

Strategic Exploration in Large Scale MDPs

Recap on Bellman Error and Bellman Operator

Bellman error of $f(s, a)$

Recap on Bellman Error and Bellman Operator

Bellman error of $f(s, a)$

$$BE(s, a) = f(s, a) - \left(r(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} f(s', a') \right)$$

Recap on Bellman Error and Bellman Operator

Bellman error of $f(s, a)$

$$BE(s, a) = f(s, a) - \left(r(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} f(s', a') \right)$$

If $BE(s, a) = 0, \forall s, a$, then $f(s, a) = Q^*(s, a), \forall s, a$

Recap on Bellman Error and Bellman Operator

Bellman error of $f(s, a)$

$$BE(s, a) = f(s, a) - \left(r(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} f(s', a') \right)$$

If $BE(s, a) = 0, \forall s, a$, then $f(s, a) = Q^*(s, a), \forall s, a$

Bellman Operator \mathcal{T} of f

Recap on Bellman Error and Bellman Operator

Bellman error of $f(s, a)$

$$BE(s, a) = f(s, a) - \left(r(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} f(s', a') \right)$$

If $BE(s, a) = 0, \forall s, a$, then $f(s, a) = Q^*(s, a), \forall s, a$

Bellman Operator \mathcal{T} of f

$$\mathcal{T}f: S \times A \rightarrow \mathbb{R}$$

Recap on Bellman Error and Bellman Operator

Bellman error of $f(s, a)$

$$BE(s, a) = f(s, a) - \left(r(s, a) + \mathbb{E}_{s' \sim P(\cdot|s, a)} \max_{a'} f(s', a') \right)$$

If $BE(s, a) = 0, \forall s, a$, then $f(s, a) = Q^*(s, a), \forall s, a$

Bellman Operator \mathcal{T} of f

$$\mathcal{T}f: S \times A \rightarrow \mathbb{R}$$

$$[\mathcal{T}f](s, a) = r(s, a) + \mathbb{E}_{s' \sim P(\cdot|s, a)} \max_a f(s', a')$$

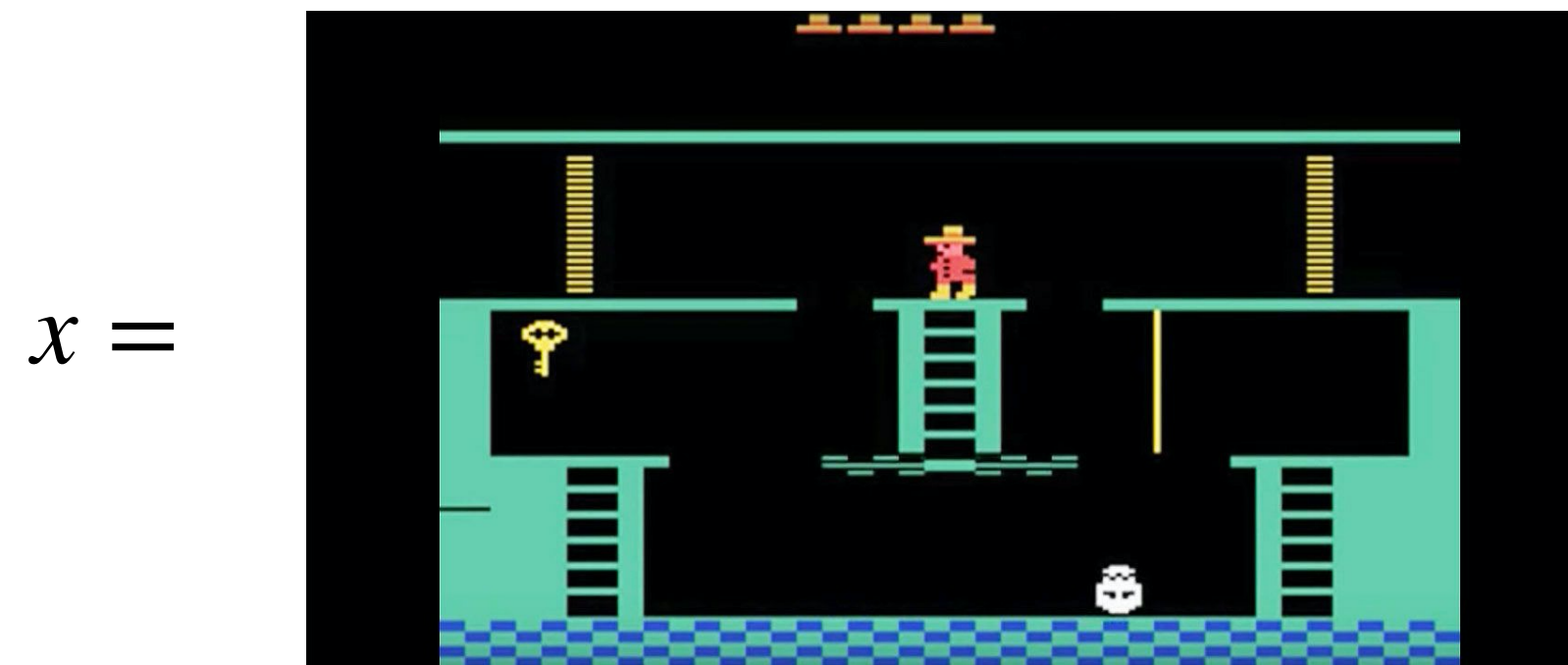
Notations

Probability of π visiting (s, a) at time step h : $d_h^\pi(s, a)$

Setting

Finite horizon episodic MDP $\{ \{X_h\}_{h=0}^H, \{A_h\}_{h=0}^{H-1}, H, x_0, r, P \}$

State space X_h is extremely large:



Not acceptable: $\text{poly}(|X|)$

Q: can we generalize using function approximation

Let's set up function class in RL setting

We will consider **Q function class** for now (and model class later)

$$\mathcal{F} \subset X \times A \mapsto [0,1]$$

Let's set up function class in RL setting

We will consider **Q function class** for now (and model class later)

$$\mathcal{F} \subset X \times A \mapsto [0,1]$$

Realizability assumption:

$$Q^* \in \mathcal{F}$$

Let's set up function class in RL setting

We will consider **Q function class** for now (and model class later)

$$\mathcal{F} \subset X \times A \mapsto [0,1]$$

Realizability assumption:

$$Q^* \in \mathcal{F}$$

Define **policy class**: $\Pi = \{ \pi : \pi(x) = \arg \max_{a \in A} f(x, a), \forall x \in X | f \in \mathcal{F} \}$

Let's set up function class in RL setting

We will consider **Q function class** for now (and model class later)

$$\mathcal{F} \subset X \times A \mapsto [0,1]$$

Realizability assumption:

$$Q^* \in \mathcal{F}$$

Define **policy class**: $\Pi = \{ \pi : \pi(x) = \arg \max_{a \in A} f(x, a), \forall x \in X | f \in \mathcal{F} \}$

i.e., each Q-approximator f induces a policy (greedy w.r.t f)

Learning Goal:

We will do PAC in this lecture rather than regret.

Given approximation error ϵ and failure prob δ ,
can we learn ϵ *near optimal policy* (i.e., $V^{\hat{\pi}} \geq V^* - \epsilon$) in # of samples scaling
poly with all relevant parameters (*here, we need poly in $\ln(|\mathcal{F}|)$*)

How to check if a Q-approximator is good?

We define **average** Bellman error below:

$$\mathcal{E}(f; \bar{f}, h) = \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right]$$

How to check if a Q-approximator is good?

We define **average** Bellman error below:

$$\mathcal{E}(f; \bar{f}, h) = \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right]$$

\bar{f} : defines roll-in distribution over x_h, a_h at stage h.

How to check if a Q-approximator is good?

We define **average** Bellman error below:

$$\mathcal{E}(f; \bar{f}, h) = \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right]$$

\bar{f} : defines roll-in distribution over x_h, a_h at stage h .

Bellman error measures consistency in one-step Bellman backup, e.g., $\mathcal{E}(Q^*; \bar{f}, h) = 0$

How to check if a Q-approximator is good?

We define **average** Bellman error below:

$$\mathcal{E}(f; \bar{f}, h) = \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right]$$

\bar{f} : defines roll-in distribution over x_h, a_h at stage h .

