# Marginalized Importance Sampling (MIS)

$$\rho_{1:H} \sim (\pi/\pi_b)^H$$

- IS: exponential variance unless $\pi \approx \pi_b$.

- FQE. (policy-eval ver. of FQ2).

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

- $\forall f \in \mathcal{F}, \ \mathcal{T}^\pi f \in \mathcal{F}.$
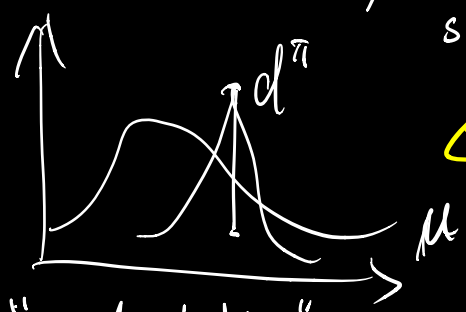
  $\forall f \in \mathcal{F}, \ \mathcal{T}^\pi f \in \mathcal{F}.$ for FQ2.

  $\Rightarrow$ w/ $\mathcal{F}$ bounded complexity. poly sample size.

  $$\| f_{k+1} - \mathcal{T}^\pi f_k \|_{2, \mu} \leq \varepsilon$$

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

  $$\Rightarrow \left\| \frac{d^\pi}{\mu} \right\|_\infty \leq C.$$

  $$\Rightarrow \left| J(\pi) - \underset{s \sim d_0}{\mathbb{E}}[f_k(s, \pi)] \right|$$

  $$\leq \text{poly}(C) \cdot \varepsilon. \quad \text{"realizability."}$$

- Problem w/ FQE: $\boxed{Q^\pi \in \mathcal{F}}$ is insufficient

- Is there method w/ $\underset{\wedge}{\text{func-approx}}$ error?
  monotone.

- Is there anything we can do. if we only have $\underline{Q^\pi \in \mathcal{F}}$

$\Rightarrow$ MIS address both questions.

  Dai, Nachum... "DICE", "M_L".

  "density estimation" ....

$MDP = (S, A, P, R, \gamma, S_0)$    $d_0$.

w.t. learn    $J(\pi) = Q^{\pi}(S_0, \pi) = \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim d^{\pi}} [r]$.

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

       $r \sim R(s,a)$

---

"Eval error lemma for $\mathcal{E}$":      TD/ Bellman error.

$\forall \mathcal{E}. \quad J(\pi) - \mathcal{E}(S_0, \pi) = \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim d^{\pi}} [r + \gamma \mathcal{E}(s', \pi) - \mathcal{E}(s, a)]$.

Proof: $d^{\pi} = (1-\gamma) \sum_{t=1}^{\infty} \gamma^{t-1} d_t^{\pi} \longrightarrow$ dist. of $(s_t, a_t)$ under policy $\pi$.

$RHS = \mathbb{E}_{(s,a) \sim d_1^{\pi}} [r + \gamma \mathcal{E}(s', \pi) - \mathcal{E}(s, a)]$
     $(s,a) \sim d_1^{\pi}$
     $s' \sim P(\cdot | s, a)$
     $a' \sim \pi(\cdot | s')$

$+ \gamma \mathbb{E}_{(s,a) \sim d_2^{\pi}} [r + \gamma \mathcal{E}(s', \pi) - \mathcal{E}(s, a)]$.

$+ \gamma^2 \quad \ldots \quad \vdots$

$\mathbb{E}[\sum_{t=1}^{\infty} \gamma^{t-1} r_t | \pi]. \quad \vdots$

$= \sum_{t=1}^{\infty} \gamma^{t-1} \mathbb{E}_{(s,a) \sim d_t^{\pi}} [r] - \mathbb{E}_{(s,a) \sim d_1^{\pi}} [\mathcal{E}(s, a)]$.

