Is Pessimism Provably Efficient for Offline RL?

Ying Jin Stanford Statistics

RL Theory Virtual Seminars April 6, 2020

Joint work with



Zhuoran Yang Princeton ORFE



Zhaoran Wang Northwestern IEMS

Backgrounds

- Episodic MDP
- Offline Learning Protocol
- Difficulty and Epistemic Uncertainty

Episodic MDP



- \blacktriangleright S: infinite state space. A: finite action space.
- Unknown reward function $r_h : S \times A \rightarrow [0, 1]$.
- Unknown transition kernel $\mathbb{P}_h(\cdot | x, a) \in \Delta(\mathcal{S})$.
- Finite horizon H: terminate when h = H.

Episodic MDP



- Policy: $\pi = {\pi_h}_{h \in [H]} : S \to \Delta(A), a_h \sim \pi_h(s_h).$
- Expected total reward: $J(\pi, x) = \mathbb{E}_{\pi}[\sum_{h=1}^{H} r_h | s_1 = x] \in [0, H].$
- Optimal policy: $\pi^*(\cdot) = \operatorname{argmax}_{\pi} J(\pi, \cdot).$

• π^* is greedy w.r.t. optimal value functions $Q^* = \{Q_h^*\}_{h \in [H]}$.

• π^* is greedy w.r.t. optimal value functions $Q^* = \{Q_h^*\}_{h \in [H]}$.

$$\pi_h^{\star}(x) = \operatorname*{argmax}_{a \in \mathcal{A}} Q_h^{\star}(x, a), \quad \forall s \in S,$$
$$Q_h^{\star}(x, a) = \underbrace{r_h(x, a) + \mathbb{E}_{s_{h+1} \sim \mathbb{P}} \left[\max_{a' \in \mathcal{A}} Q_{h+1}^{\star}(s_{h+1}, a') \mid s_h = x, a_h = a\right]}_{\mathbf{A} \sim \mathbb{P}}$$

Bellman operator \mathbb{B}_h : $\mathbb{B}_h Q_{h+1}^{\star}$

• π^* is greedy w.r.t. optimal value functions $Q^* = \{Q_h^*\}_{h \in [H]}$.

$$\pi_{h}^{\star}(x) = \underset{a \in \mathcal{A}}{\operatorname{argmax}} Q_{h}^{\star}(x, a), \quad \forall s \in S,$$

$$Q_{h}^{\star}(x, a) = \underbrace{r_{h}(x, a) + \mathbb{E}_{s_{h+1} \sim \mathbb{P}} \left[\max_{a' \in \mathcal{A}} Q_{h+1}^{\star}(s_{h+1}, a') \mid s_{h} = x, a_{h} = a \right]}_{\operatorname{Bellman operator} \mathbb{B}_{h} \colon \mathbb{B}_{h} Q_{h+1}^{\star}}$$

Bellman Equation

$$Q_h^{\star} = \mathbb{B}_h Q_{h+1}^{\star}, \quad Q_{H+1}^{\star} \equiv 0.$$

• π^* is greedy w.r.t. optimal value functions $Q^* = \{Q_h^*\}_{h \in [H]}$.

$$\pi_{h}^{\star}(x) = \underset{a \in \mathcal{A}}{\operatorname{argmax}} Q_{h}^{\star}(x, a), \quad \forall s \in S,$$

$$Q_{h}^{\star}(x, a) = \underbrace{r_{h}(x, a) + \mathbb{E}_{s_{h+1} \sim \mathbb{P}} \left[\max_{a' \in \mathcal{A}} Q_{h+1}^{\star}(s_{h+1}, a') \mid s_{h} = x, a_{h} = a \right]}_{\operatorname{Bellman operator} \mathbb{B}_{h} \colon \mathbb{B}_{h} Q_{h+1}^{\star}}$$

Bellman Equation

$$Q_h^{\star} = \mathbb{B}_h Q_{h+1}^{\star}, \quad Q_{H+1}^{\star} \equiv 0.$$

RL with function approximation:

Function class $\mathcal{F} = \{f : \mathcal{S} \times \mathcal{A} \to \mathbb{R}\}$ approximates Q_h^{\star} .

Linear functions, Neural networks, RKHS...

Offline Policy Learning Learn from Given Datasets



- Offline Data: collected a priori.
- Arbitrary trajectories: actions a_h by an offline agent (unknown rule).
- ▶ No further interactions with MDP.

Offline Policy Learning Learn from Given Datasets



- Offline Data: collected a priori.
- Arbitrary trajectories: actions a_h by an offline agent (unknown rule).
- No further interactions with MDP.
- Learning objective: performance of the learned policy

$$\mathsf{SubOpt}(\widehat{\pi}, x) = J(\pi^{\star}, x) - J(\widehat{\pi}, x),$$

where $\widehat{\pi} = \text{OfflineRL}(\mathcal{D}, \mathcal{F})$, $x \in \mathcal{S}$.

Offline Policy Learning Naive Value Iterations

By Bellman equation: approximate dynamic programming.

Naive Value Iterations

- End of Episode: $\widehat{Q}_{H+1} \leftarrow 0$.
- Dynamic Programming: $h = H, H 1, \dots, 1$,
 - **Estimate:** $\widehat{Q}_h \leftarrow \operatorname{Regress}(\mathbb{B}_h \widehat{Q}_{h+1}, \mathcal{D}, \mathcal{F}).$
 - **Optimize:** $\widehat{\pi}_h(x) \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} \widehat{Q}_h(x, a).$

Offline Policy Learning Naive Value Iterations

By Bellman equation: approximate dynamic programming.

Naive Value Iterations

- End of Episode: $\widehat{Q}_{H+1} \leftarrow 0$.
- Dynamic Programming: $h = H, H 1, \dots, 1$,
 - **Estimate**: $\widehat{Q}_h \leftarrow \operatorname{Regress}(\mathbb{B}_h \widehat{Q}_{h+1}, \mathcal{D}, \mathcal{F}).$
 - **Optimize:** $\widehat{\pi}_h(x) \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} \widehat{Q}_h(x, a).$

Hope for ...