Bellman error measures consistency in one-step Bellman backup, e.g., $\mathcal{E}(Q^*; \bar{f}, h) = 0$

Hence, any f such that $\mathcal{E}(f; \pi, h) \neq 0$, is an incorrect Q^* approximator

Optimism Led Iterative Value Function Elimination (OLIVE)

Optimism Led Iterative Value Function Elimination (OLIVE)

Initialize $\mathcal{F}_0 = \mathcal{F}$

Optimism Led Iterative Value Function Elimination (OLIVE)

Initialize $\mathcal{F}_0 = \mathcal{F}$

For $t = 0, \dots$

Optimism Led Iterative Value Function Elimination (OLIVE)

Initialize $\mathcal{F}_0 = \mathcal{F}$

For $t = 0, \dots$

$$f_t = \arg \max_{f \in \mathcal{F}_t} \left(\max_a f(x_0, a) \right)$$

Optimism Led Iterative Value Function Elimination (OLIVE)

Initialize $\mathcal{F}_0 = \mathcal{F}$

For $t = 0, \dots$

$$f_t = \arg \max_{f \in \mathcal{F}_t} \left(\max_a f(x_0, a) \right)$$

$$\text{If } \left| \widetilde{V}^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \leq \epsilon, \text{ return } \pi_{f_t}$$

Optimism Led Iterative Value Function Elimination (OLIVE)

Initialize $\mathcal{F}_0 = \mathcal{F}$

For $t = 0, \dots$

$$f_t = \arg \max_{f \in \mathcal{F}_t} \left(\max_a f(x_0, a) \right)$$

$$\text{If } \left| \widetilde{V}^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \leq \epsilon, \text{ return } \pi_{f_t}$$

Version space update:

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \widetilde{\mathcal{E}}(f, \pi_{f_t}, h) \leq \delta, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

Estimating Bellman Error under a fixed Roll-in Policy:

Given a **fixed** $\pi_{\bar{f}}$, we can evaluate all f efficiently **statistically (not computationally)**:

$$\forall f: \mathcal{E}(f; \bar{f}, h) = \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right]$$

Estimating Bellman Error under a fixed Roll-in Policy:

Given a **fixed** $\pi_{\bar{f}}$, we can evaluate all f efficiently **statistically (not computationally)**:

$$\begin{aligned}\forall f: \mathcal{E}(f; \bar{f}, h) &= \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right] \\ &= \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}, x_{h+1} \sim P(\cdot | x_h, a_h)} \left[f(x_h, a_h) - r(x_h, a_h) - \max_{a \in \mathcal{A}} f(x_{h+1}, a) \right]\end{aligned}$$

Estimating Bellman Error under a fixed Roll-in Policy:

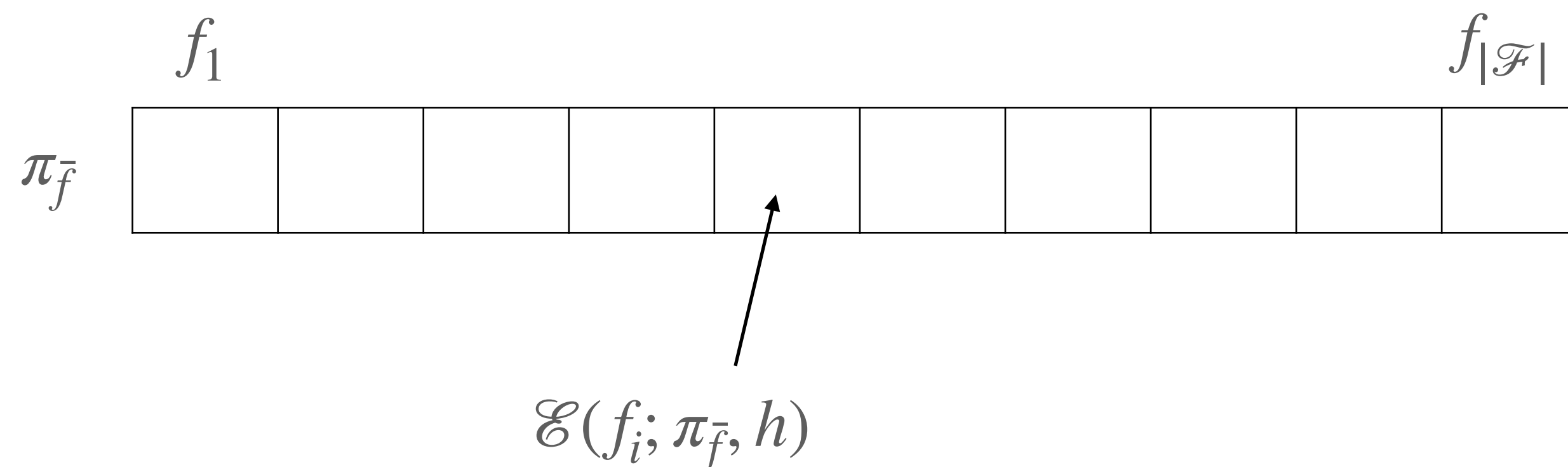
Given a **fixed** $\pi_{\bar{f}}$, we can evaluate all f efficiently **statistically (not computationally)**:

$$\begin{aligned}\forall f: \mathcal{E}(f; \bar{f}, h) &= \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right] \\ &= \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}, x_{h+1} \sim P(\cdot | x_h, a_h)} \left[f(x_h, a_h) - r(x_h, a_h) - \max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \\ &\approx \frac{1}{N} \sum_{i=1}^N \left[f(x_h^i, a_h^i) - r(x_h^i, a_h^i) - \max_{a \in \mathcal{A}} f(x_{h+1}^i, a) \right]\end{aligned}$$

Estimating Bellman Error under a fixed Roll-in Policy:

Given a **fixed** $\pi_{\bar{f}}$, we can evaluate all f efficiently **statistically (not computationally)**:

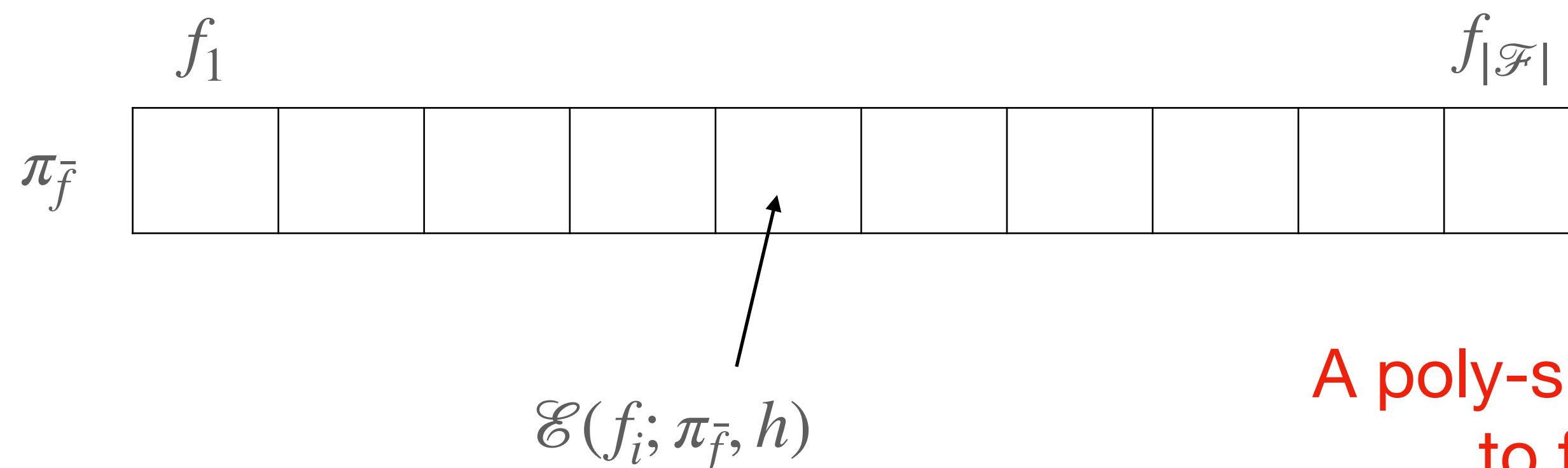
$$\begin{aligned} \forall f: \mathcal{E}(f; \bar{f}, h) &= \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right] \\ &= \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}, x_{h+1} \sim P(\cdot | x_h, a_h)} \left[f(x_h, a_h) - r(x_h, a_h) - \max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \\ &\approx \frac{1}{N} \sum_{i=1}^N \left[f(x_h^i, a_h^i) - r(x_h^i, a_h^i) - \max_{a \in \mathcal{A}} f(x_{h+1}^i, a) \right] \end{aligned}$$