$= J(\pi) - \mathcal{E}(S_0, \pi) = LHS.$      □

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

the remaing: $0 = \mathcal{E}(S_0, \pi) + \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim d^{\pi}} [\gamma \mathcal{E}(s', \pi) - \mathcal{E}(s, a)]$

"Bellman $\mathcal{E}$ for occupancy"

$d^{\pi} = \begin{array}{|c|} \hline d_1^{\pi} \\ \hline \gamma d_2^{\pi} \\ \hline \gamma^2 d_3^{\pi} \\ \hline \gamma^3 d_4^{\pi} \\ \hline \end{array} \longrightarrow \begin{array}{|c|} \hline d_1^{\pi} \\ \hline \gamma d_1^{\pi} \\ \hline \gamma^2 d_1^{\pi} \\ \hline \vdots \\ \hline \end{array}$

$\min_{\mathcal{g}}$ $\left| J(\pi) - \mathcal{g}(s_0, \pi) \right| = \left| \frac{1}{1-\gamma} \mathbb{E}_{d^\pi} \left[ r + \gamma \mathcal{g}(s', \pi) - \mathcal{g}(s,a) \right] \right| \leftarrow$

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

(arrow to $d^\pi$)

Assume: $\underline{w^{\pi/\mu} \in W}$ $\Delta$

$w^{\pi/\mu}(s,a)$
"marginalized imp. weight"
"density ratio"

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

$\Delta$

---

side. comment.

$\mathbb{E}_{\mathcal{g}} \left[ \frac{p(x)}{\mathcal{g}(x)} f(x) \right]$. $\rightarrow$ suppose.
$y = \mathcal{g}(x)$
$f(x) = f'(y)$.

$\mathbb{E}_{\mathcal{g}} \left[ \frac{p_\gamma(y)}{\mathcal{g}_\gamma(y)} f'(y) \right]$

in general:

$Var \left[ \frac{p_\gamma(y)}{\mathcal{g}_\gamma(y)} \right] \leq Var \left[ \frac{p(x)}{\mathcal{g}(x)} \right]$.

$\rho_{1:H} = \prod \frac{\pi(a_t | s_t)}{\pi(a_t | s_t)} = \frac{p^\pi(s_1, a_1, s_2, a_2 \cdots, s_H, a_H)}{p^{\pi_b}(s_1, a_1, s_2, a_2 \cdots s_H, a_H)}$

$\bigvee w^{\pi/\mu} = \frac{p^\pi(s_t, a_t)}{p^{\pi_b}(s_t, a_t)} \longleftarrow$ or arbitrary $\mu$.

in our ctx.

$X = s_1, a_1, s_2, a_2, \cdots, s_H, a_H$.

$g(x) = s_t, a_t$.

---

$$\left| J(\pi) - \mathcal{g}(s_0, \pi) \right| \leq \sup_{w \in W} \frac{1}{1-\gamma} \left| \mathbb{E}_\mu \left[ w(s,a)(r + \gamma \mathcal{g}(s', \pi) - \mathcal{g}(s,a)) \right] \right|$$

assuming $\underline{w^{\pi/\mu} \in W}$ (can relax to

$w^{\pi/\mu} \in conv(W)$.
$\Delta$

Alg: over $\mathcal{Q}$ class

$$\arg\min_{\mathcal{g} \in \mathcal{Q}} \sup_{w \in W} \left| \mathbb{E}_\mu \left[ w \cdot (r + \gamma \mathcal{g} - \mathcal{g}) \right] \right|$$

"MQL [Uehara et al '20]."

Why relaxation from $W^\pi/\mu$ to $\sup_{w \in W}$ make sense?

→ Ideally $q = Q^\pi$. → "tight relaxation"

for this func. loss: $\left[ \sup_{w \in W} \frac{1}{1-\gamma} \mathbb{E}_\mu \left[ w \cdot \left( r + \gamma q(s', \pi) - q(s,a) \right) \right] \right]$

$= 0$.   $\mathbb{E}[\cdot | s,a] = 0$.

→ MQL: 
| when $W^\pi/\mu \in W$ | valid upper bound of $|J(\pi) - q(s_0, \pi)|$ |
| when $Q^\pi \in Q$: upper bound can be minimized to $0$. |

→ "$q$ generator"   → "$w$" discriminator.

---

"MWL"   "$w$ generator"   "$q$ discriminator"

---

find $q$ : s.t.

→ $\mathbb{E}_\mu \left[ \overset{\downarrow}{w(s,a)} \cdot \left( q(s,a) - r - \gamma q(s', \pi) \right) \right] = 0$.   Bellman error.

different: $\| q - T^\pi q \|^2_{2,\mu}$.

