

OPE | IS: exponential variance. ← can we avoid?  
 "curse of horizon" FOI:  $\max_{a'} f_k(s', a')$

• FQE (policy-eval ver. of FQZ)

$$f_{k+1} \leftarrow \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{|\mathcal{D}|} \sum_{(s, a, r, s')} (f(s, a) - r - \gamma f_k(s', \pi))^2$$

Require:  $\left\{ \begin{array}{l} \forall f \in \mathcal{F}, J^\pi f \in \mathcal{F} \\ \left\| \frac{d^\pi}{\mu} \right\|_\infty \leq C. \end{array} \right.$

data:  
 $(s, a) \sim \mu,$   
 $r \sim R(s, a)$   
 $s' \sim P(\cdot | s, a).$

⇒ guarantee poly sample comp. for FQE.

• Can we make only realizability assumptions?

Maximalised IS

$$J(\pi) = \frac{1}{1-\gamma} \mathbb{E}_{\substack{(s, a) \sim d^\pi \\ r \sim R(s, a)}} [r]$$

Then: if we know:

$$w^\pi(s, a) := \frac{d^\pi(s, a)}{\mu(s, a)} \quad (\text{assume } \leq C)$$

$$J(\pi) = \frac{1}{1-\gamma} \mathbb{E}_\mu [w^\pi(s, a) \cdot r]$$

in IS:  $\frac{P^\pi(s_1, a_1, r_1, \dots, s_n, a_n, r_n)}{P^{\pi_b}(s_1, a_1, r_1, \dots, s_n, a_n, r_n)}$   $\circlearrowright r_n$

in MZS:  $\frac{P^\pi(s_n, a_n, r_n)}{P^{\pi_b}(s_n, a_n, r_n)} \cdot r_n$

$w^\pi(s,a)$  is unknown — must learn using func. approx.

MDL Try to learn  $q \approx Q^\pi$  w/ the help of  $w^\pi$

$$\forall q. J(\pi) - \mathbb{E}_{s \sim d_0} [q(s, \pi)] \\ = \frac{1}{1-\gamma} \mathbb{E}_{\substack{(s,a) \sim d^\pi \\ r \sim R(s,a) \\ s' \sim P(\cdot | s,a)}} [r + \gamma q(s', \pi) - q(s, a)]$$

$$|J(\pi) - \mathbb{E}_{s \sim d_0} [q(s, \pi)]| \\ = \frac{1}{1-\gamma} \left| \mathbb{E}_{(s,a) \sim d^\pi} [r + \gamma q(s', \pi) - q(s, a)] \right|$$

$$= \frac{1}{1-\gamma} \left| \mathbb{E}_\mu [w^\pi(s, a) \cdot (r + \gamma q(s', \pi) - q(s, a))] \right|$$

$w^\pi \in \mathcal{W}$  ← realizability.

$$\leq \max_{w \in \mathcal{W}} \frac{1}{1-\gamma} \left| \mathbb{E}_\mu [w(s, a) \cdot (r + \gamma q(s', \pi) - q(s, a))] \right|$$

$L_q(w, q)$

$$\hat{q} = \arg \min_{q \in \mathcal{Q}} \max_{w \in \mathcal{W}} L_q(w, q) \Rightarrow \hat{J}(\pi) = \mathbb{E}_{d_0} [\hat{q}]$$

consistent estimator of  $J(\pi)$

if  $\omega^\pi \in W$ ,  $Q^\pi \in Q$ .

• When  $\omega^\pi \in W$ ,  $|J(\pi) - \mathbb{E}_{d_0}[g]| \leq \max_{\omega \in W} L_g(\omega, g)$

• For  $g = Q^\pi$ ,  $\max_{\omega \in W} L_g(\omega, g) = 0$ .

(if  $Q^\pi \in Q$ ,  $\min_{g \in Q} \max_{\omega \in W} L_g(\omega, g) = 0$ )

$\Rightarrow \underset{g}{\operatorname{argmin}}$  is accurate in its OPE prediction.

Remarks:  $\omega^\pi \in W$  can be relaxed to  $\omega^\pi \in \operatorname{conv}(W)$ .