- Good estimation: $\widehat{Q}_h \approx Q_h^{\star}$.
- Good optimization:

$$\operatorname*{argmax}_{a} \widehat{Q}_{h}(x,a) \approx \operatorname*{argmax}_{a} Q_{h}^{\star}(x,a).$$

• Good actual performance:

 $J(\widehat{\pi}) \approx J(\pi^{\star}).$

Gaussian bandits

$$R_i \mid_{A_i = a} \sim N(r(a), 1), \quad a \in [K] \text{ arms.}$$

▶ *a* is pulled N(a) times in the pre-collected dataset $\mathcal{D} = \{(A_i, R_i)\}$.

Gaussian bandits

$$R_i \mid_{A_i = a} \sim N(r(a), 1), \quad a \in [K] \text{ arms}.$$

a is pulled N(a) times in the pre-collected dataset D = {(A_i, R_i)}.
Following the naive value iterations,

- Estimate by sample mean: $\widehat{\mu}(a) = \frac{\sum_{i:A_i=a} R_i}{N(a)}$.
- Optimize by greedy: $\hat{a} = \operatorname{argmax}_{a \in [K]} \hat{\mu}(a)$.

Gaussian bandits

$$R_i \mid_{A_i = a} \sim N(r(a), 1), \quad a \in [K] \text{ arms}.$$

a is pulled N(a) times in the pre-collected dataset D = {(A_i, R_i)}.
Following the naive value iterations,

• Estimate by sample mean:
$$\widehat{\mu}(a) = \frac{\sum_{i:A_i=a} R_i}{N(a)}$$
.

• Optimize by greedy:
$$\hat{a} = \operatorname{argmax}_{a \in [K]} \hat{\mu}(a)$$
.

For small N(a), a bad arm a might appear good by chance.

For small N(a), a bad arm a might appear good by chance.



• Policy $\widetilde{\pi}$ insufficiently covered by dataset \mathcal{D}

 $\Rightarrow \quad \text{Large uncertainty in our knowledge about a policy } \widetilde{\pi}.$

Policy π̃ insufficiently covered by dataset D
 ⇒ Large uncertainty in our knowledge about a policy π̃.

 Epistemic Uncertainty spuriously correlates with decision-making (greedy step). Roughly,

$$J(\widehat{\pi}) = J(\operatorname*{argmax}_{\pi} \ \widehat{J}(\pi)).$$

 \widehat{J} might be far from J for some π .

► Policy $\tilde{\pi}$ insufficiently covered by dataset \mathcal{D} ⇒ Large uncertainty in our knowledge about a policy $\tilde{\pi}$.

 Epistemic Uncertainty spuriously correlates with decision-making (greedy step). Roughly,

$$J(\widehat{\pi}) = J\left(\operatorname*{argmax}_{\pi} \, \widehat{J}(\pi)\right).$$

 \widehat{J} might be far from J for some π .



Policy *π̃* insufficiently covered by dataset *D* ⇒ Large uncertainty in our knowledge about a policy *π̃*.

 Epistemic Uncertainty spuriously correlates with decision-making (greedy step). Roughly,

$$J(\widehat{\pi}) = J\left(\operatorname*{argmax}_{\pi} \, \widehat{J}(\pi)\right).$$

 \widehat{J} might be far from J for some π .

- **•** Ruined if a bad π with large uncertainty appears to be good!
- No further interactions with MDP \Rightarrow unable to reduce uncertainty.

- **•** Ruined if a bad π with large uncertainty appears to be good!
- No further interactions with MDP \Rightarrow unable to reduce uncertainty.
- ▶ Uniform coverage for all policies? Too strong & unrealistic.

Question

Is it possible to design a provably efficient algorithm for offline RL under minimal assumptions on the dataset?

Question

Is it possible to design a provably efficient algorithm for offline RL under minimal assumptions on the dataset?

• Our solution by **Pessimism**: penalize large epistemic uncertainties.

Question

Is it possible to design a provably efficient algorithm for offline RL under minimal assumptions on the dataset?

- Our solution by **Pessimism**: penalize large epistemic uncertainties.
 - High estimated value, high uncertainty 🔅

Question

Is it possible to design a provably efficient algorithm for offline RL under minimal assumptions on the dataset?

- Our solution by **Pessimism**: penalize large epistemic uncertainties.
 - High estimated value, high uncertainty 😄
 - High estimated value, low uncertainty 😂

Question

Is it possible to design a provably efficient algorithm for offline RL under minimal assumptions on the dataset?

• Our solution by **Pessimism**: penalize large epistemic uncertainties.

- High estimated value, high uncertainty 😄
- High estimated value, low uncertainty ^(C)
- Good estimated value, low uncertainty 😂

A Solution: Pessimism for Offline Learning

- Pessimistic Value Iteration
- Why Pessimism Helps
- PEVI for Linear MDP

Pessimism for Offline Learning Pessimism Principle

Empirical success in practice:

- Pessimistic model-based [Yu et al. (2020); Kidambi et al. (2020)]
- Pessimistic value-based [Kumar et al. (2020)]
- In theory:
 - Regularized fitted Q-iterations [Liu et al. (2020)]: restrict policy class to be close to behavior policy.
 - Importance of Pessimism [Buckman et al. (2020)]

Pessimism for Offline Learning Pessimism Principle

Empirical success in practice:

- Pessimistic model-based [Yu et al. (2020); Kidambi et al. (2020)]
- Pessimistic value-based [Kumar et al. (2020)]
- In theory:
 - Regularized fitted Q-iterations [Liu et al. (2020)]: restrict policy class to be close to behavior policy.
 - Importance of Pessimism [Buckman et al. (2020)]

- This work: Pessimistic Value Iteration.
 - A principled framework for pessimism in value iterations.
 - No restriction on policy class and coverage of dataset.
 - Optimality of pessimism in the sense of information theory.