$= \mathbb{E}_\mu \left[ \left( q(s,a) - \mathbb{E}_{r, s' | s, a} \left[ r + \gamma q(s', \pi) \right] \right)^2 \right]$.

$\triangle$

LSTD.
avg Bellman
eq.
compare to
sq. ptwise.
alt. view of
MQL.

Want to learn $w$. s.t.

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

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

"Eval error lemma or w"

$$= Q^\pi(s_0, \pi) + \frac{1}{1-\gamma} \mathbb{E}[w(s,a) \cdot (\gamma Q^\pi(s', \pi) - Q^\pi(s,a))]. \quad \triangle$$

Proof: $\frac{1}{1-\gamma} \cdot \mathbb{E}_\mu [w(s,a) \cdot \underbrace{\left( r + \gamma Q^\pi(s', \pi) - Q^\pi(s,a) \right)}_{\mathbb{E}[\cdot | s,a] = 0.}] \quad = 0.$

---

find $w$. s.t. $\boxed{|J(\pi) - \frac{1}{1-\gamma} \mathbb{E}_\mu[w(s,a) \cdot r]|} < 0.$

$\boxed{\text{Assume } Q^\pi \in Q.}$ $\leq \sup_{q \in Q} \left| \underbrace{\left( q(s_0, \pi) + \frac{1}{1-\gamma} \mathbb{E}_\mu [w \cdot (\gamma q - q)] \right)}_{\triangle} \right|.$

"MWL": $\underset{w \in W}{\arg\min} \sup_{q \in Q}$ $J(\pi) \in \frac{1}{1-\gamma} \mathbb{E}_\mu [\hat{w} \cdot r] \pm \sup_q |\underline{L_W(\hat{w}, q)}|$

Relaxation is "tight": $w = w^{\pi/\mu}$

$\forall q$. $q(s_0, \pi) + \frac{1}{1-\gamma} \mathbb{E}_{d^\pi} [(\gamma q(s', \pi) - q(s,a))]$.

$= < q, \underbrace{(d_0 \times \pi)}_{} - \frac{1}{1-\gamma} d^\pi + \frac{\gamma}{1-\gamma} \underbrace{((P^\top d^\pi) \times \pi)}_{}>$

$= 0.$ $\qquad = \hat{0}.$ b/c Bellman eq for $d^\pi$

so.

$\text{if} \begin{cases} Q^\pi \in Q \Rightarrow \text{valid upper bound.} \\ w^{\pi/\mu} \in W \Rightarrow \text{minimize upper bound.} \end{cases}$

$(\to 0).$

$$J(\pi) = \frac{1}{1-\gamma} \mathbb{E}_\mu[w \cdot r] + Q^\pi(s_0, \pi) + \frac{1}{1-\gamma} \mathbb{E}_\mu[w(\gamma Q^\pi(s',\pi) - Q^\pi(s,a)].$$

$$J(\pi) = Q(s_0, \pi) + \frac{1}{1-\gamma} \mathbb{E}_{d^\pi}[r + \gamma Q(s',\pi) - Q(s,a)].$$

Define $L(w, Q) = \frac{1}{1-\gamma} \mathbb{E}_\mu[w \cdot r] + Q(s_0, \pi) + \frac{1}{1-\gamma} \mathbb{E}_\mu[w(\gamma Q(s',\pi) - Q(s,a)].$

$$J(\pi) = L(w, Q^\pi) = L(w^{\pi/\mu}, Q) \quad \forall w, Q.$$

$Q^\pi \in Q$ then:

$\forall w. \quad J(\pi) = L(w, Q^\pi) \leq \sup_{Q \in Q} L(w, Q)$

$\geq \inf_{Q \in Q} L(w, Q).$

$J(\pi) \in [\inf_Q L(w, Q), \sup_Q L(w, Q)] \quad \forall w.$

$\Rightarrow \sup_{w \in W} \inf_{Q \in Q} L(w, Q) \leq J(\pi) \leq \inf_{w \in W} \sup_{Q \in Q} L(w, Q) \checkmark.$

$w^{\pi/\mu} \in W$

$\forall Q. \quad \inf_{w \in W} L(w, Q) \leq J(\pi) = L(w^{\pi/\mu}, Q) \leq \sup_{w \in W} L(w, Q).$

$\Rightarrow \sup_{Q \in Q} \inf_{w \in W} L(w, Q) \leq J(\pi) \leq \inf_{Q \in Q} \sup_{w \in W} L(w, Q). \checkmark$

Sion's minimax thm: b/c $L(w, Q)$ is convex-concave in $w/Q$. and $W$ & $Q$ are convex.