Estimating Bellman Error under a fixed Roll-in Policy:

Given a **fixed** $\pi_{\bar{f}}$, we can evaluate all f efficiently **statistically (not computationally)**:

$$\begin{aligned} \forall f : \mathcal{E}(f; \bar{f}, h) &= \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right] \\ &= \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}, x_{h+1} \sim P(\cdot | x_h, a_h)} \left[f(x_h, a_h) - r(x_h, a_h) - \max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \\ &\approx \frac{1}{N} \sum_{i=1}^N \left[f(x_h^i, a_h^i) - r(x_h^i, a_h^i) - \max_{a \in \mathcal{A}} f(x_{h+1}^i, a) \right] \end{aligned}$$



A poly-size dataset allows us to fill up all entries

OLIVE Revisit (ignoring statistical error for simplicity)

Initialize $\mathcal{F}_0 = \mathcal{F}$

For $t = 0, \dots$

$$f_t = \arg \max_{f \in \mathcal{F}_t} \max_a f(x_0, a)$$

$$\text{If } \left| V^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \leq \epsilon, \text{ return } \pi_{f_t}$$

Version space update:

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \mathcal{E}(f; \pi_{f_t}, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

OLIVE Revisit (ignoring statistical error for simplicity)

Initialize $\mathcal{F}_0 = \mathcal{F}$

For $t = 0, \dots$

$$f_t = \arg \max_{f \in \mathcal{F}_t} \max_a f(x_0, a)$$

$$\text{If } \left| V^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \leq \epsilon, \text{ return } \pi_{f_t}$$

Version space update:

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \mathcal{E}(f; \pi_{f_t}, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

1. Upon termination we succeed (due to optimism)

OLIVE Revisit (ignoring statistical error for simplicity)

Initialize $\mathcal{F}_0 = \mathcal{F}$

For $t = 0, \dots$

$$f_t = \arg \max_{f \in \mathcal{F}_t} \max_a f(x_0, a)$$

$$\text{If } \left| V^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \leq \epsilon, \text{ return } \pi_{f_t}$$

Version space update:

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \mathcal{E}(f; \pi_{f_t}, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

1. Upon termination we succeed (due to optimism)
2. If not terminate, we make non-trivial progress

OLIVE Revisit (ignoring statistical error for simplicity)

Initialize $\mathcal{F}_0 = \mathcal{F}$

For $t = 0, \dots$

$$f_t = \arg \max_{f \in \mathcal{F}_t} \max_a f(x_0, a)$$

If $\left| V^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \leq \epsilon$, return π_{f_t}

Version space update:

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \mathcal{E}(f; \pi_{f_t}, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

1. Upon termination we succeed (due to optimism)
2. If not terminate, we make non-trivial progress
3. Total # of such non-trivial progress is bounded

Quality of Returned Policy upon Termination:

$$f_t = \arg \max_{f \in \mathcal{F}_t} \max_a f(s_0, a)$$
$$\text{If } \left| V^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \leq \epsilon, \text{ return } \pi_{f_t}$$

Quality of Returned Policy upon Termination:

$$f_t = \arg \max_{f \in \mathcal{F}_t} \max_a f(s_0, a)$$
$$\text{If } \left| V^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \leq \epsilon, \text{ return } \pi_{f_t}$$

Note $Q^* \in \mathcal{F}_t, \forall t$,

(we only eliminate things that are obviously wrong & Q^* has zero bellman error everywhere)

Quality of Returned Policy upon Termination:

$$f_t = \arg \max_{f \in \mathcal{F}_t} \max_a f(s_0, a)$$
$$\text{If } \left| V^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \leq \epsilon, \text{ return } \pi_{f_t}$$

Note $Q^* \in \mathcal{F}_t, \forall t$,

(we only eliminate things that are obviously wrong & Q^* has zero bellman error everywhere)

$$\max_a f_t(x_0, a) \geq \max_a Q^*(x_0, a) = V^*(x_0)$$

Quality of Returned Policy upon Termination:

$$f_t = \arg \max_{f \in \mathcal{F}_t} \max_a f(s_0, a)$$
$$\text{If } \left| V^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \leq \epsilon, \text{ return } \pi_{f_t}$$

Note $Q^* \in \mathcal{F}_t, \forall t$,

(we only eliminate things that are obviously wrong & Q^* has zero bellman error everywhere)

$$\max_a f_t(x_0, a) \geq \max_a Q^*(x_0, a) = V^*(x_0)$$

$$\left| V^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \leq \epsilon \Rightarrow V^{\pi_{f_t}} \geq \max_a f_t(x_0, a) - \epsilon \geq V^*(x_0) - \epsilon$$

Quality of Returned Policy upon Termination:

$$f_t = \arg \max_{f \in \mathcal{F}_t} \max_a f(s_0, a)$$
$$\text{If } \left| V^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \leq \epsilon, \text{ return } \pi_{f_t}$$

Note $Q^* \in \mathcal{F}_t, \forall t$,

(we only eliminate things that are obviously wrong & Q^* has zero bellman error everywhere)

$$\max_a f_t(x_0, a) \geq \max_a Q^*(x_0, a) = V^*(x_0)$$

$$\left| V^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \leq \epsilon \Rightarrow V^{\pi_{f_t}} \geq \max_a f_t(x_0, a) - \epsilon \geq V^*(x_0) - \epsilon$$

Optimism ensures that once termination happens, we are done!

No termination means we found a bad Q^* -approximator:

Claim [performance difference lemma]:

$$\max_a f_t(x_0, a) - V^{\pi_{f_t}}(x_0) = \sum_{h=0}^{H-1} \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{f_t}}} \left[f_t(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_a f_t(x_{h+1}, a) \right] \right]$$

Proof: a straight telescoping sum

No termination means we found a bad Q^* -approximator:

Claim [performance difference lemma]:

$$\max_a f_t(x_0, a) - V^{\pi_{f_t}}(x_0) = \sum_{h=0}^{H-1} \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{f_t}}} \left[f_t(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_a f_t(x_{h+1}, a) \right] \right]$$

Proof: a straight telescoping sum

$$f_t(x_0, \pi_{f_t}(x_0)) - r(x_0, \pi_{f_t}(x_0)) - \mathbb{E}_{x_1 \sim P(\cdot | x_0, \pi_{f_t}(x_0))} \max_a f_t(x_1, a)$$

No termination means we found a bad Q^* -approximator:

Claim [performance difference lemma]:

$$\max_a f_t(x_0, a) - V^{\pi_{f_t}}(x_0) = \sum_{h=0}^{H-1} \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{f_t}}} \left[f_t(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_a f_t(x_{h+1}, a) \right] \right]$$

Proof: a straight telescoping sum

$$\begin{aligned} & f_t(x_0, \pi_{f_t}(x_0)) - r(x_0, \pi_{f_t}(x_0)) - \mathbb{E}_{x_1 \sim P(\cdot | x_0, \pi_{f_t}(x_0))} \max_a f_t(x_1, a) \\ & + \mathbb{E}_{x_1 \sim d_1^{\pi_{f_t}}} \left[f_t(x_1, \pi_{f_t}(x_1)) - r(x_1, \pi_{f_t}(x_1)) - \mathbb{E}_{x_2 \sim P(\cdot | x_1, \pi_{f_t}(x_1))} \max_a f_t(x_2, a) \right] \end{aligned}$$

No termination means we found a bad Q^* -approximator:

Claim [performance difference lemma]:

$$\max_a f_t(x_0, a) - V^{\pi_{f_t}}(x_0) = \sum_{h=0}^{H-1} \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{f_t}}} \left[f_t(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_a f_t(x_{h+1}, a) \right] \right]$$