Algorithm: Pessimistic Value Iterations (General Form)

 $\blacktriangleright \quad \text{Estimate: } \overline{Q}_h \leftarrow \text{Regress}(\mathbb{B}_h \widehat{Q}_{h+1}, \mathcal{D}, \mathcal{F}).$



Algorithm: Pessimistic Value Iterations (General Form)

 $\blacktriangleright \quad \text{Estimate: } \overline{Q}_h \leftarrow \text{Regress}(\mathbb{B}_h \widehat{Q}_{h+1}, \mathcal{D}, \mathcal{F}).$

Uncertainty quantification (UQ): w.h.p.

$$\left|\overline{Q}_{h} - (\mathbb{B}_{h}\widehat{Q}_{h+1})\right| \leq \Gamma_{h}, \quad \forall h \in [H].$$



Algorithm: Pessimistic Value Iterations (General Form)

- Estimate: $\overline{Q}_h \leftarrow \operatorname{Regress}(\mathbb{B}_h \widehat{Q}_{h+1}, \mathcal{D}, \mathcal{F}).$
- Uncertainty quantification (UQ): w.h.p. $\left|\overline{Q}_{h} (\mathbb{B}_{h}\widehat{Q}_{h+1})\right| \leq \Gamma_{h}$
- Construct pessimistic value function

 \mathbb{B}_h

$$\widehat{Q}_{h}(x,a) = \underbrace{\overline{Q}_{h}(x,a)}_{\text{VI}} \underbrace{-\Gamma_{h}(x,a)}_{\text{penalty}}$$

$$\widehat{Q}_{h+1}$$

 Q_h

Algorithm: Pessimistic Value Iterations (General Form)

• Estimate: $\overline{Q}_h \leftarrow \operatorname{Regress}(\mathbb{B}_h \widehat{Q}_{h+1}, \mathcal{D}, \mathcal{F}).$

Uncertainty quantification (UQ): w.h.p.

$$\left|\overline{Q}_{h} - (\mathbb{B}_{h}\widehat{Q}_{h+1})\right| \leq \frac{\Gamma_{h}}{\Lambda}, \quad \forall h \in [H].$$

Construct pessimistic value function

$$\widehat{Q}_{h}(x,a) = \underbrace{\overline{Q}_{h}(x,a)}_{\text{VI}} \underbrace{-\Gamma_{h}(x,a)}_{\text{penalty}}$$

• Optimize: $\widehat{\pi}_h(x) = \operatorname{argmax}_{a \in \mathcal{A}} \widehat{Q}_h(x, a).$

Why Pessimism Helps? Boundedness of Evaluation Error¹



¹Adapted from Lemma 5.1 in (JYW'20)

Why Pessimism Helps? Boundedness of Evaluation Error¹



Define the model evaluation error

$$\iota_h(x,a) = (\mathbb{B}_h \widehat{Q}_{h+1})(x,a) - \widehat{Q}_h(x,a).$$

¹Adapted from Lemma 5.1 in (JYW'20)
Why Pessimism Helps? Boundedness of Evaluation Error¹



Define the model evaluation error

$$\iota_h(x,a) = (\mathbb{B}_h \widehat{Q}_{h+1})(x,a) - \widehat{Q}_h(x,a).$$

By pessimistic construction, the model evaluation error satisfy

$$\mathbb{B}_{h}\widehat{Q}_{h+1} \in \left[\widehat{Q}_{h}, \widehat{Q}_{h} + 2\Gamma_{h}\right]$$

$$\Rightarrow \quad 0 \leq \iota_{h}(s_{h}, a_{h}) \leq 2\Gamma_{h}(s_{h}, a_{h}).$$

¹Adapted from Lemma 5.1 in (JYW'20)

Why Pessimism Helps? Decomposition of Suboptimality²

We have the decomposition of suboptimality

$$\begin{split} \mathsf{SubOpt}(\widehat{\pi}; x) &= \underbrace{-\sum_{h=1}^{H} \mathbb{E}_{\widehat{\pi}} \left[\iota_h(s_h, a_h) \, \big| \, s_1 = x \right]}_{(\mathbf{i}): \; \mathbf{Spurious \; Correlation}} + \underbrace{\sum_{h=1}^{H} \mathbb{E}_{\pi^\star} \left[\iota_h(s_h, a_h) \, \big| \, s_1 = x \right]}_{(\mathbf{i}): \; \mathbf{Intrinsic \; Uncertainty}} \\ &+ \underbrace{\sum_{h=1}^{H} \mathbb{E}_{\pi^\star} \left[\langle \widehat{Q}_h(s_h, \cdot), \pi_h^\star(\cdot \, | \, s_h) - \widehat{\pi}_h(\cdot \, | \, s_h) \rangle_{\mathcal{A}} \, \big| \, s_1 = x \right]}_{\mathbf{i}}. \end{split}$$

(iii): Optimization Error

²Adapted from Lemma 3.1 in (JYW'20)

By pessimistic construction, the model evaluation error satisfy

 $0 \leq \iota_h(s_h, a_h) \leq 2\Gamma_h(s_h, a_h).$

► (i) spurious correlation is always non-positive

$$-\sum_{h=1}^{H} \mathbb{E}_{\widehat{\pi}} \left[\iota_h(s_h, a_h) \, \big| \, s_1 = x \right] \le 0.$$

³Adapted from Theorem 4.2 in (JYW'20)

- ► (i) spurious correlation is always non-positive
- ▶ Greedy policy ensures (iii) optimization error is always non-positive

$$\sum_{h=1}^{H} \mathbb{E}_{\pi^*} \left[\langle \widehat{Q}_h(s_h, \cdot), \pi_h^*(\cdot \mid s_h) - \widehat{\pi}_h(\cdot \mid s_h) \rangle_{\mathcal{A}} \mid s_1 = x \right] \le 0.$$

³Adapted from Theorem 4.2 in (JYW'20)

