

DPE IS: exp variance. (unless $\pi \approx \pi_v$).

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

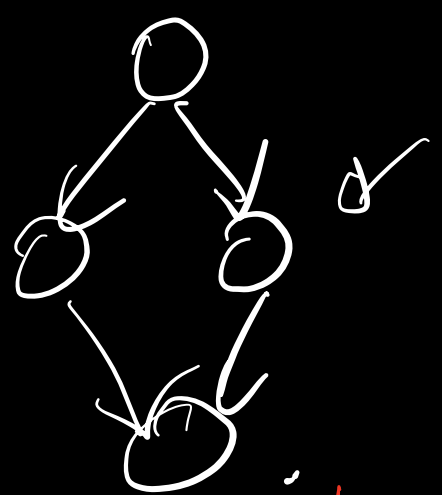
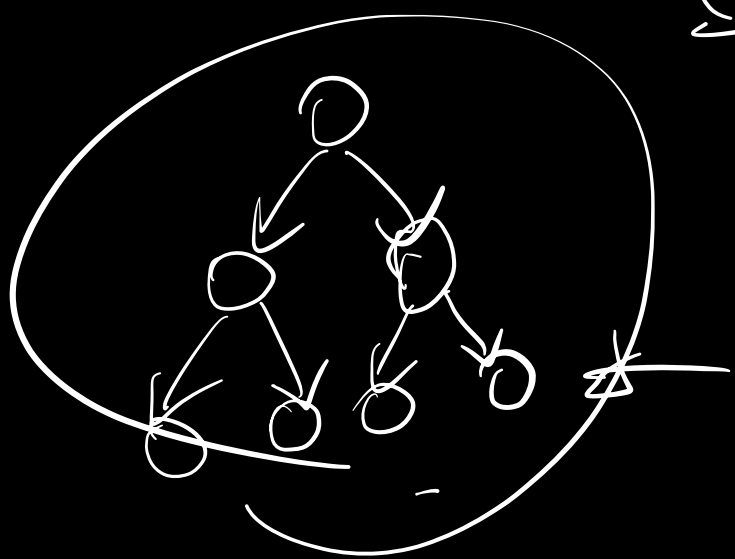
Require $\{ T^\pi f \in \mathcal{F} \ \forall f \in \mathcal{F}$

$\| \frac{d^\pi}{\mu} \|_{\infty} \leq C$

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

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

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

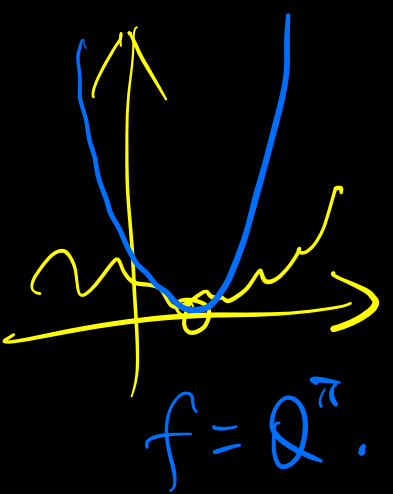


$\mathbb{Q}^\pi \in \mathcal{F}$

Setup: data $\sim \mu$ $w^\pi \in W$, $w^\pi(s,a) = \frac{d^\pi(s,a)}{\mu(s,a)}$
 Goal: estimate $J(\pi) = \mathbb{E}_\mu [w^\pi \cdot r] / (1-\gamma)$.

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

$$L(f) = | \cdot | = \frac{1}{1-\gamma} | \mathbb{E}_{d^\pi} [(f - \mathcal{T}^\pi f)] |$$



$$= \frac{1}{1-\gamma} | \mathbb{E}_\mu \left[\frac{d^\pi(s,a)}{\mu(s,a)} (f - \mathcal{T}^\pi f) \right] |$$

$$\leq \frac{1}{1-\gamma} \max_{w \in W} | \mathbb{E}_\mu [w \cdot (f - \mathcal{T}^\pi f)] |$$

$$= \frac{1}{1-\gamma} \max_{w \in W} | \mathbb{E}_{(s,a,r,s')} [w(s,a) (f(s,a) - r - \gamma f(s,a))] |$$

$\sim \mu$

Alg: $\hat{f} = \arg \min_{f \in \mathcal{F}} \max_{w \in W}$

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

$$f: \forall w \in W, \mathbb{E}_\mu [w \cdot (f - \mathcal{T}^\pi f)] = 0.$$

$$\forall (s, a) (f - \mathcal{T}^\pi f)(s, a) = 0.$$

$$\left| \frac{1}{1-\gamma} \mathbb{E}_\mu \left[\frac{d}{\mu} R \right] - J(\pi) \right|$$

$$w = \frac{d^\pi}{\mu}$$

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

$$\hat{w} = \underset{w \in W}{\operatorname{argmin}} \max_{f \in \mathcal{F}} | \int f |$$

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

$$\mathcal{O} \left(\|\hat{w}\|_\infty \left(F \sqrt{\frac{1}{n} \ln \frac{k}{\delta}} \right) \right)$$