Proof: a straight telescoping sum

$$\begin{aligned} & f_t(x_0, \pi_{f_t}(x_0)) - r(x_0, \pi_{f_t}(x_0)) - \mathbb{E}_{x_1 \sim P(\cdot | x_0, \pi_{f_t}(x_0))} \max_a f_t(x_1, a) \\ & + \mathbb{E}_{x_1 \sim d_1^{\pi_{f_t}}} \left[f_t(x_1, \pi_{f_t}(x_1)) - r(x_1, \pi_{f_t}(x_1)) - \mathbb{E}_{x_2 \sim P(\cdot | x_1, \pi_{f_t}(x_1))} \max_a f_t(x_2, a) \right] \\ & + \mathbb{E}_{x_2 \sim d_2^{\pi_{f_t}}} \left[f_t(x_2, \pi_{f_t}(x_2)) - r(x_2, \pi_{f_t}(x_2)) - \mathbb{E}_{x_3 \sim P(\cdot | x_2, \pi_{f_t}(x_2))} \max_a f_t(x_3, a) \right] \end{aligned}$$

No termination means we found a bad Q^* -approximator:

Claim [performance difference lemma]:

$$\max_a f_t(x_0, a) - V^{\pi_{f_t}}(x_0) = \sum_{h=0}^{H-1} \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{f_t}}} \left[f_t(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_a f_t(x_{h+1}, a) \right] \right]$$

Proof: a straight telescoping sum

$$\begin{aligned} & f_t(x_0, \pi_{f_t}(x_0)) - r(x_0, \pi_{f_t}(x_0)) - \mathbb{E}_{x_1 \sim P(\cdot | x_0, \pi_{f_t}(x_0))} \max_a f_t(x_1, a) \\ & \quad + \mathbb{E}_{x_1 \sim d_1^{\pi_{f_t}}} \left[f_t(x_1, \pi_{f_t}(x_1)) - r(x_1, \pi_{f_t}(x_1)) - \mathbb{E}_{x_2 \sim P(\cdot | x_1, \pi_{f_t}(x_1))} \max_a f_t(x_2, a) \right] \\ & \quad \quad \quad + \mathbb{E}_{x_2 \sim d_2^{\pi_{f_t}}} \left[f_t(x_2, \pi_{f_t}(x_2)) - r(x_2, \pi_{f_t}(x_2)) - \mathbb{E}_{x_3 \sim P(\cdot | x_2, \pi_{f_t}(x_2))} \max_a f_t(x_3, a) \right] \\ & \quad \quad \quad \dots \\ & \quad \quad \quad = f_t(x_0, \pi_{f_t}(x_0)) - \mathbb{E} \left[\sum_{h=0}^{H-1} r_h \right] \end{aligned}$$

No termination means we found a bad Q^* -approximator:

$$\max_a f_t(x_0, a) - V^{\pi_{f_t}}(x_0) = \sum_{h=0}^{H-1} \mathbb{E}_{x_h \sim d_h^{\pi_{f_t}}} \left[f_t(x_h, \pi_{f_t}(x_h)) - r(x_h, \pi_{f_t}(x_h)) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_a f_t(x_{h+1}, a) \right] \right]$$

If we do not terminate, i.e.,

$$\left| V^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \geq \epsilon,$$

then:

No termination means we found a bad Q^* -approximator:

$$\max_a f_t(x_0, a) - V^{\pi_{f_t}}(x_0) = \sum_{h=0}^{H-1} \mathbb{E}_{x_h \sim d_h^{\pi_{f_t}}} \left[f_t(x_h, \pi_{f_t}(x_h)) - r(x_h, \pi_{f_t}(x_h)) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_a f_t(x_{h+1}, a) \right] \right]$$

If we do not terminate, i.e.,

$$\left| V^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \geq \epsilon,$$

then:

$$\max_a f_t(x_0, a) - V^{\pi_{f_t}}(x_0) = \sum_{h=0}^{H-1} \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{f_t}}} \left[f_t(x_h, a_h) - r(x_h, \pi_{f_t}(x_h)) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_a f_t(x_{h+1}, a) \right] \right] \geq \epsilon$$

No termination means we found a bad Q^* -approximator:

$$\max_a f_t(x_0, a) - V^{\pi_{f_t}}(x_0) = \sum_{h=0}^{H-1} \mathbb{E}_{x_h \sim d_h^{\pi_{f_t}}} \left[f_t(x_h, \pi_{f_t}(x_h)) - r(x_h, \pi_{f_t}(x_h)) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_a f_t(x_{h+1}, a) \right] \right]$$

If we do not terminate, i.e.,

$$\left| V^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \geq \epsilon,$$

then:

$$\max_a f_t(x_0, a) - V^{\pi_{f_t}}(x_0) = \sum_{h=0}^{H-1} \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{f_t}}} \left[f_t(x_h, a_h) - r(x_h, \pi_{f_t}(x_h)) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_a f_t(x_{h+1}, a) \right] \right] \geq \epsilon$$

$$\Rightarrow \exists h, \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{f_t}}} \left[f_t(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_a f_t(x_{h+1}, a) \right] \right] \geq \epsilon/H$$

We need to argue how many episodes we have before termination

$$\forall h : \mathcal{E}_h \in \mathbb{R}^{|\mathcal{F}| \times |\mathcal{F}|}$$

f

π_f		$\mathcal{E}_{f, \pi_f, h}$			

Rank of this Matrix is defined as Bellman Rank

Progress on Value Function Elimination

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \mathcal{E}(f; \pi_f, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

$\pi_{f_{t-1}}$	= 0	= 0	= 0	= 0	= 0

Progress on Value Function Elimination

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \mathcal{E}(f; \pi_{f_t}, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

If we don't terminate at t:

			f_t			
π_{f_t}			$\neq 0$			
$\pi_{f_{t-1}}$	$= 0$	$= 0$	$= 0$	$= 0$	$= 0$	$= 0$
				$\mathcal{E}(f_t; \pi_{f_t}, h) \geq \epsilon/H$		

Progress on Value Function Elimination

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \mathcal{E}(f; \pi_{f_t}, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

If we don't terminate at t:

π_{f_t}			$\neq 0$		
$\pi_{f_{t-1}}$	$= 0$	$= 0$	$= 0$	$= 0$	$= 0$
				$\mathcal{E}(f_t; \pi_{f_t}, h) \geq \epsilon/H$	

Progress on Value Function Elimination

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \mathcal{E}(f; \pi_{f_t}, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

If we don't terminate at t:

π_{f_t}		$\neq 0$		$\neq 0$	
$\pi_{f_{t-1}}$	$= 0$	$= 0$	$= 0$	$= 0$	$= 0$
			$\mathcal{E}(f_t; \pi_{f_t}, h) \geq \epsilon/H$		

Progress on Value Function Elimination

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \mathcal{E}(f; \pi_f, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

If we don't terminate at t:

π_{f_t}		$\neq 0$		$\neq 0$	
$\pi_{f_{t-1}}$	$= 0$	$= 0$	$= 0$	$= 0$	$= 0$
			$\mathcal{E}(f_t; \pi_{f_t}, h) \geq \epsilon/H$		

Progress on Value Function Elimination

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \mathcal{E}(f; \pi_f, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

If we still cannot terminate at t+1:

		f_{t+1}	f_t			
$\pi_{f_{t+1}}$		$\neq 0$				
π_{f_t}			$\neq 0$		$\neq 0$	
$\pi_{f_{t-1}}$	$= 0$	$= 0$	$= 0$	$= 0$	$= 0$	$= 0$
				$\mathcal{E}(f_t; \pi_{f_t}, h) \geq \epsilon/H$		

If we don't terminate at t:

Progress on Value Function Elimination

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \mathcal{E}(f; \pi_f, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

If we still cannot terminate at t+1:

	f_{t+1}	f_t			
$\pi_{f_{t+1}}$	$\neq 0$				
π_{f_t}	$= 0$	$\neq 0$		$\neq 0$	
$\pi_{f_{t-1}}$	$= 0$	$= 0$	$= 0$	$= 0$	$= 0$
			$\mathcal{E}(f_t; \pi_{f_t}, h) \geq \epsilon/H$		

