

Data:  $(s, r) \stackrel{iid}{\sim} \mu$ ,  $r \sim R(s, a)$ ,  $s' \sim P(\cdot | (s, a))$   $\uparrow$   $\mathcal{F}_{k-1}$

FQI:  $f_k \leftarrow \underset{f \in \mathcal{F}}{\operatorname{argmin}} \sum_{(s, a, r, s')} (f(s, a) - r - \gamma \max_{a'} f(s', a'))^2$

Assumptions: ①  $\mathcal{T}f \in \mathcal{F}, \forall f \in \mathcal{F}$ .  $\uparrow$   $\pi = \pi_{f_k}$   
 ②  $\frac{d_t^\pi(s, a)}{\mu(s, a)} \leq C, \forall s, a, t, \pi$

FQE:  $f_k \leftarrow \underset{f \in \mathcal{F}}{\operatorname{argmin}} \sum_{(s, a, r, s')} (f(s, a) - r - \gamma \underbrace{f(s', \pi)}_{f_{k-1}})^2$

$J(\pi) = \mathbb{E}_{s \sim d_0} [V''(s)] = \mathbb{E}_{s \sim d_0} [Q''(s, \pi)]$

$\hat{J}(\pi) = \mathbb{E}_{s \sim d_0} [f_k(s, \pi)]$

Assum: ①  $\mathcal{T}''f \in \mathcal{F}, \forall f \in \mathcal{F}$

②  $\frac{d_t^\pi(s, a)}{\mu(s, a)} \leq C, \forall s, a, t$

$\mathcal{L}_D(f; f') = \frac{1}{|D|} \sum_{(s, a, r, s') \in D} (f(s, a) - r - \gamma \underbrace{f'(s', \pi)}_{f_{k-1}})^2$

$f_k \leftarrow \underset{f \in \mathcal{F}}{\operatorname{argmin}} \mathcal{L}_D(f; f_{k-1})$

$\mathcal{L}_\mu(f; f') = \mathbb{E}_D [\mathcal{L}_D(f; f')]$

Concentration analysis

define.

$$\mathbb{E}_\mu[\|\cdot\|^2] = \|\cdot\|_{2,\mu}^2$$

$$\mathbb{E}_\mu[(f_k - T^\pi f_{k-1})^2] =: \|f_k - T^\pi f_{k-1}\|_{2,\mu}^2$$

Fact:  $\|f_k - T^\pi f_{k-1}\|_{2,\mu}^2 = \underbrace{\mathcal{L}_\mu(f_k; f_{k-1})}_{\Delta} - \underbrace{\mathcal{L}_\mu(T^\pi f_{k-1}; f_{k-1})}_{\Delta}$

Hoeffding:  $|\mathcal{L}_D(f; f') - \mathcal{L}_\mu(f; f')| \leq \underline{\underline{\epsilon}}$   
 $\forall f, f' \in \mathcal{F}$ .  $\epsilon \propto \sqrt{\frac{1}{n} \log \frac{|\mathcal{F}|}{\delta}}$

$$\leq \mathcal{L}_D(f_k; f_{k-1}) - \mathcal{L}_D(T^\pi f_{k-1}; f_{k-1})$$

$f \in \mathbb{R}^{S \times A} + 2\epsilon$

$$T^\pi f_{k-1} \in \mathcal{F}$$

$$\leq 2\epsilon$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\|x\|_2^2$$

$$= x_1^2 + x_2^2 + x_3^2$$

$$\mu = \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right)$$

$$\|x\|_{2, \mu} = \sqrt{\mathbb{E}_{i \sim \mu} [x_i^2]} = \sqrt{\frac{x_1^2 + x_2^2 + x_3^2}{3}}$$

$$\|x\|_1 \geq \|x\|_2 \geq \|x\|_\infty$$

$$\|x\|_{1, \mu} \leq \|x\|_{2, \mu} \leq \|x\|_\infty$$

$\|x\|_{p, \mu}$   $\leftarrow$  distribution.  
 $p=1, 2, \infty$

"nu"

$$\|x\|_{p, \nu} \leq C^{1/p} \|x\|_{p, \mu} \quad \text{where} \quad \frac{\nu(i)}{\mu(i)} \leq C \quad \forall i$$

$$\|x\|_{p, \nu}^p = \mathbb{E}_{i \sim \nu} [ |x_i|^p ]$$

$$= \sum_i \nu(i) |x_i|^p$$

$$= \sum_i \frac{\nu(i)}{\mu(i)} \mu(i) |x_i|^p$$

$$\leq C \|x\|_{p, \mu}^p$$

$$|\hat{J}(\pi) - J(\pi)| \quad d_0 = d_1$$

$$= \left| \mathbb{E}_{s \sim d_0} [f_k(s, \pi)] - \mathbb{E}_{s \sim d_0} [Q^\pi(s, \pi)] \right|$$

