

Batch RL + Fitted Q-iteration

- a) Coverage Assumptions
- b) Function Approx Assumptions

MPP: $S, A, P, R: S \times A \rightarrow [0, 1], \gamma \in [0, 1], d_0 \in \Delta(S)$

$V_{max} := \frac{1}{1-\gamma}$

Data: $\{(x_i, a_i, r_i, x'_i)\}_{i=1}^n \quad (x_i, a_i) \overset{i.i.d.}{\sim} \mu \in \Delta(S \times A)$
 $r_i \sim R(x_i, a_i), x'_i \sim P(\cdot | x_i, a_i)$

Goal: Output $\hat{\pi}$ s.t. with prob $1-\delta$: $V(\hat{\pi}) \geq V^* - \epsilon$

Example: $\gamma=0, d_0 = \delta_{x_0} \Rightarrow$ MAB problem

$\rightarrow \mu = d_0 \circ \delta_{a^*}$, cannot learn reward for $a' \neq a^*$
 \Rightarrow Cannot find π^*
 \Rightarrow Must make coverage assumption

Assmp 1: (Unif concentrability) $\exists C_{unif} < \infty$ s.t. $\forall x, a: \frac{1}{\mu(x, a)} \leq C_{unif}$
 $\Rightarrow |S| \times |A| \leq C_{unif} \Rightarrow$ Tabular Setting

Prop 1: Under Assmp 1, \exists an alg s.t.

$V^* - V(\hat{\pi}) \leq \tilde{O} \left(\frac{1}{(1-\gamma)^2} \left(\sqrt{\frac{C_{unif}}{n}} + V_{max} \gamma (|S| \sqrt{\frac{C_{unif}}{n}}) \right) \right)$

(Note: $\sqrt{\frac{C_{unif}}{n}}$ is labeled $= \Delta^2$ and $V_{max} \gamma (|S| \sqrt{\frac{C_{unif}}{n}})$ is labeled $= \Delta^1$)

\uparrow hiding $\log(|S||A|/\delta)$ factors

Sketch: n samples, $\frac{n}{\text{Conf}}$ samples from each (x, a)

$$\hat{R}(x, a) = \frac{\sum_i \mathbb{1}\{(x_i, a_i) = (x, a)\} r_i}{\sum_i \mathbb{1}\{(x_i, a_i) = (x, a)\}}$$

$$\hat{P}(x' | x, a) = \frac{\sum_i \mathbb{1}\{(x_i, a_i, x'_i) = (x, a, x')\}}{\sum_i \mathbb{1}\{(x_i, a_i) = (x, a)\}}$$

Concentration

$$\forall x, a: |\hat{R}(x, a) - R(x, a)| \leq \sqrt{\frac{\text{Conf}}{n}} =: \Delta_R \leftarrow \text{Hoeffding + Union}$$

$$\forall x, a: \|\hat{P}(\cdot | x, a) - P(\cdot | x, a)\|_{TV} \leq \sqrt{\frac{\text{Conf}}{n}} =: \Delta_P$$

$\hat{\pi}$, opt in (\hat{P}, \hat{R}) \rightarrow PDL

$$V^* - V(\hat{\pi}) = \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{d_{h, \hat{\pi}}} \left[Q^*(x, \pi^*(x)) - Q^*(x, \hat{\pi}(x)) \right]$$

- $\oplus \hat{Q}(x, \pi^*(x)) \leftarrow$
- $\ominus \hat{Q}(x, \hat{\pi}(x)) \leftarrow$
- $\hat{Q}(x, \hat{\pi}(x)) \geq \hat{Q}(x, \pi^*(x))$

$$\approx \sum_h \gamma^h \mathbb{E}_{d_{h, \hat{\pi}}} \left[|Q^*(x, \pi^*(x)) - \hat{Q}(x, \pi^*(x))| + |\hat{Q}(x, \hat{\pi}(x)) - Q^*(x, \hat{\pi}(x))| \right]$$

$$\mathbb{E}_v \left[|Q^*(x, a) - \hat{Q}(x, a)| \right] \leq \mathbb{E}_v \left[|R(x, a) - \hat{R}(x, a)| \right] \leq \Delta_R$$