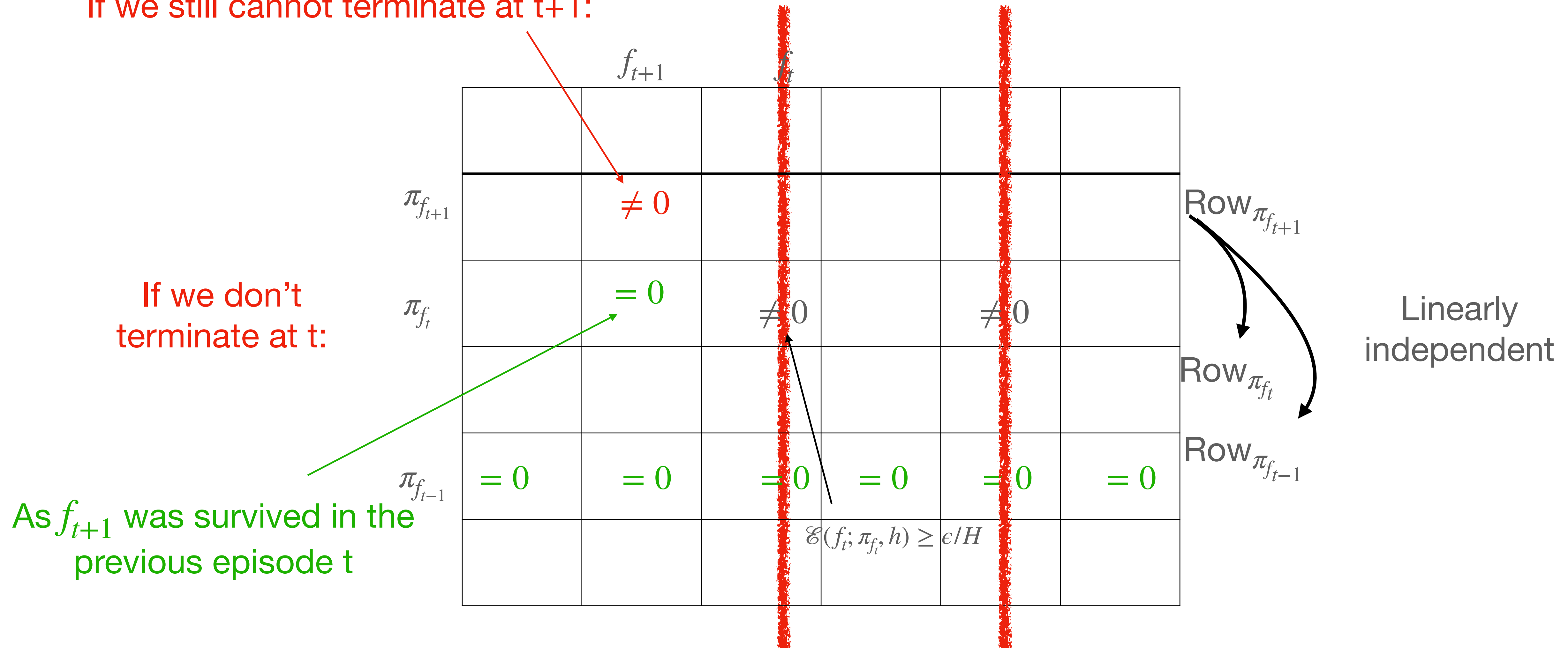
If we don't terminate at t:

As f_{t+1} was survived in the previous episode t

Progress on Value Function Elimination

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \mathcal{E}(f; \pi_f, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

If we still cannot terminate at t+1:



Progress on Value Function Elimination

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \mathcal{E}(f; \pi_f, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

$$\forall h : \mathcal{E}_h \in \mathbb{R}^{|\mathcal{F}| \times |\mathcal{F}|}$$

f

π_f		$\mathcal{E}_{f; \pi_f, h}$			

Progress on Value Function Elimination

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \mathcal{E}(f; \pi_f, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

$$\forall h : \mathcal{E}_h \in \mathbb{R}^{|\mathcal{F}| \times |\mathcal{F}|}$$

f

π_f		$\mathcal{E}_{f; \pi_f, h}$			

Every episode, we identify a row that is linearly independent of all previous rows we found!

Progress on Value Function Elimination

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \mathcal{E}(f; \pi_f, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

$$\forall h : \mathcal{E}_h \in \mathbb{R}^{|\mathcal{F}| \times |\mathcal{F}|}$$

f

π_f		$\mathcal{E}_{f; \pi_f, h}$			

Every episode, we identify a row that is linearly independent of all previous rows we found!

Then we must terminate in # of iterations at most (**Rank H**)

OLIVE Revisit (ignoring statistical error for simplicity)

Initialize $\mathcal{F}_0 = \mathcal{F}$

For $t = 0, \dots$

$$f_t = \arg \max_{f \in \mathcal{F}_t} \max_a f(x_0, a)$$

$$\text{If } \left| V^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \leq \epsilon, \text{ return } \pi_{f_t}$$

Version space update:

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \mathcal{E}(f; \pi_{f_t}, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

1. Upon termination we succeed
2. If not terminate, we make non-trivial progress
3. Total # of such non-trivial progress is bounded

OLIVE Revisit (ignoring statistical error for simplicity)

Initialize $\mathcal{F}_0 = \mathcal{F}$

1. Requires samples (poly in $1/\epsilon$)

For $t = 0, \dots$

$$f_t = \arg \max_{f \in \mathcal{F}_t} \max_a f(x_0, a)$$

$$\text{If } \left| V^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \leq \epsilon, \text{ return } \pi_{f_t}$$

Version space update:

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \mathcal{E}(f; \pi_{f_t}, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

1. Upon termination we succeed

2. If not terminate, we make non-trivial progress

3. Total # of such non-trivial progress is bounded

OLIVE Revisit (ignoring statistical error for simplicity)

Initialize $\mathcal{F}_0 = \mathcal{F}$

1. Requires samples (poly in $1/\epsilon$)

For $t = 0, \dots$

$$f_t = \arg \max_{f \in \mathcal{F}_t} \max_a f(x_0, a)$$

$$\text{If } \left| V^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \leq \epsilon, \text{ return } \pi_{f_t}$$

Version space update:

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \mathcal{E}(f; \pi_{f_t}, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

2. Requires samples (poly in $1/\epsilon, \ln(|\mathcal{F}|)$)
(needs to hold for all f)

1. Upon termination we succeed

2. If not terminate, we make non-trivial progress

3. Total # of such non-trivial progress is bounded

OLIVE Revisit (ignoring statistical error for simplicity)

Initialize $\mathcal{F}_0 = \mathcal{F}$

1. Requires samples (poly in $1/\epsilon$)

For $t = 0, \dots$

$$f_t = \arg \max_{f \in \mathcal{F}_t} \max_a f(x_0, a)$$

$$\text{If } \left| V^{\pi_{f_t}} - \max_a f_t(x_0, a) \right| \leq \epsilon, \text{ return } \pi_{f_t}$$

Version space update:

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F}_t : \mathcal{E}(f; \pi_{f_t}, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

2. Requires samples (poly in $1/\epsilon, \ln(|\mathcal{F}|)$)
(needs to hold for all f)

1. Upon termination we succeed

2. If not terminate, we make non-trivial progress

3. Total # of such non-trivial progress is bounded

$$\text{Poly} \left(H, \frac{1}{\epsilon}, \text{Rank}, \ln(|\mathcal{F}|) \right)$$

Low Bellman Rank Example

$$\mathcal{E}(f; \bar{f}, h) = \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right] = \langle \xi(f), \eta(\bar{f}) \rangle$$

Low Bellman Rank Example

$$\mathcal{E}(f; \bar{f}, h) = \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right] = \langle \xi(f), \eta(\bar{f}) \rangle$$

1. Tabular MDP:

Low Bellman Rank Example

$$\mathcal{E}(f; \bar{f}, h) = \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right] = \langle \xi(f), \eta(\bar{f}) \rangle$$

1. Tabular MDP:

$$\mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right] = \sum_{x, a \in X_h \times A_h} d_h^{\pi_{\bar{f}}}(s, a) \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right]$$

Low Bellman Rank Example

$$\mathcal{E}(f; \bar{f}, h) = \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right] = \langle \xi(f), \eta(\bar{f}) \rangle$$