- ► (i) spurious correlation is always non-positive
- ▶ Greedy policy ensures (iii) optimization error is always non-positive
- A clean suboptimality bound

$$\mathsf{SubOpt}(\widehat{\pi}; x) \leq \sum_{h=1}^{H} \mathbb{E}_{\pi^*} \left[\iota_h(s_h, a_h) \, \big| \, s_1 = x \right]$$
$$\leq 2 \sum_{h=1}^{H} \mathbb{E}_{\pi^*} \left[\Gamma_h(s_h, a_h) \, \big| \, s_1 = x \right]$$

³Adapted from Theorem 4.2 in (JYW'20)

- ► (i) spurious correlation is always non-positive
- ▶ Greedy policy ensures (iii) optimization error is always non-positive
- A clean suboptimality bound

$$\mathsf{SubOpt}(\widehat{\pi}; x) \le 2\sum_{h=1}^{H} \mathbb{E}_{\pi^*} \left[\Gamma_h(s_h, a_h) \, \big| \, s_1 = x \right]$$

- Only depends on the trajectory of π^{\ast}
- Pessimism eliminates spurious correlation.

³Adapted from Theorem 4.2 in (JYW'20)

- ► (i) spurious correlation is always non-positive
- ▶ Greedy policy ensures (iii) optimization error is always non-positive
- A clean suboptimality bound

$$\mathsf{SubOpt}(\widehat{\pi}; x) \le 2\sum_{h=1}^{H} \mathbb{E}_{\pi^*} \left[\Gamma_h(s_h, a_h) \, \big| \, s_1 = x \right]$$

- Only depends on the trajectory of π^{*}
- Pessimism eliminates spurious correlation.

Question

How to construct the uncertainty quantifier Γ_h ?

³Adapted from Theorem 4.2 in (JYW'20)

Instantiation of PEVI Warm-up: i.i.d. Tabular Case

Assume \mathcal{D} consists of K i.i.d. trajectories from behavior policy π^b .

Assume concentrability coefficient

$$\sup_{x,a,h} \frac{\nu_h^\star(x,a)}{\nu_h^b(x,a)} \le \kappa^\star.$$

Instantiation of PEVI Warm-up: i.i.d. Tabular Case

Assume \mathcal{D} consists of K i.i.d. trajectories from behavior policy π^b .

Assume concentrability coefficient

$$\sup_{x,a,h} \frac{\nu_h^\star(x,a)}{\nu_h^b(x,a)} \le \kappa^\star.$$

Uncertainty quantifier

$$\Gamma_h(x,a) \propto N_h(x,a)^{-1/2}$$

Instantiation of PEVI Warm-up: i.i.d. Tabular Case

Assume \mathcal{D} consists of K i.i.d. trajectories from behavior policy π^b .

Assume concentrability coefficient

$$\sup_{x,a,h} \frac{\nu_h^\star(x,a)}{\nu_h^b(x,a)} \le \kappa^\star.$$

Uncertainty quantifier

$$\Gamma_h(x,a) \propto N_h(x,a)^{-1/2}$$

$$\mathsf{SubOpt}(\widehat{\pi}; x) \leq \sqrt{S^2 A} \cdot H^2 \sqrt{\kappa^* / K}.$$

Instantiation of PEVI Linear MDP

Definition (Linear MDP)

We say an episodic MDP (S, A, H, \mathbb{P}, r) is a linear MDP with a known feature map $\phi : S \times A \to \mathbb{R}^d$ if there exist d unknown (signed) measures $\mu_h = (\mu_h^{(1)}, \dots, \mu_h^{(d)})$ over S and an unknown vector $\theta_h \in \mathbb{R}^d$ such that

$$\mathbb{P}_{h}(x' \mid x, a) = \langle \phi(x, a), \mu_{h}(x') \rangle,$$

$$\mathbb{E}[r_{h}(s_{h}, a_{h}) \mid s_{h} = x, a_{h} = a] = \langle \phi(x, a), \theta_{h} \rangle$$

for all $(x, a, x') \in S \times A \times S$ at each step $h \in [H]$. Here we assume $\|\phi(x, a)\| \leq 1$ for all $(x, a) \in S \times A$ and $\max\{\|\mu_h(S)\|, \|\theta_h\|\} \leq \sqrt{d}$ at each step $h \in [H]$, where $\|\mu_h(S)\| = \int_S \|\mu_h(x)\| dx$.

• Linearity of Bellman update: $\mathbb{B}_h \widehat{Q}_{h+1} = \phi^\top \widehat{\theta}_h$ for some $\widehat{\theta}_h \in \mathbb{R}^d$.

• Linear function approximation $\mathcal{F} = \{ f_{\theta}(x, a) = \phi(x, a)^{\top} \theta, \ \theta \in \mathbb{R}^d \}.$

Instantiation of PEVI Linear MDP

Algorithm: PEVI for Linear MDP

- Estimate: $\overline{Q}_h(x,a) = \phi(x,a)^\top \widehat{\theta}_h$ via ridge regression.
- Uncertainty quantification

$$\Gamma_h(x,a) \asymp dH \cdot \left(\phi(x,a)^\top \Lambda_h^{-1} \phi(x,a)\right)^{1/2},$$

where Λ_h is the augmented sample covariance matrix of $\phi(s_h, a_h)$. Pessimistic value function

$$\widehat{Q}_h(x,a) = \phi(x,a)^\top \widehat{\theta}_h - c \cdot dH \cdot \left(\phi(x,a)^\top \Lambda_h^{-1} \phi(x,a)\right)^{1/2}$$

• Optimize: $\widehat{\pi}_h(x) = \operatorname{argmax}_{a \in \mathcal{A}} \widehat{Q}_h(x, a).$

Instantiation of PEVI - Linear MDP

Compliance Assumption

Assumption: Compliance

Let $\mathbb{P}_{\mathcal{D}}$ be the joint distribution of the dataset $\mathcal{D} = \{(x_h^{\tau}, a_h^{\tau}, r_h^{\tau})\}_{\tau, h=1}^{K, H}$. We say \mathcal{D} is compliant with an MDP $(\mathcal{S}, \mathcal{A}, H, \mathbb{P}, r)$ if

$$\mathbb{P}_{\mathcal{D}}\left(r_{h}^{\tau} = r', x_{h+1}^{\tau} = x' \mid \{(x_{h}^{j}, a_{h}^{j})\}_{j=1}^{\tau}, \{(r_{h}^{j}, x_{h+1}^{j})\}_{j=1}^{\tau-1}\right)$$
$$= \mathbb{P}\left(r_{h} = r', s_{h+1} = x' \mid s_{h} = x_{h}^{\tau}, a_{h} = a_{h}^{\tau}\right)$$

for all $r' \in [0,1]$, $x' \in S$, $h \in [H]$, $\tau \in [K]$. Here \mathbb{P} is taken with respect to the underlying MDP.

