

Importance Sampling (IS) | Want to estimate  $\mathbb{E}_{x \sim p} [f(x)]$ .

▷  $x \sim p$ ,  $f(x)$  is Monte-Carlo estimate. ↵

▷ What if we only have  $x \sim q$ ?

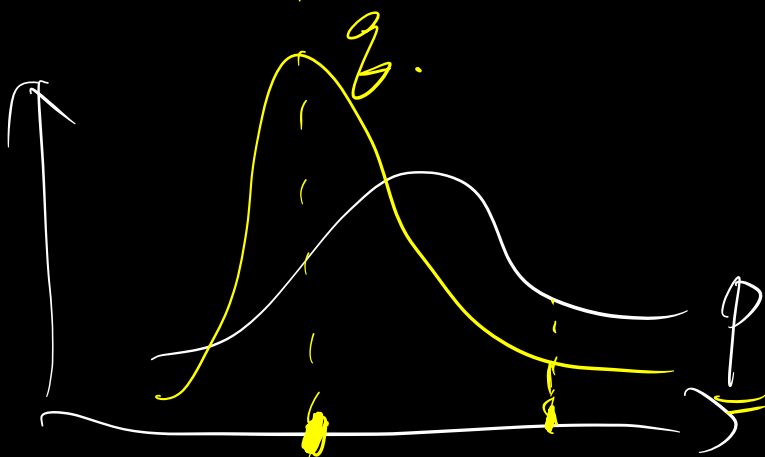
IS:  $x \sim q$ ,  $\frac{p(x)}{q(x)} \cdot f(x)$  (assume  $\forall x, \text{ s.t. } p(x) > 0, q(x) > 0$ )

$$\mathbb{E}_{x \sim q} \left[ \frac{p(x)}{q(x)} f(x) \right] = \sum_{x \in X} q(x) \cdot \frac{p(x)}{q(x)} f(x)$$

$$= \sum_x p(x) f(x) = \mathbb{E}_{x \sim p} [f(x)].$$

$\frac{p(x)}{q(x)}$  : "density ratio" / "importance weight".

$$\mathbb{E}_{x \sim q} \left[ \frac{p(x)}{q(x)} \right] = \sum_{x \in X} q(x) \cdot \frac{p(x)}{q(x)} = 1.$$



IW "weighting"

IPS:

inverse propensity scores.

Application to bandits.

behavior policy.

CB:  $x \sim d_0$ ,  $a \sim \pi_b(\cdot | x)$ ,  $r \sim R(\cdot | x, a)$ .

Off-policy evaluation: how to evaluate  $\pi$ .

(i.e. estimate  $\mathbb{E}[r | \pi]$  using data from  $\pi_b \neq \pi$ .

target/eval policy.

$(x, a, r) \sim p \Leftrightarrow x \sim d_0, a \sim \pi, r \sim R(\cdot | x, a)$ .

$(x, a, r) \sim q \Leftrightarrow x \sim d_0, a \sim \pi_b, r \sim R(\cdot | x, a)$ .

$$\mathbb{E}[r | \pi] = \mathbb{E}_{(x, a, r) \sim p} [r].$$

$$= \mathbb{E}_{(x, a, r) \sim q} \left[ \frac{p(x, a, r)}{q(x, a, r)} \cdot r \right].$$

$$\frac{p(x, a, r)}{q(x, a, r)} = \frac{p(x) \cdot p(a|x) \cdot p(r|x, a)}{q(x) \cdot q(a|x) \cdot q(r|x, a)}$$

$$\begin{aligned}
 &= \frac{\cancel{d_0(x)} \cdot \pi(a|x) \cdot \cancel{R(r|x,a)}}{\cancel{d_0(x)} \cdot \pi_b(a|x) \cdot \cancel{R(r|x,a)}} \\
 &= \frac{\pi(a|x)}{\pi_b(a|x)}
 \end{aligned}$$

$\Rightarrow$  IS estimator for DPG:  $\frac{\pi(a|x)}{\pi_b(a|x)} \cdot \gamma$ .

To apply IS in practice:

record data:  $(x, a, r, \pi_b(a|x))$ .

Special Case

logging/proposal prob.

$\pi_b(a|x) = 1/|A|$ . (unif policy)

$\pi(x)$  is deterministic.

IS estimator:  $\frac{\mathbb{I}[\pi(x)=a]}{1/|A|} \cdot \gamma$

• Assume  $r \in [0, 1]$ ,  $\pi \in [0, |A|]$ .

• In expectation, only  $n/|A|$  data pts "useful"

$$\cdot \left\{ (x^{(i)}, a^{(i)}, y^{(i)}) \right\}_{i=1}^n$$

$$\frac{1}{n} \sum_{i=1}^n \frac{\mathbb{I}[\pi(x^{(i)}) = a^{(i)}]}{|A|} \cdot y^{(i)}$$

$$= \frac{1}{n/|A|} \sum_{i: \pi(x^{(i)}) = a^{(i)}} y^{(i)}$$

expected size of  $\{i: \pi(x^{(i)}) = a^{(i)}\}$   
(instead of actual size).