1. Tabular MDP:

$$\begin{aligned} \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right] &= \sum_{x, a \in X_h \times A_h} d_h^{\pi_{\bar{f}}}(s, a) \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right] \\ &= \langle d_h^{\pi_{\bar{f}}}, f(\cdot, \cdot) - r(\cdot, \cdot) - \mathbb{E}_{x_{h+1} \sim P(\cdot | \cdot, \cdot)} \max_a f(x_{h+1}, a) \rangle \end{aligned}$$

Low Bellman Rank Example

$$\mathcal{E}(f; \bar{f}, h) = \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right] = \langle \xi(f), \eta(\bar{f}) \rangle$$

1. Tabular MDP:

$$\begin{aligned} \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right] &= \sum_{x, a \in X_h \times A_h} d_h^{\pi_{\bar{f}}}(s, a) \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right] \\ &= \langle d_h^{\pi_{\bar{f}}}, f(\cdot, \cdot) - r(\cdot, \cdot) - \mathbb{E}_{x_{h+1} \sim P(\cdot | \cdot, \cdot)} \max_a f(x_{h+1}, a) \rangle \end{aligned}$$

Rank at most $|X| |A|$

2. Linear MDPs:

$$P(\cdot | x, a) = \mu^\star \phi(x, a), \quad r(x, a) = \theta^\star \phi(x, a)$$

$$Q^\star(x, a) = (w^\star)^\top \phi(x, a)$$

2. Linear MDPs:

$$P(\cdot | x, a) = \mu^\star \phi(x, a), \quad r(x, a) = \theta^\star \phi(x, a)$$

$$Q^\star(x, a) = (w^\star)^\top \phi(x, a)$$

$$\mathcal{F} = \{w^\top \cdot \phi(x, a) : w \in \mathbb{R}^d, \|w\|_2 \leq W\}$$

2. Linear MDPs:

$$P(\cdot | x, a) = \mu^\star \phi(x, a), \quad r(x, a) = \theta^\star \phi(x, a)$$

$$Q^\star(x, a) = (w^\star)^\top \phi(x, a)$$

$$\mathcal{F} = \{w^\top \cdot \phi(x, a) : w \in \mathbb{R}^d, \|w\|_2 \leq W\}$$

$$\mathbb{E}_{x_h, a_h \sim d_h^{\pi_f}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right]$$

2. Linear MDPs:

$$P(\cdot | x, a) = \mu^\star \phi(x, a), \quad r(x, a) = \theta^\star \phi(x, a)$$

$$Q^\star(x, a) = (w^\star)^\top \phi(x, a)$$

$$\mathcal{F} = \{w^\top \cdot \phi(x, a) : w \in \mathbb{R}^d, \|w\|_2 \leq W\}$$

$$\begin{aligned} & \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right] \\ &= \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[w^\top \phi(x_h, a_h) - (\theta^\star)^\top \phi(x_h, a_h) - \phi(x_h, a_h)^\top (\mu^\star)^\top \left(\max_a f(\cdot, a) \right) \right] \end{aligned}$$

2. Linear MDPs:

$$P(\cdot | x, a) = \mu^\star \phi(x, a), \quad r(x, a) = \theta^\star \phi(x, a)$$

$$Q^\star(x, a) = (w^\star)^\top \phi(x, a)$$

$$\mathcal{F} = \{w^\top \cdot \phi(x, a) : w \in \mathbb{R}^d, \|w\|_2 \leq W\}$$

$$\begin{aligned} & \mathbb{E}_{x_h, a_h \sim d_h^{\pi_f}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right] \\ &= \mathbb{E}_{x_h, a_h \sim d_h^{\pi_f}} \left[w^\top \phi(x_h, a_h) - (\theta^\star)^\top \phi(x_h, a_h) - \phi(x_h, a_h)^\top (\mu^\star)^\top \left(\max_a f(\cdot, a) \right) \right] \\ &= \mathbb{E}_{x_h, a_h \sim d_h^{\pi_f}} \left[\left(w + \theta^\star + (\mu^\star)^\top \left(\max_a f(\cdot, a) \right) \right)^\top \phi(x_h, a_h) \right] \end{aligned}$$

2. Linear MDPs:

$$P(\cdot | x, a) = \mu^\star \phi(x, a), \quad r(x, a) = \theta^\star \phi(x, a)$$

$$Q^\star(x, a) = (w^\star)^\top \phi(x, a)$$

$$\mathcal{F} = \{w^\top \cdot \phi(x, a) : w \in \mathbb{R}^d, \|w\|_2 \leq W\}$$

$$\begin{aligned} & \mathbb{E}_{x_h, a_h \sim d_h^{\pi_f}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right] \\ &= \mathbb{E}_{x_h, a_h \sim d_h^{\pi_f}} \left[w^\top \phi(x_h, a_h) - (\theta^\star)^\top \phi(x_h, a_h) - \phi(x_h, a_h)^\top (\mu^\star)^\top \left(\max_a f(\cdot, a) \right) \right] \\ &= \mathbb{E}_{x_h, a_h \sim d_h^{\pi_f}} \left[\left(w + \theta^\star + (\mu^\star)^\top \left(\max_a f(\cdot, a) \right) \right)^\top \phi(x_h, a_h) \right] \\ &= \left\langle \left(w + \theta^\star + (\mu^\star)^\top \left(\max_a f(\cdot, a) \right) \right), \mathbb{E}_{x_h, a_h \sim d_h^{\pi_f}} [\phi(x_h, a_h)] \right\rangle \end{aligned}$$

2. Linear MDPs:

$$P(\cdot | x, a) = \mu^\star \phi(x, a), \quad r(x, a) = \theta^\star \phi(x, a)$$

$$Q^\star(x, a) = (w^\star)^\top \phi(x, a)$$

$$\mathcal{F} = \{w^\top \cdot \phi(x, a) : w \in \mathbb{R}^d, \|w\|_2 \leq W\}$$

$$\begin{aligned} & \mathbb{E}_{x_h, a_h \sim d_h^{\pi_f}} \left[f(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} f(x_{h+1}, a) \right] \right] \\ &= \mathbb{E}_{x_h, a_h \sim d_h^{\pi_f}} \left[w^\top \phi(x_h, a_h) - (\theta^\star)^\top \phi(x_h, a_h) - \phi(x_h, a_h)^\top (\mu^\star)^\top \left(\max_a f(\cdot, a) \right) \right] \\ &= \mathbb{E}_{x_h, a_h \sim d_h^{\pi_f}} \left[\left(w + \theta^\star + (\mu^\star)^\top \left(\max_a f(\cdot, a) \right) \right)^\top \phi(x_h, a_h) \right] \\ &= \left\langle \left(w + \theta^\star + (\mu^\star)^\top \left(\max_a f(\cdot, a) \right) \right), \mathbb{E}_{x_h, a_h \sim d_h^{\pi_f}} [\phi(x_h, a_h)] \right\rangle \end{aligned}$$

Rank at most d

2. Linear Function with Bellman Completeness

$$\mathcal{F} = \{w^\top \cdot \phi(x, a) : w \in \mathbb{R}^d, \|w\|_2 \leq W\}$$

$$r(x, a) = \theta^\star \phi(x, a),$$

and for any linear function $f(x, a) := w^\top \phi(x, a)$,

$$\text{we have } \mathcal{T}f(x, a) = (w')^\top \phi(x, a)$$

2. Linear Function with Bellman Completeness

$$\mathcal{F} = \{w^\top \cdot \phi(x, a) : w \in \mathbb{R}^d, \|w\|_2 \leq W\}$$

$$r(x, a) = \theta^\star \phi(x, a),$$

and **for any linear function** $f(x, a) := w^\top \phi(x, a)$,

$$\text{we have } \mathcal{T}f(x, a) = (w')^\top \phi(x, a)$$

(Go and verify that linear MDP is a special instance of this setting)

2. Linear Function with Bellman Completeness

$$\mathcal{F} = \{w^\top \cdot \phi(x, a) : w \in \mathbb{R}^d, \|w\|_2 \leq W\}$$

$$r(x, a) = \theta^\star \phi(x, a),$$

and for any linear function $f(x, a) := w^\top \phi(x, a)$,
we have $\mathcal{T}f(x, a) = (w')^\top \phi(x, a)$