$\hat{Q}(x, a) = R + \hat{P}V(x)$
 $\hat{Q}(x, a) = R + \hat{P}\hat{V}$

$$\mathbb{E}_v \left[\gamma \left| \mathbb{E}_{x' \sim P_{x, a}} V^*(x') - \mathbb{E}_{x' \sim \hat{P}_{x, a}} \hat{V}(x') \right| \right]$$

$$\mathbb{E}_v \left[\left| \mathbb{E}_{x' \sim P_{x, a}} V^*(x') - \hat{V}(x) \right| \right] + \gamma \mathbb{E}_v \left[\left| \mathbb{E}_{\hat{P}} \hat{V}(x') - \mathbb{E}_{\hat{P}} \hat{V}(x') \right| \right]$$

$$P \mathbb{E}_{x' \sim P_{x'}} [| \mathbb{E} V^*(x') - \hat{V}(x') | + \gamma \mathbb{E}_{\hat{P}} [| \mathbb{E} \hat{V}(x') - \hat{V}(x') |]$$

$$V^*(x) = \max_a Q^*(x, a)$$

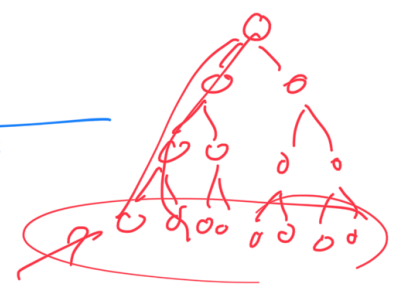
$$\hat{V}(x) = \max_a \hat{Q}(x, a)$$

$$| V^*(x) - \hat{V}(x) | \leq | Q^*(x, \hat{\pi}(x)) - \hat{Q}(x, \hat{\pi}(x)) |$$

$$\text{error} \leq \Delta_P + V_{\max} \delta \Delta_P + \gamma \cdot \text{error} \Rightarrow \text{error} \leq \frac{1}{1-\gamma} (\Delta_P + V_{\max} \delta \Delta_P)$$

Assmp 2: (Concentrability) $\exists C < \infty$ s.t.,

$$\forall \text{ nonstationary } \pi, h, x, a \quad \frac{d_h^\pi(x, a)}{\mu(x, a)} \leq C$$



$$R: S \times A \rightarrow [0, V_{\max}]$$

FQI: Initialize $f_0 \in F$ arbitrarily

Iterate:

$$f_k \leftarrow \underset{f \in F}{\text{argmin}} \sum_{i=1}^n \left(f(x_i, a_i) - \underbrace{r_i - \gamma \max_{a'} f_{k-1}(x_i, a')}_{\text{GP}}$$

Output π_{f_k} for the best K

Bellman op: $\mathcal{T}f: (x, a) \mapsto \mathbb{E} [r + \gamma \max_{a'} f(x', a') | x, a]$

Assmp 3: (Completeness) $\mathcal{T}F \subseteq F$

$$\hookrightarrow = \{ \mathcal{T}f: f \in F \}$$

$$\Rightarrow Q^* \in F$$

$$Q^* = \lim_{t \rightarrow \infty} \mathcal{T}^t f_0 \quad \text{for any } f_0$$

$|F| < \infty$

Thm: The k^{th} iterate of FQI satisfies w.p. $\geq 1-\delta$

$$V^* - V(\pi_k) \leq O \left(\frac{V_{\max}}{(1-\gamma)^2} \left(\frac{C \log(|F|/\delta)}{1-\gamma} \right) + \frac{2\gamma^k V_{\max}}{(1-\gamma)} \right)$$

Thm: The k^{th} iterate of FQI satisfies w.p. $\geq 1-\delta$

$$V^* - V(\pi_k) \leq O\left(\frac{V_{max}}{(1-\gamma)^2} \sqrt{\frac{C \log(|F|/\delta)}{n}}\right) + \frac{2\gamma^k V_{max}}{(1-\gamma)}$$

$\rightarrow 0$ as $n \rightarrow \infty$ opt. term

Proof: $\|f\|_v^2 = \mathbb{E}_{x \sim \nu} f(x,a)^2 \quad v \in \Delta(S \times A)$

Concentration: w.p. $1-\delta$

$$\forall k \quad \|f_k - \mathbb{T}f_{k-1}\|_{\mu}^2 \leq O\left(\frac{V_{max}^2 \log(|F|/\delta)}{n}\right) =: \Delta_{SQ}$$

$$V^* - V(\pi_k) = \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{d_h, \pi_k} [Q^*(x, \pi_k(x)) - Q^+(x, \pi_k(x))]$$

$$\leq \sum_{h=0}^{\infty} \gamma^h \|Q^* - f_k\|_{2, d_h, \pi_k, \pi_k^*} + \gamma^h \|Q^* - f_k\|_{2, d_h, \pi_k, \pi_k}$$

$$\|Q^* - f_k\|_{2, \nu} \leq \|Q^* - \mathbb{T}f_{k-1}\|_{2, \nu} + \|\mathbb{T}f_{k-1} - f_k\|_{2, \nu}$$

Assump: $\forall x,a, \frac{v(x,a)}{\mu(x)} \leq C$

$$\leq \gamma \sqrt{\mathbb{E}_v \left[\left| \mathbb{E}_{x \sim \nu} \max_a Q^*(x,a) - \max_a f_{k-1}(x,a) \right|^2 \right]} \leq \sqrt{C} \cdot \sqrt{\mathbb{E}_v \left[\left| \mathbb{E}_{x \sim \nu} \max_a Q^*(x,a) - \max_a f_{k-1}(x,a) \right|^2 \right]} = \sqrt{C} \|\mathbb{T}f_{k-1} - f_k\|_{2, \nu} \leq \sqrt{C \cdot \Delta_{SQ}}$$

