{E}(Q^\pi, \pi; \mathcal{D}) \leq c V_{\max}^2 \frac{L}{n}$$

simultaneously for all  $\pi \in \Pi$  with probability at least  $1 - \delta$ . Hence  $Q^\pi \in \mathcal{F}_{\pi, \varepsilon_r}$  for all  $\pi$ .  $\square$

### Step 2: Bellman error control for version space members.

**Lemma 9** (Bellman error bound). *Under Assumptions 2–3, on the same high-probability event, for all  $f \in \mathcal{F}_{\pi, \varepsilon_r}$ :*

$$\mathbb{E}_\mu[(f - \mathcal{T}^\pi f)^2] \leq c_2^2 V_{\max}^2 L/n.$$

*Previously established fast-rate ingredient.* Again, this is the same operator-agnostic squared-loss fast-rate template from Lectures 7 and 8. For any fixed  $f \in \mathcal{F}_{\pi, \varepsilon_r}$ , completeness gives  $g = \mathcal{T}^\pi f \in \mathcal{F}$ , and the version-space condition says that the empirical excess loss of  $f$  relative to the Bellman target  $g$  is at most  $\varepsilon_r$ . Applying the same self-bounding Bernstein argument as in Lecture 8 to this excess loss yields  $\mathbb{E}_\mu[(f - \mathcal{T}^\pi f)^2] \leq c_2^2 V_{\max}^2 L/n$ .  $\square$

**Step 3: Bellman error to decision loss.** This is exactly the Performance Difference Identity from Lecture 6, specialized to the policy-specific operator  $\mathcal{T}^\pi$  and the point-mass initial distribution at  $s_0$ .

**Lemma 10** (Bellman error to decision loss). *For any  $g : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  and any policy  $\pi$ :*

$$g(s_0, \pi) - J(\pi) = \frac{1}{1 - \gamma} \mathbb{E}_{d_\pi}[g(s, a) - (\mathcal{T}^\pi g)(s, a)].$$

*Proof.* Since  $(\mathcal{T}^\pi g)(s, a) = R(s, a) + \gamma (\mathcal{P}^\pi g)(s, a)$ , we have

$$\frac{1}{1 - \gamma} \mathbb{E}_{d_\pi}[g - \mathcal{T}^\pi g] = \frac{1}{1 - \gamma} \mathbb{E}_{d_\pi}[g - R - \gamma \mathcal{P}^\pi g].$$

Using  $d_\pi(s, a) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t d_\pi^{(t)}(s, a)$  where  $d_\pi^{(t)}$  is the step- $t$  occupancy, this becomes

$$\sum_{t=0}^{\infty} \gamma^t (\mathbb{E}_{d_\pi^{(t)}}[g] - \mathbb{E}_{d_\pi^{(t)}}[R] - \gamma \mathbb{E}_{d_\pi^{(t+1)}}[g]).$$

The  $g$ -terms telescope to  $\mathbb{E}_{d_\pi^{(0)}}[g] = g(s_0, \pi)$ , while the reward terms sum to  $J(\pi)$ . Hence

$$\frac{1}{1 - \gamma} \mathbb{E}_{d_\pi}[g - \mathcal{T}^\pi g] = g(s_0, \pi) - J(\pi).$$

□

**Step 4: From pessimism to Bellman error.** Compared with the Bellman-error-to-decision-loss lemmas from Lectures 6 and 7, BCP only needs the comparator-side Bellman error: pessimistic selection already replaces the learned-policy-side term.

For a comparator policy  $\pi$ , let

$$f_{\pi, \min} := \operatorname{argmin}_{f \in \mathcal{F}_{\pi, \varepsilon_r}} f(s_0, \pi),$$

which is well-defined on the event of Lemma 8 because  $Q^\pi \in \mathcal{F}_{\pi, \varepsilon_r}$ .

**Lemma 11** (Pessimism reduces suboptimality to Bellman error). *On the event of Lemma 8, for any comparator  $\pi \in \Pi$ :*

$$J(\pi) - J(\hat{\pi}) \leq \frac{1}{1 - \gamma} \mathbb{E}_{d_\pi}[\mathcal{T}^\pi f_{\pi, \min} - f_{\pi, \min}] \leq \frac{1}{1 - \gamma} \mathbb{E}_{d_\pi}[|f_{\pi, \min} - \mathcal{T}^\pi f_{\pi, \min}|].$$

*Proof.* Since  $Q^{\hat{\pi}} \in \mathcal{F}_{\hat{\pi}, \varepsilon_r}$  (Lemma 8),

$$J(\hat{\pi}) = Q^{\hat{\pi}}(s_0, \hat{\pi}) \geq \min_{f \in \mathcal{F}_{\hat{\pi}, \varepsilon_r}} f(s_0, \hat{\pi}).$$