• Weighted LS / self-normalized.

replace  $n/|A|$  by  $|\{i: \pi(x^{(i)}) = a^{(i)}\}|$ .

In the general case:

replace  $n$  by  $\sum_{i=1}^n \frac{\pi(a|x)}{\pi_b(a|x)}$ .

Variance: suppose  $r \in [0, 1]$ .

$$\frac{\mathbb{I}[\pi(x)=a]}{1/|A|} \cdot r \in [0, |A|]$$

for r.v. supported on  $[0, |A|]$ .

Var  $\sim O(|A|^2)$  in worst case

$$\begin{aligned} & \text{Var} \left[ \frac{\mathbb{I}[\pi(x)=a]}{1/|A|} \cdot r \right] \\ & \leq \mathbb{E} \left[ \frac{\mathbb{I}[\pi(x)=a]}{1/|A|^2} \cdot r^2 \right] \\ & \leq |A| \cdot \mathbb{E} \left[ \frac{\mathbb{I}[\pi(x)=a]}{1/|A|} \right] \cdot 1 = |A|. \end{aligned}$$

Further special case  $r \equiv \text{const}$ .

we know:  $\forall \pi, \mathbb{J}(\pi) = \text{const}$

this should be "easy".

MC estimator:  $O$  variance.

IS estimator:  $\Theta(|A|)$  variance.

How to improve IS?

(i) WLS -

$$(ii) \text{Var}(IS) \leq \mathbb{E}\left[\frac{\mathbb{I}(X(x)=a)}{1/|A|^2} \cdot r^2\right].$$

to reduce variance: center reward around 0.

in the special case:  $\left\{ \begin{array}{l} \text{shift reward by const.} \\ \text{DPE} \\ \text{add const back.} \end{array} \right.$

General case: Doubly Robust (DR)

Given  $\hat{Q}: X \times A \rightarrow \mathbb{R}$ . such that

$$\hat{Q}(x, a) \approx \mathbb{E}[r | x, a].$$

$$\text{DR: } \hat{Q}(x, \pi) + \frac{\pi(a|x)}{\pi_b(a|x)} \cdot (r - \hat{Q}(x, a))$$

Claim: DR is unbiased.

$$\mathbb{E}_{\pi_b} \left[ \hat{Q}(x, \pi) - \frac{\pi(a|x)}{\pi_b(a|x)} \hat{Q}(x, a) \right] = 0.$$

$\uparrow \sim \pi_b$

IS is special case where  $\hat{Q} \equiv 0$ .

# Application to multi-step RL

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

$$\tau = (s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_H, a_H, r_H).$$

Collect data w/  $a_t \sim \pi_b$ .

to eval policy  $\pi$ ,  $\mathbb{E} \left[ \sum_{n=1}^H \gamma^{n-1} r_n \mid \pi \right]$ .

## Derivation of IS

Define  $\tau \sim p$  as induced by  $\pi$ .

$\tau \sim q$  --- by  $\pi_b$ .

$$\mathbb{E}_{\tau \sim p} \left[ \sum_{n=1}^H \gamma^{n-1} r_n \right] = \mathbb{E}_{\tau \sim q} \left[ \underbrace{\frac{p(\tau)}{q(\tau)}}_{\text{ratio}} \cdot \sum_{n=1}^H \gamma^{n-1} r_n \right]$$

$$\frac{p(\tau)}{q(\tau)} = \frac{\cancel{d_0(s_1)} \cdot \pi(a_1|s_1) \cdot \cancel{P(s_2|s_1, a_1)} \cdot \pi(a_2|s_2) \cdots}{\cancel{d_0(s_1)} \cdot \pi_b(a_1|s_1) \cdot \cancel{P(s_2|s_1, a_1)} \cdot \pi_b(a_2|s_2) \cdots}$$

$$= \frac{\pi(a_1|s_1)}{\pi_b(a_1|s_1)} \cdots \frac{\pi(a_H|s_H)}{\pi_b(a_H|s_H)}$$

$$=: \rho_1 \cdot \rho_2 \cdots \rho_H = \rho_{1:H}$$

$$\Rightarrow \text{IS estimator: } \rho_{1:H} = \sum_{h=1}^H \gamma^{h-1} r_h$$


---

Special case:  $\pi$  is deterministic  
 $\pi_b$  is unif. rand.

$$\rho_{1:H} = \begin{cases} 0, & \exists h, \pi(s_h) \neq a_h \\ |A|^H, & \text{o.w.} \end{cases} \quad \left( \text{w.p. } 1/|A|^H \right)$$


---

"Curse of horizon"

In general, blow up is  $\sim$

$$\left( \max_{s,a} \frac{\pi(a|s)}{\pi_b(a|s)} \right)^H$$



therefore, IS only works when either  
 $H$  is small, or,  $\pi \approx \pi_b$ .

---

## Per-step IS

$$\mathbb{E} \left[ \sum_{h=1}^H \gamma^{h-1} r_h \mid \pi \right] = \sum_{h=1}^H \gamma^{h-1} \mathbb{E}[r_h \mid \pi]$$

Idea: estimate  $\mathbb{E}[r_h \mid \pi]$  for  $h \in [H]$

when estimating  $\mathbb{E}[r_h \mid \pi]$ .

pretend traj is  $(s_1, a_1, r_1, \dots, s_h, a_h, r_h)$ .