- Only require that \mathcal{D} evolves according to the MDP.
- Minimal assumptions on actions a^t_h: allow for arbitrarily collected data.
 - i.i.d. trajectories from a behavior policy \checkmark
 - sequentially adjusted actions $a_h^\tau \in \sigma(\{x_{h+1}^j, r_h^j\}_{j < \tau})$ \checkmark

Theorem 4.4 (JYW'20)

If ${\mathcal D}$ is compliant with the underlying MDP, then w.h.p,

$$\mathsf{SubOpt}\big(\widehat{\pi};x\big) \leq c \cdot dH \sum_{h=1}^{H} \mathbb{E}_{\pi^{\star}} \Big[\big(\phi(s_h, a_h)^{\top} \Lambda_h^{-1} \phi(s_h, a_h)\big)^{1/2} \, \Big| \, s_1 = x \Big]$$

up to logarithm factors of d, H, K.

Theorem 4.4 (JYW'20)

If ${\mathcal D}$ is compliant with the underlying MDP, then w.h.p,

$$\mathsf{SubOpt}\big(\widehat{\pi};x\big) \leq c \cdot dH \sum_{h=1}^{H} \mathbb{E}_{\pi^{\star}} \Big[\big(\phi(s_h, a_h)^{\top} \Lambda_h^{-1} \phi(s_h, a_h)\big)^{1/2} \, \Big| \, s_1 = x \Big]$$

up to logarithm factors of d, H, K.

▶ Minimal-assumption guarantee: only require compliance of *D*.

Theorem 4.4 (JYW'20)

If ${\mathcal D}$ is compliant with the underlying MDP, then w.h.p,

$$\mathsf{SubOpt}\big(\widehat{\pi};x\big) \leq c \cdot dH \sum_{h=1}^{H} \mathbb{E}_{\pi^{\star}} \Big[\big(\phi(s_h, a_h)^{\top} \Lambda_h^{-1} \phi(s_h, a_h)\big)^{1/2} \, \Big| \, s_1 = x \Big].$$

up to logarithm factors of d, H, K.

- ▶ Minimal-assumption guarantee: only require compliance of *D*.
- Oracle property: only depends on how well π* is covered no requirement on coverage of all trajectories.

Theorem 4.4 (JYW'20)

If ${\mathcal D}$ is compliant with the underlying MDP, then w.h.p,

$$\mathsf{SubOpt}(\widehat{\pi}; x) \leq c \cdot dH \sum_{h=1}^{H} \mathbb{E}_{\pi^{\star}} \left[\left(\phi(s_h, a_h)^{\top} \Lambda_h^{-1} \phi(s_h, a_h) \right)^{1/2} \middle| s_1 = x \right].$$

up to logarithm factors of d, H, K.

- ▶ Minimal-assumption guarantee: only require compliance of *D*.
- Oracle property: only depends on how well π* is covered no requirement on coverage of all trajectories.
- Data-dependent upper bound: (offline) data is what it is.

Instantiation of PEVI - Linear MDP Suboptimality Upper Bound: Special Cases

- ▶ Well-explored case: ≈ uniformly good coverage of all policies.
- ► The suboptimality achieves √1/K rate if D consist of K i.i.d. trajectories from behavior policy π
 and

$$\lambda_{\min} \left(\mathbb{E}_{\bar{\pi}} [\phi(s_h, a_h) \phi(s_h, a_h)^\top] \right) \ge c \quad \text{for some } c > 0.$$

Instantiation of PEVI - Linear MDP Suboptimality Upper Bound: Special Cases

Essentially-explored case: ≈ good coverage of optimal policy.
 The suboptimality achieves √1/K rate if

$$\Lambda_h \succeq \mathbf{I} + c \cdot K \cdot \mathbb{E}_{\pi^*} \left[\phi(s_h, a_h) \phi(s_h, a_h)^\top \, \big| \, s_1 = x \right].$$

Instantiation of PEVI - **Linear MDP** Suboptimality Upper Bound: Special Cases

Essentially-explored case: \approx good coverage of optimal policy.

• The suboptimality achieves $\sqrt{1/K}$ rate if

$$\Lambda_h \succeq \mathbf{I} + c \cdot K \cdot \mathbb{E}_{\pi^*} \big[\phi(s_h, a_h) \phi(s_h, a_h)^\top \, \big| \, s_1 = x \big].$$

Question

Is coverage of optimal π^* the essential information in \mathcal{D} ?

Is Pessimism Efficient?

Minimax Lower Bounds for Linear MDP

• Answer: Coverage of optimal π^* is the essential information in \mathcal{D} .

Pessimism is (nearly) minimax optimal in linear setting.

Answer: Coverage of optimal π* is the essential information in D.
 Pessimism is (nearly) minimax optimal in linear setting.

Minimax Optimality in Linear MDP

▶ Upper bound: pessimistic policy $\hat{\pi}$ and compliant $\mathcal{D} \sim \mathcal{M}$,

$$\mathsf{SubOpt}\big(\mathcal{M}, \widehat{\pi}; x\big) \leq c \cdot dH \sum_{h=1}^{H} \mathbb{E}_{\pi^{\star}} \Big[\big(\phi(s_h, a_h)^{\top} \Lambda_h^{-1} \phi(s_h, a_h)\big)^{1/2} \,\Big| \, s_1 = x \Big].$$

Answer: Coverage of optimal π* is the essential information in D.
 Pessimism is (nearly) minimax optimal in linear setting.