By the optimality of  $\hat{\pi}$  in (8),

$$\min_{f \in \mathcal{F}_{\hat{\pi}, \varepsilon_r}} f(s_0, \hat{\pi}) \geq \min_{f \in \mathcal{F}_{\pi, \varepsilon_r}} f(s_0, \pi) = f_{\pi, \min}(s_0, \pi).$$

Therefore,

$$J(\pi) - J(\hat{\pi}) \leq Q^\pi(s_0, \pi) - f_{\pi, \min}(s_0, \pi).$$

Applying Lemma 10 with  $g = f_{\pi, \min}$  gives

$$f_{\pi, \min}(s_0, \pi) - J(\pi) = \frac{1}{1 - \gamma} \mathbb{E}_{d_\pi}[f_{\pi, \min} - \mathcal{T}^\pi f_{\pi, \min}].$$

Rearranging yields

$$J(\pi) - f_{\pi, \min}(s_0, \pi) = \frac{1}{1 - \gamma} \mathbb{E}_{d_\pi} [\mathcal{T}^\pi f_{\pi, \min} - f_{\pi, \min}].$$

Combining with the previous display, we obtain

$$J(\pi) - J(\hat{\pi}) \leq \frac{1}{1 - \gamma} \mathbb{E}_{d_\pi} [\mathcal{T}^\pi f_{\pi, \min} - f_{\pi, \min}].$$

The final inequality follows from the triangle inequality.  $\square$

## Main Result

**Theorem 12** (BCP suboptimality). *Under Assumptions 2–3, let  $\hat{\pi}$  be the output of (8) with  $\varepsilon = \varepsilon_r = c V_{\max}^2 L/n$  and  $L = \log(|\mathcal{F}||\Pi|/\delta)$ . For any comparator  $\pi \in \Pi$  with  $\mathcal{C}(d_\pi; \mu, \mathcal{F}, \pi) \leq C_2$ , with probability at least  $1 - \delta$ :*

$$J(\pi) - J(\hat{\pi}) \leq O\left(\frac{V_{\max} \sqrt{C_2}}{1 - \gamma} \sqrt{\frac{L}{n}}\right).$$

*Proof.* On the event of Lemma 8, the function  $f_{\pi, \min}$  belongs to  $\mathcal{F}_{\pi, \varepsilon_r}$ , so Lemma 11 gives

$$J(\pi) - J(\hat{\pi}) \leq \frac{1}{1 - \gamma} \mathbb{E}_{d_\pi} [|f_{\pi, \min} - \mathcal{T}^\pi f_{\pi, \min}|].$$

By the definition of  $\mathcal{C}$ ,

$$\mathbb{E}_{d_\pi} [|f_{\pi, \min} - \mathcal{T}^\pi f_{\pi, \min}|] \leq \sqrt{\mathcal{C}(d_\pi; \mu, \mathcal{F}, \pi) \cdot \mathbb{E}_\mu [(f_{\pi, \min} - \mathcal{T}^\pi f_{\pi, \min})^2]} \leq c_2 \sqrt{C_2} V_{\max} \sqrt{L/n},$$

where the last step uses Lemma 9. Combining the two displays completes the proof.  $\square$

**BCP suboptimality (exact realizability + completeness):**

$$J(\pi) - J(\hat{\pi}) \leq O\left(\frac{V_{\max} \sqrt{C_2}}{1 - \gamma} \sqrt{\frac{\log(|\mathcal{F}||\Pi|/\delta)}{n}}\right) \quad \text{when } \mathcal{C}(d_\pi; \mu, \mathcal{F}, \pi) \leq C_2.$$

**Key features of BCP.**

- **No explicit density ratio.** BCP requires only the *function-class-aware* distribution shift coefficient  $\mathcal{C}$ , not the crude density ratio  $\|d_\pi/\mu\|_\infty$ .
- **Bounded degradation.** Setting  $\pi = \pi_b$  (behavior policy) with  $\mu = d_{\pi_b}$  gives  $\mathcal{C} = 1$ , and  $J(\pi_b) - J(\hat{\pi}) \leq \tilde{O}\left(\frac{V_{\max}}{1-\gamma} \sqrt{1/n}\right)$ : BCP is competitive with the data-collecting policy.
- **Competing with  $\pi^*$ .** Setting  $\pi = \pi^*$  recovers  $J(\pi^*) - J(\hat{\pi}) \leq \tilde{O}\left(\frac{V_{\max}\sqrt{C_2}}{1-\gamma} \sqrt{1/n}\right)$  under  $\mathcal{C}(d_{\pi^*}; \mu, \mathcal{F}, \pi^*) \leq C_2$ .
- **Adaptive bias-variance tradeoff.** The general Theorem 3.1 in [Xie et al. \(2021\)](#) allows approximate realizability ( $\varepsilon_{\mathcal{F}} > 0$ ) and completeness ( $\varepsilon_{\mathcal{F}, \mathcal{F}} > 0$ ), and decomposes the bound into on-support and off-support errors; the algorithm automatically adapts to the best decomposition.