$\Rightarrow \sup_Q \inf_w L(w, Q) = \inf_w \sup_Q L(w, Q).$

Why misspecification? $\quad W^{\pi/\mu} \notin W \quad\quad Q^\pi \notin Q$

$$\boxed{\frac{d^\pi}{\mu}} \quad\quad \mu = d^{\pi_b}$$

---

How about learning? $\quad$ For OPE, we $\wedge$ learn $Q^\pi$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ try to.

For learning, we try to learn $Q^*$. $\quad \boxed{\max_\pi Q^*(s',a')}$

$\Rightarrow \forall w. \quad \mathbb{E}_\mu \left[ w(s,a) \cdot \left( Q^*(s,a) - r - \gamma V_{Q^*}(s') \right) \right] = 0.$

"MABO" $\Downarrow$
$\underset{g \in Q}{\arg\min} \; \boxed{\underset{w \in W}{\max}} \; \left| \mathbb{E}_\mu \left[ w \cdot \left( g(s,a) - r - \gamma V_g(s') \right) \right] \right|.$

$\quad\quad \hookrightarrow$ output $\pi_g$. $\leftarrow$ when near-optimal?

Func-approx: ① $\quad Q^* \in Q$.

$\quad\quad\quad\quad$ ② $\quad \forall \pi_g \;\; s.t. \;\; g \in Q, \quad \dfrac{d^{\pi_g}}{\mu} \in W.$

$\quad\quad\quad\quad\quad\quad\quad \forall g. \quad\quad\quad\quad \underset{\uparrow}{r + \gamma g(s',\pi) - g(s,a)}$

Lemma:

$$J(\pi^*) - J(\pi_g) \le \frac{1}{1-\gamma}\left( \mathbb{E}_{d^{\pi^*}}[Tg - g] + \mathbb{E}_{d^{\pi_g}}[g - Tg] \right)$$

Proof: Recall eval error Lemma:

$$\forall \pi, \quad J(\pi) - g(s_0, \pi) = \frac{1}{1-\gamma} \mathbb{E}_{d^\pi}\left[ r + \gamma g(s',\pi) - g(s,a) \right],$$

$$J(\pi^*) - J(\pi_g) = J(\pi^*) - g(s_0, \pi_g) + g(s_0, \pi_g) - J(\pi_g).$$

$$\le J(\pi^*) - g(s_0, \pi^*) \quad + \quad g(s_0, \pi_g) - J(\pi_g).$$

$\quad\quad\quad\quad \downarrow$ eval $\pi^*$ use $g$. $\quad\quad\quad\quad\quad\quad\quad\quad\quad \downarrow$ eval $\pi_g$ use $g$.

$$= \frac{1}{1-\gamma} \mathbb{E}_{d^{\pi^*}}\left[ r + \gamma g(s', \pi^*) - g(s,a) \right] - \frac{1}{1-\gamma} \mathbb{E}_{d^{\pi_g}}\left[ r + \gamma g(s',\pi_g) - g(s,a) \right]$$

$$\leq \frac{1}{1-\gamma} \mathbb{E}_{d^{\pi^*}}\left[r + \gamma g(s', \pi_g) - g(s,a)\right] - \frac{1}{1-\gamma} \mathbb{E}_{d^{\pi_g}}\left[\mathcal{T}_g - g\right].$$

$$= \frac{1}{1-\gamma} \mathbb{E}_{d^{\pi^*}}\left[\mathcal{T}_g - g\right] - \frac{1}{1-\gamma} \mathbb{E}_{d^{\pi_g}}\left[\mathcal{T}_g - g\right]. \quad \checkmark$$

$$\underset{w}{\arg\min} \underset{g}{\max} \left[\cdot - \frac{\ell}{\cdot}\right]$$