Minimax Optimality in Linear MDP

• Upper bound: pessimistic policy $\widehat{\pi}$ and compliant $\mathcal{D} \sim \mathcal{M}$,

$$\mathsf{SubOpt}(\mathcal{M},\widehat{\pi};x) \leq c \cdot dH \sum_{h=1}^{H} \mathbb{E}_{\pi^{\star}} \Big[\big(\phi(s_h, a_h)^{\top} \Lambda_h^{-1} \phi(s_h, a_h) \big)^{1/2} \, \Big| \, s_1 = x \Big].$$

• Lower bound: for any offline learning algorithm $Algo(\cdot)$,

$$\sup_{\mathcal{M},\mathcal{D}} \mathbb{E}_{\mathcal{D}} \left[\frac{\mathsf{SubOpt}(\mathcal{M}, \mathsf{Algo}(\mathcal{D}); x)}{\sum_{h=1}^{H} \mathbb{E}_{\pi^{\star}} \left[\left(\phi(s_h, a_h)^{\top} \Lambda_h^{-1} \phi(s_h, a_h) \right)^{1/2} \middle| s_1 = x \right]} \right] \ge c.$$

• Answer: Coverage of optimal π^* is the essential information in \mathcal{D} .

Pessimism is (nearly) minimax optimal in linear setting.

Minimax Optimality in Linear MDP

• Lower bound: for any offline learning algorithm $Algo(\cdot)$,

$$\sup_{\mathcal{M},\mathcal{D}} \mathbb{E}_{\mathcal{D}}\left[\frac{\mathsf{SubOpt}(\mathcal{M},\mathsf{Algo}(\mathcal{D});x)}{\sum_{h=1}^{H} \mathbb{E}_{\pi^{\star}}\left[\left(\phi(s_{h},a_{h})^{\top}\Lambda_{h}^{-1}\phi(s_{h},a_{h})\right)^{1/2} \middle| s_{1}=x\right]}\right] \geq c.$$

- Dependence on true MDP *M* and its optimal policy π^{*}.
- Essential Hardness in \mathcal{D} : how well (sample covariance) Λ_h covers π^* .

Thank you!

Proof of Minimax Lower Bound Construction of Hard Instance

• A subclass of linear MDP $M(p_1, p_2, p_3)$.



- Actions $\mathcal{A} = \{b_1, b_2, \dots, b_A\}$, states $\mathcal{S} = \{x_0, x_1, x_2\}$.
- Initial state s₁ ≡ x₀.
- Transition at the first step P₁(x₁ | x₀, b_j) = p_j, p_j = p₃ for j ≥ 3.
- Absorbing states x_1, x_2 , $\mathbb{P}_h(x_i | x_i, a) = 1$ for i = 1, 2 and all $a \in \mathcal{A}$.
- Deterministic rewards r_h(x₁, a) = 1 and r_h(x₂, a) = 0 for all a ∈ A.

Pre-determined actions in \mathcal{D} : take action b_j for n_j times.

• Two instances $\mathcal{M}_1 = M(p^*, p, p)$, $\mathcal{M}_2 = M(p, p^*, p)$, $p < p^*$.

• Optimal policy for \mathcal{M}_1 takes b_1 , while optimal policy for \mathcal{M}_2 takes b_2 .

• Two instances $\mathcal{M}_1 = M(p^*, p, p)$, $\mathcal{M}_2 = M(p, p^*, p)$, $p < p^*$.

• Optimal policy for \mathcal{M}_1 takes b_1 , while optimal policy for \mathcal{M}_2 takes b_2 .

• Suboptimality of any policy π on the two MDPs are

$$\mathsf{SubOpt}\big(\mathcal{M}_{\ell}, \pi; x_0\big) = (p^* - p)(H - 1)\big(1 - \pi_1(b_{\ell} \mid x_0)\big)$$

• 2-point argument: any policy makes mistake either on \mathcal{M}_1 or on \mathcal{M}_2 .

• Two instances $\mathcal{M}_1 = M(p^*, p, p)$, $\mathcal{M}_2 = M(p, p^*, p)$, $p < p^*$.

• Optimal policy for \mathcal{M}_1 takes b_1 , while optimal policy for \mathcal{M}_2 takes b_2 .

• Suboptimality of any policy π on the two MDPs are

$$\mathsf{SubOpt}\big(\mathcal{M}_{\ell}, \pi; x_0\big) = (p^* - p)(H - 1)\big(1 - \pi_1(\mathbf{b}_{\ell} \mid x_0)\big)$$

• 2-point argument: any policy makes mistake either on \mathcal{M}_1 or on \mathcal{M}_2 .

Reduction to Testing:

$$\max_{\ell \in \{1,2\}} \sqrt{n_{\ell}} \cdot \mathbb{E}_{\mathcal{D} \sim \mathcal{M}_{\ell}} \left[\mathsf{SubOpt} \left(\mathcal{M}_{\ell}, \mathsf{Algo}(\mathcal{D}); x_{0} \right) \right]$$

$$\geq \frac{\sqrt{n_{1}n_{2}}}{\sqrt{n_{1}} + \sqrt{n_{2}}} \cdot (p^{*} - p) \cdot (H - 1)$$

$$\times \left(\mathbb{E}_{\mathcal{D} \sim \mathcal{M}_{1}} \left[1 - \pi_{1}(b_{1} \mid x_{0}) \right] + \mathbb{E}_{\mathcal{D} \sim \mathcal{M}_{2}} \left[\pi_{1}(b_{1} \mid x_{0}) \right] \right)$$

$$\geq \frac{\sqrt{n_{1}n_{2}}}{\sqrt{n_{1}} + \sqrt{n_{2}}} \cdot (p^{*} - p) \cdot (H - 1) \left(1 - \mathsf{TV}(\mathbb{P}_{\mathcal{D} \sim \mathcal{M}_{1}}, \mathbb{P}_{\mathcal{D} \sim \mathcal{M}_{2}}) \right)$$

Reduction to Testing:

$$\max_{\ell \in \{1,2\}} \sqrt{n_{\ell}} \cdot \mathbb{E}_{\mathcal{D} \sim \mathcal{M}_{\ell}} \left[\mathsf{SubOpt} \left(\mathcal{M}_{\ell}, \mathtt{Algo}(\mathcal{D}); x_0 \right) \right] \geq c \cdot (H-1),$$

with a careful choice of n_1, n_2 and p, p^* .