**A Remark on Computation**

The constrained BCP selection (8) is stated as  $\arg \max_{\pi \in \Pi} \min_{f \in \mathcal{F}_{\pi, \varepsilon}} f(s_0, \pi)$ , which is *information-theoretic* in nature—a direct implementation would search over  $\Pi$  and, for each  $\pi$ , solve a constrained minimization over  $\mathcal{F}$ . It is worth contrasting this with the online setting.

**Online vs. offline computation.** The main contrast is not that online RL is impossible while offline RL is easy. Rather, in our current theory line, online methods with general function approximation are primarily *value-based* and centered on the Bellman *optimality* operator  $\mathcal{T}$ . As discussed in Lecture 8, the  $\max_{a' \in \mathcal{A}}$  inside  $\mathcal{T}$  makes Bellman-error minimization a moving-target, non-convex problem. By contrast, Bellman-consistent pessimism is naturally paired with more *policy-based* relaxations: for a fixed policy  $\pi$ , the critic uses the policy-specific operator  $\mathcal{T}^\pi$ , so the inner Bellman-consistency objective is a regression problem rather than an optimization involving the  $\max_a$  nonlinearity.

On top of this critic, one can then use standard policy-optimization primitives:

- **PSPI** (Section 4 of [Xie et al., 2021](#)) replaces (8) by a regularized objective and updates the actor via policy mirror descent / soft policy iteration (cf. Lecture 10).
- **ATAC** ([Cheng et al., 2022](#)) is a later actor-critic style implementation of offline pessimism based on adversarially trained critics.

This computational gap is further amplified by exploration: online optimism must couple value estimation with data collection (cf. GOLF in Lecture 11), whereas offline pessimism works with a fixed dataset and can be implemented through a regression-plus-policy-optimization loop.

The sample-complexity rates currently achieved by these implementable algorithms are not known to match the information-theoretic rate of (8), so we state no specific bound here; see the respective papers for the precise guarantees.

## Summary

### Summary of offline RL guarantees:

Algorithm	Suboptimality	Key Technique
Naive greedy	Can be $\Omega(H)$ (term (i) uncontrolled)	No pessimism
PEVI (linear)	$\tilde{O}(dH) \cdot \sum_h \mathbb{E}_{\pi^*} [\ \phi\ _{\Lambda^{-1}}]$	Ellipsoidal penalty
BCP (general FA)	$O\left(\frac{V_{\max} \sqrt{C_2}}{1-\gamma} \sqrt{\frac{\log  \mathcal{F}   \Pi }{n}}\right)$	Version space
Lower bound	$\Omega\left(\sum_h \mathbb{E}_{\pi^*} [\ \phi\ _{\Lambda^{-1}}]\right)$	Info-theoretic

**Online vs. Offline RL:**

	Online RL (Lec 11)	Offline RL (this lecture)
<b>Data</b>	Adaptive (agent explores)	Fixed dataset (no interaction)
<b>Principle</b>	Optimism	Pessimism
<b>Q-update</b>	$\hat{Q} = \hat{\mathcal{T}}\hat{V} + \text{bonus}$	$\hat{Q} = \hat{\mathcal{T}}\hat{V} - \text{penalty}$
<b>Effect</b>	Explores uncertain regions	Avoids uncertain regions
<b>Perf. depends on</b>	Bonuses along <i>learned</i> policy	Penalties along <i>optimal</i> policy
<b>General FA</b>	GOLF: global optimism	BCP: global pessimism

**Key takeaways.**

- Offline RL is fundamentally different from online RL: without the ability to collect new data, the agent cannot explore, and *optimism* can be catastrophically wrong.
- The *pessimism principle* resolves this by penalizing uncertain state-action pairs, ensuring the learned policy avoids regions where the dataset provides insufficient coverage.
- The suboptimality decomposition (Lemma 1) reveals three sources of error: the pessimism-controlled term, the optimism-controlled term, and the optimization error. Pessimism eliminates the pessimism-controlled term, leaving only the optimism-controlled term.
- PEVI achieves *minimax optimal* suboptimality for linear MDPs, with a data-dependent bound that exhibits the *oracle property*: the suboptimality depends only on coverage of the *optimal* policy's trajectory.
- Bellman-consistent pessimism (BCP) extends the pessimism principle to general function approximation via version spaces, paralleling the GOLF algorithm's extension of optimism in the online setting.
- The online-offline duality is remarkably clean: the same mathematical objects (bonus = penalty, confidence set = version space, optimism = pessimism) appear in both settings with only a *sign change*. Understanding one setting deeply illuminates the other.

**References**

Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, volume 24, 2011.

Ching-An Cheng, Tengyang Xie, Nan Jiang, and Alekh Agarwal. Adversarially trained

actor critic for offline reinforcement learning. In *International Conference on Machine Learning*, 2022.

Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline RL? In *International Conference on Machine Learning*, pages 5084–5096, 2021.

Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. Bellman-consistent pessimism for offline reinforcement learning. In *Advances in Neural Information Processing Systems*, volume 34, 2021.