(Go and verify that linear MDP is a special instance of this setting)

$$\mathbb{E}_{x_h, a_h \sim d_h^{\pi_f}} \left[w^\top \phi(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} w^\top \phi(x_{h+1}, a) \right] \right]$$

2. Linear Function with Bellman Completeness

$$\mathcal{F} = \{w^\top \cdot \phi(x, a) : w \in \mathbb{R}^d, \|w\|_2 \leq W\}$$

$$r(x, a) = \theta^\star \phi(x, a),$$

and for any linear function $f(x, a) := w^\top \phi(x, a)$,
we have $\mathcal{T}f(x, a) = (w')^\top \phi(x, a)$

(Go and verify that linear MDP is a special instance of this setting)

$$\begin{aligned} & \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[w^\top \phi(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} w^\top \phi(x_{h+1}, a) \right] \right] \\ &= \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[w^\top \phi(x_h, a_h) - (w')^\top \phi(x_h, a_h) \right] \end{aligned}$$

2. Linear Function with Bellman Completeness

$$\mathcal{F} = \{w^\top \cdot \phi(x, a) : w \in \mathbb{R}^d, \|w\|_2 \leq W\}$$

$$r(x, a) = \theta^\star \phi(x, a),$$

and for any linear function $f(x, a) := w^\top \phi(x, a)$,
we have $\mathcal{T}f(x, a) = (w')^\top \phi(x, a)$

(Go and verify that linear MDP is a special instance of this setting)

$$\mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[w^\top \phi(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} w^\top \phi(x_{h+1}, a) \right] \right]$$

$$= \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[w^\top \phi(x_h, a_h) - (w')^\top \phi(x_h, a_h) \right]$$

$$= \left\langle \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[\phi(x_h, a_h) \right], w - w' \right\rangle$$

2. Linear Function with Bellman Completeness

$$\mathcal{F} = \{w^\top \cdot \phi(x, a) : w \in \mathbb{R}^d, \|w\|_2 \leq W\}$$

$$r(x, a) = \theta^\star \phi(x, a),$$

and for any linear function $f(x, a) := w^\top \phi(x, a)$,
we have $\mathcal{T}f(x, a) = (w')^\top \phi(x, a)$

(Go and verify that linear MDP is a special instance of this setting)

$$\mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[w^\top \phi(x_h, a_h) - r(x_h, a_h) - \mathbb{E}_{x_{h+1} \sim P(\cdot | x_h, a_h)} \left[\max_{a \in \mathcal{A}} w^\top \phi(x_{h+1}, a) \right] \right]$$

$$= \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[w^\top \phi(x_h, a_h) - (w')^\top \phi(x_h, a_h) \right]$$

Rank at most d

$$= \left\langle \mathbb{E}_{x_h, a_h \sim d_h^{\pi_{\bar{f}}}} \left[\phi(x_h, a_h) \right], w - w' \right\rangle$$

Other Examples it captures:

Low rank MDP (requires a small modification in definition)

Reactive Predictive State Representation (PSRs)

Reactive POMDP

Examples it does not capture:

Factored MDPs

Linear Q^* (?)

Model-based RL: Function Approximation in Model

Let's set up function class in Model-based RL Setting

We consider a model class \mathcal{P}

$$\mathcal{P} \subset X \times A \mapsto \Delta(X)$$

Let's set up function class in Model-based RL Setting

We consider a model class \mathcal{P}

$$\mathcal{P} \subset X \times A \mapsto \Delta(X)$$

1. **Realizability** assumption:

$$P^* \in \mathcal{P}$$

Let's set up function class in Model-based RL Setting

We consider a model class \mathcal{P}

$$\mathcal{P} \subset X \times A \mapsto \Delta(X)$$

1. **Realizability** assumption:

$$P^* \in \mathcal{P}$$

2. We also assume reward function r is known

Let's set up function class in Model-based RL Setting

We consider a model class \mathcal{P}

$$\mathcal{P} \subset X \times A \mapsto \Delta(X)$$

1. **Realizability** assumption:

$$P^* \in \mathcal{P}$$

2. We also assume reward function r is known

3. Define **Optimal-Planner (OP)**: $\text{OP}(P, r) = \succ (\pi_P^*, V_P^*, Q_P^*)$

Let's set up function class in Model-based RL Setting

We consider a model class \mathcal{P}

$$\mathcal{P} \subset X \times A \mapsto \Delta(X)$$

1. **Realizability** assumption:

$$P^* \in \mathcal{P}$$

2. We also assume reward function r is known

Computational oracle, no
real-world samples needed

3. Define **Optimal-Planner (OP)**: $OP(P, r) = \succ (\pi_P^*, V_P^*, Q_P^*)$

Let's set up function class in Model-based RL Setting

We consider a model class \mathcal{P}

$$\mathcal{P} \subset X \times A \mapsto \Delta(X)$$

1. **Realizability** assumption:

$$P^* \in \mathcal{P}$$

2. We also assume reward function r is known

Computational oracle, no
real-world samples needed

3. Define **Optimal-Planner (OP)**: $OP(P, r) = \succ (\pi_P^*, V_P^*, Q_P^*)$

4. Induced Policy, Value function, and Q function class:

$$\Pi := \{\pi_P^* : P \in \mathcal{P}\}, \mathcal{V} = \{V_P^* : P \in \mathcal{P}\}, \mathcal{Q} = \{Q_P^* : P \in \mathcal{P}\}$$

Let's set up function class in Model-based RL Setting

We consider a model class \mathcal{P}

$$\mathcal{P} \subset X \times A \mapsto \Delta(X)$$

1. **Realizability** assumption:

$$P^* \in \mathcal{P}$$

2. We also assume reward function r is known

Computational oracle, no
real-world samples needed

3. Define **Optimal-Planner (OP)**: $OP(P, r) = \succ (\pi_P^*, V_P^*, Q_P^*)$

4. Induced Policy, Value function, and Q function class:

$$\Pi := \{\pi_P^* : P \in \mathcal{P}\}, \mathcal{V} = \{V_P^* : P \in \mathcal{P}\}, \mathcal{Q} = \{Q_P^* : P \in \mathcal{P}\}$$

Indeed you can
run OLIVE w/ \mathcal{Q} !

How to check if a model approximator is good?

Witness model error: $\mathcal{E}(P; \pi, h) = \max_{f \in \mathcal{F}} \mathbb{E}_{x_h, a_h \sim d_h^\pi} \left[\mathbb{E}_{x_{h+1} \sim P_{x_h, a_h}} f(x_h, a_h, x_{h+1}) - \mathbb{E}_{x_{h+1} \sim P_{x_h, a_h}^*} f(x_h, a_h, x_{h+1}) \right]$

Witness function (or aka discriminators): $\mathcal{F} \subset X \times A \times X \mapsto \mathbb{R}$

How to check if a model approximator is good?

$$\mathcal{E}(P^*; \pi, h) = 0, \forall \pi, h$$

Witness model error: $\mathcal{E}(P; \pi, h) = \max_{f \in \mathcal{F}} \mathbb{E}_{x_h, a_h \sim d_h^\pi} \left[\mathbb{E}_{x_{h+1} \sim P_{x_h, a_h}} f(x_h, a_h, x_{h+1}) - \mathbb{E}_{x_{h+1} \sim P_{x_h, a_h}^*} f(x_h, a_h, x_{h+1}) \right]$

Witness function (or aka discriminators): $\mathcal{F} \subset X \times A \times X \mapsto \mathbb{R}$

How to check if a model approximator is good?

$$\mathcal{E}(P^*; \pi, h) = 0, \forall \pi, h$$

Witness model error: $\mathcal{E}(P; \pi, h) = \max_{f \in \mathcal{F}} \mathbb{E}_{x_h, a_h \sim d_h^\pi} \left[\mathbb{E}_{x_{h+1} \sim P_{x_h, a_h}} f(x_h, a_h, x_{h+1}) - \mathbb{E}_{x_{h+1} \sim P_{x_h, a_h}^*} f(x_h, a_h, x_{h+1}) \right]$