Reduction to Testing:

$$\max_{\ell \in \{1,2\}} \sqrt{n_{\ell}} \cdot \mathbb{E}_{\mathcal{D} \sim \mathcal{M}_{\ell}} \Big[\mathsf{SubOpt} \big(\mathcal{M}_{\ell}, \mathtt{Algo}(\mathcal{D}); x_0 \big) \Big] \geq c \cdot (H-1),$$

with a careful choice of n_1, n_2 and p, p^* .

At the same time, w.h.p.,

$$\sum_{h=1}^{H} \mathbb{E}_{\pi^{*,\ell},\mathcal{M}_{\ell}} \left[\left(\phi(s_h, a_h)^{\top} \Lambda_h^{-1} \phi(s_h, a_h) \right)^{1/2} \middle| s_1 = x_0 \right] \approx (H-1)/\sqrt{n_{\ell}}.$$

Reduction to Testing:

$$\max_{\ell \in \{1,2\}} \sqrt{n_{\ell}} \cdot \mathbb{E}_{\mathcal{D} \sim \mathcal{M}_{\ell}} \Big[\mathsf{SubOpt} \big(\mathcal{M}_{\ell}, \mathtt{Algo}(\mathcal{D}); x_0 \big) \Big] \ge c \cdot (H-1),$$

with a careful choice of n_1, n_2 and p, p^* .

At the same time, w.h.p.,

$$\sum_{h=1}^{H} \mathbb{E}_{\pi^{*,\ell},\mathcal{M}_{\ell}} \left[\left(\phi(s_h, a_h)^{\top} \Lambda_h^{-1} \phi(s_h, a_h) \right)^{1/2} \middle| s_1 = x_0 \right] \approx (H-1)/\sqrt{n_{\ell}}.$$

Final lower bound

$$\sup_{\mathcal{M},\mathcal{D}} \mathbb{E}_{\mathcal{D}}\left[\frac{\mathsf{SubOpt}(\mathcal{M},\mathsf{Algo}(\mathcal{D});x)}{\sum_{h=1}^{H} \mathbb{E}_{\pi^{\star}}\left[\left(\phi(s_{h},a_{h})^{\top}\Lambda_{h}^{-1}\phi(s_{h},a_{h})\right)^{1/2} \middle| s_{1}=x\right]}\right] \geq c.$$

Proof of Upper Bound for Linear MDP Formula of Uncertainty Quantifier

▶ Bellman update $\mathbb{B}_h \widehat{Q}_{h+1}(x, a) = \phi(x, a)^\top w_h$ for some $w_h \in \mathbb{R}^d$.

Proof of Upper Bound for Linear MDP Formula of Uncertainty Quantifier

▶ Bellman update $\mathbb{B}_h \widehat{Q}_{h+1}(x, a) = \phi(x, a)^\top w_h$ for some $w_h \in \mathbb{R}^d$. ▶ Ridge estimator $\overline{Q}_h(x, a) = \phi(x, a)^\top \widehat{w}_h$, with

$$\begin{split} \widehat{w}_{h} &= \Lambda_{h}^{-1} \Big(\sum_{\tau=1}^{K} \phi(x_{h}^{\tau}, a_{h}^{\tau}) \cdot \left(r_{h}^{\tau} + \widehat{V}_{h+1}(x_{h+1}^{\tau}) \right) \Big), \\ \text{where} \quad \widehat{V}_{h+1}(x) &= \max_{a \in \mathcal{A}} \widehat{Q}_{h+1}(x, a), \\ \Lambda_{h} &= \sum_{\tau=1}^{K} \phi(x_{h}^{\tau}, a_{h}^{\tau}) \phi(x_{h}^{\tau}, a_{h}^{\tau})^{\top} + \lambda \cdot \mathbf{I}, \end{split}$$

Proof of Upper Bound for Linear MDP Formula of Uncertainty Quantifier

▶ Bellman update $\mathbb{B}_h \widehat{Q}_{h+1}(x, a) = \phi(x, a)^\top w_h$ for some $w_h \in \mathbb{R}^d$. ▶ Ridge estimator $\overline{Q}_h(x, a) = \phi(x, a)^\top \widehat{w}_h$, with

$$\begin{split} \hat{w}_{h} &= \Lambda_{h}^{-1} \Big(\sum_{\tau=1}^{K} \phi(x_{h}^{\tau}, a_{h}^{\tau}) \cdot \left(r_{h}^{\tau} + \hat{V}_{h+1}(x_{h+1}^{\tau}) \right) \Big), \\ \text{where} \quad \hat{V}_{h+1}(x) &= \max_{a \in \mathcal{A}} \hat{Q}_{h+1}(x, a), \\ \Lambda_{h} &= \sum_{\tau=1}^{K} \phi(x_{h}^{\tau}, a_{h}^{\tau}) \phi(x_{h}^{\tau}, a_{h}^{\tau})^{\top} + \lambda \cdot \mathbf{I}, \end{split}$$

Uncertainty quantifier chosen as

$$\Gamma_h(x,a) = \beta \cdot \left(\phi(x,a)^\top \Lambda_h^{-1} \phi(x,a)\right)^{1/2},$$

where the constant β is to be specified.
• Validity of Γ_h : w.h.p. for all (x, a) and all $h \in [H]$,

$$\left|\mathbb{B}_{h}\widehat{Q}_{h+1}(x,a) - \overline{Q}_{h}(x,a)\right| \leq \Gamma_{h}(x,a) = \beta \cdot \left(\phi(x,a)^{\top} \Lambda_{h}^{-1} \phi(x,a)\right)^{1/2}$$