**MWL** Goal: learn  $w \approx \omega^\pi$

i.e. learn  $w$  to minimize

$$\left| \boxed{J(\pi)} - \frac{1}{1-\gamma} \mathbb{E}_\mu[w \cdot r] \right| \quad (*) \quad \mathbb{E}[\cdot | s, a] = R(s, a)$$

$$= \left| \boxed{\mathbb{E}_{s \sim d_0}[Q^\pi(s, \pi)]} + \frac{1}{1-\gamma} \mathbb{E}_\mu \left[ w(s, a) (\gamma Q^\pi(s', \pi) - Q^\pi(s, a)) \right] \right|$$

$$\boxed{Q^\pi \in Q} \leq \max_{g \in Q} \left| \mathbb{E}_{s \sim d_0}[g(s, \pi)] + \frac{1}{1-\gamma} \mathbb{E}_\mu[w(s, a)(\gamma g(s', \pi) - g(s, a))] \right|$$

$$= \max_{q \in \mathcal{Q}} L_w(w, q).$$

Bellman flow error of  $\mu, w$ .

Verify this is tight upper bound:

$$\forall q, L_w(w^\pi, q)$$

$$= \left| \mathbb{E}_{s \sim d_0} [q(s, \pi)] + \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim d^\pi} [\gamma q(s', \pi) - q(s, a)] \right|$$

$$= 0. \quad (\text{Bellman flow eq for } d^\pi).$$

$$-q(s, a) \quad w/ \quad (s, a) \sim \sum_{t=1}^{\infty} \gamma^{t-1} d_t^\pi$$

$$+q(s, a) \quad w/ \quad (s, a) \sim d_1^\pi$$

$$+q(s, a) \quad w/ \quad (s, a) \sim \gamma P(d^\pi) \times \pi$$

$$= \sum_{t=2}^{\infty} \gamma^{t-1} d_t^\pi$$

MIS for learning good policy?  $\left[ \mathbb{E}_\mu [w \cdot (q - T_q)] \right]$

Learn  $\underline{q} \approx Q^*$  s.t.  $\pi_q$  is a good policy.  $\uparrow$

$$(*) \underset{q \in \mathcal{Q}}{\text{argmin}} \max_{w \in W} \frac{1}{1-\gamma} \left| \mathbb{E}_\mu \left[ w(s, a) \cdot \left( q(s, a) - \gamma - \gamma \max_{a'} q(s', a') \right) \right] \right|$$

When does this work?  $\underline{Q^* \in \mathcal{Q}}$ .

$$J(\pi^*) - J(\pi_q) \leq \frac{1}{1-\gamma} \left( \mathbb{E}_{d^{\pi^*}} [T_q - q] + \mathbb{E}_{d^{\pi_q}} [q - T_q] \right).$$

↓  
can't guarantee  
 $q = T_q$   
on all states.

Implies  $\Rightarrow$  (\*) learns good policy if  $w^{\pi_q} \in \mathcal{W}, \forall q \in \mathcal{Q}$ .

Proof: Recall  $\forall \pi, q, J(\pi) - \mathbb{E}_{d_0} [q(s, \pi)] = \frac{1}{1-\gamma} \mathbb{E}_{d^\pi} [Tq - q]$ .

$$\begin{aligned} J(\pi^*) - J(\pi_q) &\leq J(\pi^*) - \mathbb{E}_{s \sim d_0} [q(s, \pi^*)] + \mathbb{E}_{d_0} [q(s, \pi_q)] \\ &= \frac{1}{1-\gamma} \mathbb{E}_{d^{\pi^*}} [T^{\pi^*} q - q] + \mathbb{E}_{d_0} [q(s, \pi_q)] - J(\pi_q) \\ &+ \frac{1}{1-\gamma} \mathbb{E}_{d^{\pi_q}} [q - Tq]. \end{aligned}$$

$q \approx Q^* \begin{cases} q \approx Q^{\pi^*} \\ q \approx Q^{\pi_q} \end{cases}$

$$\leq \frac{1}{1-\gamma} \left( \mathbb{E}_{d^{\pi^*}} [Tq - q] + \mathbb{E}_{d^{\pi_q}} [q - Tq] \right).$$

Note: if  $Q^* \in \mathcal{Q}$  and  $\forall q \in \mathcal{Q}, w^{\pi_q} \in \mathcal{W}$ .

$$\begin{aligned} &|\mathbb{E}_{d^{\pi^*}} [Tq - q] + \mathbb{E}_{d^{\pi_q}} [q - Tq]| \\ &\leq 2 \cdot \max_{w \in \mathcal{W}} |\mathbb{E}_\mu [w \cdot (q - Tq)]|. \end{aligned}$$