Witness function (or aka discriminators): $\mathcal{F} \subset X \times A \times X \mapsto \mathbb{R}$

Integral Probability Metric (IPM): given two distributions $P_1 \in \Delta(X), P_2 \in \Delta(X)$:

$$\text{IPM}_{\mathcal{F}} = \max_{f \in \mathcal{F}} \mathbb{E}_{x \sim P_1} f(x) - \mathbb{E}_{x \sim P_2} f(x)$$

How to check if a model approximator is good?

$$\mathcal{E}(P^*; \pi, h) = 0, \forall \pi, h$$

Witness model error: $\mathcal{E}(P; \pi, h) = \max_{f \in \mathcal{F}} \mathbb{E}_{x_h, a_h \sim d_h^\pi} \left[\mathbb{E}_{x_{h+1} \sim P_{x_h, a_h}} f(x_h, a_h, x_{h+1}) - \mathbb{E}_{x_{h+1} \sim P_{x_h, a_h}^*} f(x_h, a_h, x_{h+1}) \right]$

Witness function (or aka discriminators): $\mathcal{F} \subset X \times A \times X \mapsto \mathbb{R}$

Integral Probability Metric (IPM): given two distributions $P_1 \in \Delta(X), P_2 \in \Delta(X)$:

$$\text{IPM}_{\mathcal{F}} = \max_{f \in \mathcal{F}} \mathbb{E}_{x \sim P_1} f(x) - \mathbb{E}_{x \sim P_2} f(x)$$

$\mathcal{F} = \{f : \|f\|_\infty \leq 1\} \Rightarrow$ Total Variation Distance;

$\mathcal{F} = \{f : f \text{ is 1-Lipschitz}\} \Rightarrow$ Wasserstein Distance;

How to check if a model approximator is good?

$$\mathcal{E}(P^*; \pi, h) = 0, \forall \pi, h$$

Witness model error:
$$\mathcal{E}(P; \pi, h) = \max_{f \in \mathcal{F}} \mathbb{E}_{x_h, a_h \sim d_h^\pi} \left[\mathbb{E}_{x_{h+1} \sim P_{x_h, a_h}} f(x_h, a_h, x_{h+1}) - \mathbb{E}_{x_{h+1} \sim P_{x_h, a_h}^*} f(x_h, a_h, x_{h+1}) \right]$$

Witness function (or aka discriminators): $\mathcal{F} \subset X \times A \times X \mapsto \mathbb{R}$

Integral Probability Metric (IPM): given two distributions $P_1 \in \Delta(X), P_2 \in \Delta(X)$:

$$\text{IPM}_{\mathcal{F}} = \max_{f \in \mathcal{F}} \mathbb{E}_{x \sim P_1} f(x) - \mathbb{E}_{x \sim P_2} f(x)$$

$\mathcal{F} = \{f : \|f\|_\infty \leq 1\} \Rightarrow$ Total Variation Distance;

$\mathcal{F} = \{f : f \text{ is 1-Lipschitz}\} \Rightarrow$ Wasserstein Distance;

Assumption (for analysis simplicity): $\mathcal{V} \subseteq \mathcal{F}$, where recall $\mathcal{V} = \{V_P^* : P \in \mathcal{P}\}$

Optimism Led Iterative Model Elimination (OLIM)

Initialize $\mathcal{P}_0 = \mathcal{P}$

For $t = 0, \dots$

$$P_t = \arg \max_{P \in \mathcal{P}_t} V_P^\star(x_0)$$

$$\pi_t := \pi_{P_t}^\star$$

If $\left| V^{\pi_t} - V_{P_t}^\star \right| \leq \epsilon$, return π_t

Version space update:

$$\mathcal{P}_{t+1} = \left\{ P \in \mathcal{P}_t : \mathcal{E}(P; \pi_t, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

Optimism Led Iterative Model Elimination (OLIM)

Initialize $\mathcal{P}_0 = \mathcal{P}$

For $t = 0, \dots$

$$P_t = \arg \max_{P \in \mathcal{P}_t} V_P^\star(x_0)$$

$$\pi_t := \pi_{P_t}^\star$$

If $\left| V^{\pi_t} - V_{P_t}^\star \right| \leq \epsilon$, return π_t

Version space update:

$$\mathcal{P}_{t+1} = \left\{ P \in \mathcal{P}_t : \mathcal{E}(P; \pi_t, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

1. Terminate means we succeed

Optimism Led Iterative Model Elimination (OLIM)

Initialize $\mathcal{P}_0 = \mathcal{P}$

For $t = 0, \dots$

$$P_t = \arg \max_{P \in \mathcal{P}_t} V_P^\star(x_0)$$

$$\pi_t := \pi_{P_t}^\star$$

If $\left| V^{\pi_t} - V_{P_t}^\star \right| \leq \epsilon$, return π_t

Version space update:

$$\mathcal{P}_{t+1} = \left\{ P \in \mathcal{P}_t : \mathcal{E}(P; \pi_t, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

1. Terminate means we succeed

2. If not terminate, we make progress

Optimism Led Iterative Model Elimination (OLIM)

Initialize $\mathcal{P}_0 = \mathcal{P}$

For $t = 0, \dots$

$$P_t = \arg \max_{P \in \mathcal{P}_t} V_P^\star(x_0)$$

$$\pi_t := \pi_{P_t}^\star$$

If $\left| V^{\pi_t} - V_{P_t}^\star \right| \leq \epsilon$, return π_t

Version space update:

$$\mathcal{P}_{t+1} = \left\{ P \in \mathcal{P}_t : \mathcal{E}(P; \pi_t, h) = 0, \forall h \in \{0, 1, \dots, H-1\} \right\}$$

1. Terminate means we succeed
2. If not terminate, we make progress
3. Total progress is upper bounded

Progress on Model Elimination

MDP + \mathcal{P} + \mathcal{F} determines the following matrix: $\mathcal{E} \in \mathbb{R}^{|\mathcal{P}| \times |\mathcal{P}|}$

... P ...

.									
.									
.									
π_P^\star			$\mathcal{E}_{P; \pi_P^\star, h}$						

Rank of this Matrix = Witness Rank

Sample complexity of Optimism Led Iterative Model Elimination:

Under some assumption of discriminators \mathcal{F} (i.e., $\mathcal{V} \in \mathcal{F}$),
Witness rank \leq Bellman Rank (\mathcal{Q})

Sample complexity of Optimism Led Iterative Model Elimination:

$$\text{Poly} \left(H, \frac{1}{\epsilon}, \text{Witness-Rank}, \ln(|\mathcal{P}| |\mathcal{F}|) \right)$$

Under some assumption of discriminators \mathcal{F} (i.e., $\mathcal{V} \in \mathcal{F}$),
Witness rank \leq Bellman Rank (\mathcal{Q})

Comparison of Witness Rank and Bellman Rank (More Broadly, Model-based Versus Model-free)

Comparison of Witness Rank and Bellman Rank (More Broadly, Model-based Versus Model-free)

Let us compare under a specific setup:

We start with model class \mathcal{P} , we convert models to Q functions $\mathcal{Q} = \{Q_P^\star : P \in \mathcal{P}\}$

Comparison of Witness Rank and Bellman Rank (More Broadly, Model-based Versus Model-free)

Let us compare under a specific setup:

We start with model class \mathcal{P} , we convert models to Q functions $\mathcal{Q} = \{Q_P^* : P \in \mathcal{P}\}$

$\mathcal{P} \Rightarrow$ Witness Rank and OLIME

$\mathcal{Q} \Rightarrow$ Bellman Rank and OLIVE

Comparison of Witness Rank and Bellman Rank (More Broadly, Model-based Versus Model-free)

Let us compare under a specific setup:

We start with model class \mathcal{P} , we convert models to Q functions $\mathcal{Q} = \{Q_P^* : P \in \mathcal{P}\}$

$\mathcal{P} \Rightarrow$ Witness Rank and OLIME

$\mathcal{Q} \Rightarrow$ Bellman Rank and OLIVE

Theorem[Exponential Separation]:

\exists MDPs (Factored MDP!) and realizable model class \mathcal{P} , s.t:

Witness-Rank(\mathcal{P}) is exponentially smaller than Bellman-rank(\mathcal{Q}) in horizon H