$$= \left| \mathbb{E}_{s \sim d_0} \left[ f_k(s, \pi) - \left( \mathcal{T}^\pi f_{k-1} \right)(s, \pi) \right. \right. \\ \left. \left. + \left( \mathcal{T}^\pi f_{k-1} \right)(s, \pi) \right] - \mathbb{E}_{s \sim d_0} \left[ \left( \mathcal{T}^\pi Q^\pi \right)(s, \pi) \right] \right|$$

$$\leq \left| \mathbb{E}_{s \sim d_0} \left[ \left( f_k(s, a) - \left( \mathcal{T}^\pi f_{k-1} \right)(s, a) \right) \right] \right|$$

$a \sim \pi(\cdot | s)$  (I)

$$+ \left| \mathbb{E}_{s \sim d_0} \left[ R(s, \pi) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi)} [f_{k-1}(s', \pi)] \right] \right. \\ \left. - \mathbb{E}_{s \sim d_0} \left[ R(s, \pi) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi)} [Q^\pi(s', \pi)] \right] \right|$$

(II)

$$(I) \leq \| f_k - \mathcal{T}^\pi f_{k-1} \|_1, d_0 \times \pi.$$

$$\leq \| f_k - \mathcal{T}^\pi f_{k-1} \|_2, d_0 \times \pi.$$

$(s, a) \sim d_1^\pi$

$$\leq \| f_k - \mathcal{T}^\pi f_{k-1} \|_{2, \mu} \times \sqrt{C}.$$

$$(II) = \left| \mathbb{E}_{s \sim d_0} \left[ R(s, \pi) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi)} \left[ f_{k+1}(s', \pi) \right] \right] \right.$$

$$\left. - \mathbb{E}_{s \sim d_0} \left[ R(s, \pi) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi)} \left[ Q^\pi(s', \pi) \right] \right] \right|$$

$$\Rightarrow \gamma \left| \mathbb{E}_{s' \sim d_2^\pi} \left[ \underbrace{f_{k+1}(s', \pi)}_{= \hat{c}} - \underbrace{Q^\pi(s', \pi)}_{= Q(\gamma^k V_{\max})} \right] \right|$$

$$\mathbb{E}_{s \sim d_0} [f_K(s, \pi_1)] - J(\pi) = \left( \sum_{t=1}^K \gamma^{t-1} \mathbb{E}_{d_t^\pi} [f_{K-t+1} - \mathcal{T}^{\pi_{t+1}} f_{K-t}] \right) - \mathbb{E}[\sum_{t=K+1}^{\infty} \gamma^{t-1} r_t | \pi, d_0]$$

$$\{f_1, f_2, \dots, f_k\}, \quad \pi_1, \pi_2, \dots, \pi_{k+1} = \pi.$$

$$\mathbb{E}_{s \sim d_0} [f(s, \pi)] - J(\pi) = \frac{1}{1-\gamma} \mathbb{E}_{d^\pi} [f - \mathcal{T}^\pi f].$$

FOI:  $\|f_k - \mathcal{T}f_{k-1}\|_{2,\mu} \leq \varepsilon. \quad \checkmark$

$$J(\pi) - J(\hat{\pi}) \leq \sum_{t=1}^K \gamma^{t-1} \left( \mathbb{E}_{d_t^{\pi}} [\mathcal{T}f_{K-t} - f_{K-t+1}] + \mathbb{E}_{d_t^{\hat{\pi}}} [f_{K-t+1} - \mathcal{T}f_{K-t}] \right) + \gamma^K V_{\max}.$$

$\pi, \hat{\pi} \dots \quad \hat{\pi} = \{\pi_t\}_{t=1:k}, \quad \pi_t = \pi_{f_{K-t+1}}.$

$$\hat{\pi} = \pi_{f_k}.$$

$$J(\pi^*) - J(\pi_{f_k})$$

$$= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi_{f_k}}} [V^*(s) - Q^*(s, \pi_{f_k})].$$

$$\leq \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi_{f_k}}} [V^*(s) - V^{\hat{\pi}}(s)].$$

$$= \frac{1}{1-\gamma} J(\pi^*)_{d^{\pi_{f_k}}} - J(\hat{\pi})_{d^{\pi_{f_k}}}.$$