The difference is decomposed into

$$\mathbb{B}_{h}\widehat{Q}_{h+1})(x,a) - \bar{Q}_{h}(x,a) = \phi(x,a)^{\top}(w_{h} - \widehat{w}_{h})$$

$$= \underbrace{\phi(x,a)^{\top}w_{h} - \phi(x,a)^{\top}\Lambda_{h}^{-1}\left(\sum_{\tau=1}^{K}\phi(x_{h}^{\tau},a_{h}^{\tau}) \cdot (\mathbb{B}_{h}\widehat{Q}_{h+1})(x_{h}^{\tau},a_{h}^{\tau})\right)}_{(i)}_{(i)}$$

$$- \underbrace{\phi(x,a)^{\top}\Lambda_{h}^{-1}\left(\sum_{\tau=1}^{K}\phi(x_{h}^{\tau},a_{h}^{\tau}) \cdot (r_{h}^{\tau} + \widehat{V}_{h+1}(x_{h+1}^{\tau}) - (\mathbb{B}_{h}\widehat{Q}_{h+1})(x_{h}^{\tau},a_{h}^{\tau}))\right)}_{(i)}.$$
(ii)

Boundedness of w_h : since $\widehat{Q}_{h+1} \in [0, H-h]$, it holds that

 $\|\widehat{w}_h\| \le H\sqrt{Kd/\lambda}.$

Boundedness of w_h : since $\widehat{Q}_{h+1} \in [0, H-h]$, it holds that

 $\|\widehat{w}_h\| \le H\sqrt{Kd/\lambda}.$

The first term bounded as

$$\left| (\mathbf{i}) \right| = \lambda \cdot \left| \phi(x, a)^{\top} \Lambda_h^{-1} w_h \right| \le H \sqrt{d\lambda} \cdot \sqrt{\phi(x, a)^{\top} \Lambda_h^{-1} \phi(x, a)}.$$

▶ Boundedness of w_h : since $\widehat{Q}_{h+1} \in [0, H-h]$, it holds that $\|\widehat{w}_h\| \leq H\sqrt{Kd/\lambda}.$

The first term bounded as

$$\left| (\mathbf{i}) \right| = \lambda \cdot \left| \phi(x, a)^{\top} \Lambda_h^{-1} w_h \right| \le H \sqrt{d\lambda} \cdot \sqrt{\phi(x, a)^{\top} \Lambda_h^{-1} \phi(x, a)}.$$

The second term bounded as

$$\begin{split} \left| (\mathrm{ii}) \right| &\leq \Big\| \sum_{\tau=1}^{K} \phi(x_h^{\tau}, a_h^{\tau}) \cdot \epsilon_h^{\tau}(\hat{V}_{h+1}) \Big\|_{\Lambda_h^{-1}} \cdot \sqrt{\phi(x, a)^{\top} \Lambda_h^{-1} \phi(x, a)} \\ \text{where } \epsilon_h^{\tau}(V) &= r_h^{\tau} + V(x_{h+1}^{\tau}) - \mathbb{E} \big[r_h(s_h, a_h) + V(s_{h+1}) \, \big| \, s_h = x_h^{\tau}, a_h = a_h^{\tau} \big]. \end{split}$$

Uniform concentration:

$$\left| (\mathsf{ii}) \right| \leq \sup_{V_{h+1} \in \mathcal{V}_{h+1}} \left\| \sum_{\tau=1}^{K} \phi(x_h^{\tau}, a_h^{\tau}) \cdot \epsilon_h^{\tau}(V_{h+1}) \right\|_{\Lambda_h^{-1}} \cdot \sqrt{\phi(x, a)^{\top} \Lambda_h^{-1} \phi(x, a)}.$$

• Supremum over function class \mathcal{V}_{h+1} with the form

$$V_{h+1}(x) = \max_{a \in \mathcal{A}} \left\{ \min\left\{ \phi(x, a)^\top \theta - \beta \cdot \sqrt{\phi(x, a)^\top \Sigma^{-1} \phi(x, a)}, H - h \right\}^+ \right\},\$$

for appropriately bounded $\theta \in \mathbb{R}^d$, $\beta \in \mathbb{R}$ and $\Sigma \in \mathbb{R}^{d \times d}$.

Uniform concentration:

$$\left| (\mathsf{ii}) \right| \leq \sup_{V_{h+1} \in \mathcal{V}_{h+1}} \left\| \sum_{\tau=1}^{K} \phi(x_h^{\tau}, a_h^{\tau}) \cdot \epsilon_h^{\tau}(V_{h+1}) \right\|_{\Lambda_h^{-1}} \cdot \sqrt{\phi(x, a)^{\top} \Lambda_h^{-1} \phi(x, a)}.$$

• Supremum over function class \mathcal{V}_{h+1} with the form

$$V_{h+1}(x) = \max_{a \in \mathcal{A}} \left\{ \min \left\{ \phi(x, a)^\top \theta - \beta \cdot \sqrt{\phi(x, a)^\top \Sigma^{-1} \phi(x, a)}, H - h \right\}^+ \right\},\$$

for appropriately bounded $\theta \in \mathbb{R}^d$, $\beta \in \mathbb{R}$ and $\Sigma \in \mathbb{R}^{d \times d}$.

The second term is bounded as

$$\left| \mathsf{(ii)} \right| \leq \beta/2 \cdot \sqrt{\phi(x,a)^\top \Lambda_h^{-1} \phi(x,a)}, \quad \beta = c \cdot dH \cdot \mathsf{PolyLog}(d,H,K)$$

- Concentration of self-normalized process for a single V_{h+1} . (Only compliance of \mathcal{D} is needed for the concentration.)
- ε -covering of linear function class \mathcal{V}_{h+1} for uniform concentration.