$\pi = \text{argmax}_a \{ \max_{x \sim \nu} Q^*(x,a), \max_{x \sim \nu} f_{k-1}(x,a) \}$

$$\leq \gamma \|Q^* - f_{k-1}\|_{v, \pi}$$

$$\rightarrow \|Q^* - f_k\|_{2, \nu} \leq \frac{1}{1-\gamma} \sqrt{C \cdot \Delta_{SQ}} + \gamma^k V_{max}$$

Concentration:

Completeness

holds in Linear MDP

Concentrability:

If f is linear in ϕ , then

$$\|f\|_v^2 = f^T \Sigma_v f \quad \Sigma_v = \mathbb{E}_{x \sim \nu} \phi \phi^T$$

$$L \leq C \|f\|_{\mu}^2$$

$$f^T \Sigma_v f = f^T \sum_{\mu}^{-1/2} \sum_{\mu}^{-1/2} \sum_{\mu}^{-1/2} \sum_{\mu}^{-1/2} f$$

$$\text{tr} \left(\sum_{\mu}^{-1/2} \sum_{\mu}^{-1/2} \sum_{\mu}^{-1/2} \sum_{\mu}^{-1/2} f f^T \right) \leq \lambda_{\max} \left(\sum_{\mu}^{-1/2} \sum_{\mu}^{-1/2} \right) \cdot f^T \sum_{\mu} f$$

"Linear" Concentrability

$$\forall \pi, h, \lambda_{\max} \left(\sum_{\pi, h} \left(\sum_{\mu} \right) \right) \leq C_{lm}$$

Thm (Solk): linear completeness + linear concentrability \Leftarrow

$$V^* - V(\pi_k) \leq \text{poly} \left(C_{lm}, d, \frac{1}{n} \right) + O \left(\frac{r^k V_{\max}}{nr} \right)$$

Completeness

Very strong, holds in Linear MDP in LQR

Non-monotonic

$$G \geq F$$

$$\gamma F \leq F \Rightarrow \gamma G \leq G \Leftarrow$$

$G^* \in F$ is monotone

Realizability only: $[X_{1:n}, J_{1:n}]$

$$V^* - V(\hat{\pi}) \leq n^{-1/4} \cdot \log(F/S)$$

Very strong concentrability assumption

$$\mu(a|s) \geq \frac{1}{C_A} \quad \frac{P(s'|s, a)}{\mu(s')} \leq C_S$$

$\hookrightarrow \Rightarrow$ our concentrability assumption

Shan, Russo, Wg, Den Foster

Realizability only not enough w/ "weak" concentrability

\Leftarrow linear, "linear" concentrability

\rightarrow Need $\exp(H)$ samples $\exp(H)$

minimax approach

$$\text{argmin}_{f \in F} \max_{w \in W} |L(f, w)| \hookrightarrow \mathbb{E}_{\mu} \left[w(x_{t^c}) \left(\underbrace{f(x_{t^c})}_{d''(x_{t^c})} - \underbrace{r - r_{\min}}_a \underbrace{f(x_{t^c}, a)} \right) \right]$$

Horizon Dependence

$\frac{1}{(T-r)^2}$ dependence, $\frac{1}{nr}$ possible

Bellman Residual Minimization

minimax approach

$$\arg\min_{f \in \mathcal{F}} \max_{w \in \mathcal{W}} |L(f, w)|$$

$$\hookrightarrow = \mathbb{E}_{\mu} \left[w(x_i) \left(f(x_i) - \underbrace{r - r \max_{a'} f(x_i, a)}_{\frac{d''(x_i)}{\mu(x_i)}} \right) \right]$$

Horizon Dependence

$\frac{1}{(1-\gamma)^2}$ dependence, $\frac{1}{1-\gamma}$ possible

Bellman Residual Minimization

$$\arg\min_{f \in \mathcal{F}} \max_{g \in \mathcal{F}} \mathbb{E}_{\mu} \left(f - r - \gamma \max_{a'} f(x_i, a) \right)^2 - \mathbb{E}_{\mu} \left(g - r - \gamma \max_{a'} f \right)^2$$

noise level of 'f's problem'

$$V(\pi) - V(\pi_f) \leq \frac{1}{1-\gamma} \left(\mathbb{E}_{\mu} \left[\underbrace{\pi_f - f}_{f_L - T f_{L-1}} \right] + \mathbb{E}_{\mu} \left[f - \pi_f \right] \right)$$

What if we have poor coverage?

- Obviously cannot match V^*

Model based \Rightarrow realizability only is ok
(but model-based realizability is v. strong!)

Model-free \Rightarrow realizability only seems insufficient
(but much weaker)