$\Rightarrow$  (step-wise) IS:  $\sum_{h=1}^H \gamma^{h-1} \cdot \underbrace{p_{1:h}}_{\Delta} r_h$ .

DR

rewrite step-wise IS.

recursively:

Define  $\underline{V_0} := 0$ , then  $\forall h < H$ .

$$V_{H-h+1} = P_h (r_n + \gamma V_{H-h})$$

"remaining  
time steps"

Claim:  $V_H = \text{step-wise IS}$ .

$$\begin{aligned} V_H &= P_1 (\underline{r_1} + \gamma P_2 (\underline{r_2} + \gamma P_3 (\dots))) \\ &= P_1 r_1 + \gamma P_{1:2} r_2 + \gamma^2 P_{1:3} r_3 \\ &\quad + \dots + \gamma^{H-1} P_{1:H} r_H. \end{aligned}$$

$$V_{H-h+1} = P_h (r_n + \gamma V_{H-h})$$

$\underline{V_{H-h+1}^\pi(s_n)}$ 
 $\frac{\pi(a_n|s_n)}{\pi_b(a_n|s_n)}$ 
 $Q_{H-h+1}^\pi(s_n, a_n)$ 
 $\frac{V_{H-h}^\pi(s_{h+1})}{\pi_b}$

Give  $\hat{Q} \approx Q^\pi$ .

$$V_{H-h}^{\text{DR}} := \underbrace{\hat{Q}(s_h, \pi) + P_h(V_h + \gamma V_{H-h}^{\text{DR}} - \hat{Q}(s_h, a_h))}_{\text{cancel in expectation.}}$$

Policy Gradient |  $\{\pi_\theta : \theta \in \Theta\}$      $\max_{\theta} J(\pi_\theta) := \mathbb{E}_{d_\theta} [V^{\pi_\theta}(s)]$

(S)GD: compute (unbiased estim) of  $\nabla_{\theta} J(\pi_\theta)$ .

Question: can we estimate  $\nabla_{\theta} J(\pi_\theta)$  using data alone?

Answer: yes if we have data trajectories from  $\pi_\theta$ .

Derivation | Let  $\nabla J(\pi) := \nabla_{\theta} J(\pi_\theta)$ .

Goal:  $\nabla J(\pi) = \mathbb{E}_{\text{traj} \sim \pi} [\text{function of traj}]$   
 "policy gradient"

"REINFORCE":

$$\begin{aligned} \nabla J(\pi) &= \nabla_{\theta} \sum_{\tau} \underbrace{R(\tau)}_{\text{const.}} \cdot P^{\pi_\theta}(\tau) \\ &= \sum_{\tau} R(\tau) \cdot \nabla P^{\pi}(\tau) \\ &= \sum_{\tau} R(\tau) \cdot P^{\pi}(\tau) \cdot \nabla \log P^{\pi}(\tau) \end{aligned}$$

Let  $\tau = (s_1, a_1, r_1, \dots, s_H, a_H, r_H)$

$$R(\tau) = \sum_{h=1}^H \gamma^{h-1} r_h$$

$$P^{\pi}(\tau) = d_0(s_1) \cdot \pi(a_1 | s_1) \cdot P(s_2 | s_1, a_1) \cdot \pi(a_2 | s_2) \dots$$

$$= \sum_{\tau} R(\tau) \cdot P^{\pi}(\tau)$$

$$\nabla \log(\cancel{d_0(s_1)} \cdot \pi(a_1|s_1) \cdot P(s_2|s_1, a_1) \cdot \cancel{\pi(a_2|s_2)} \cdot \dots \cdot P(s_H|s_{H-1}, a_{H-1}) \cdot \cancel{\pi(a_H|s_H)})$$

$$\pi(a_1|s_1) \cdot \dots \cdot P(s_H|s_{H-1}, a_{H-1}) \cdot \pi(a_H|s_H)$$

$$= \sum_{\tau} P^{\pi}(\tau) \cdot R(\tau) \cdot \sum_{h=1}^H \nabla \log \pi(a_h|s_h)$$

$$= \mathbb{E}_{\tau \sim \pi} \left[ \underbrace{\left( \sum_{h=1}^H \gamma^{h-1} r_h \right)}_{\Delta} \cdot \underbrace{\left( \sum_{h=1}^H \nabla \log \pi(a_h|s_h) \right)}_{\Delta} \right]$$

"REINFORCE"

Remark: 1.  $\nabla_{\theta} J(\pi_{\theta}) = \lim_{\Delta \theta \rightarrow 0}$

$$\frac{J(\pi_{\theta+\Delta\theta}) - J(\pi_{\theta})}{\Delta\theta}$$

assumes scalar  $\theta$  but can be generalized.

$\Rightarrow$  "REINFORCE" =  $\lim_{\Delta\theta \rightarrow 0}$

$$\frac{\overset{\text{traj-}}{IS}(\pi_{\theta+\Delta\theta}) - IS(\pi_{\theta})}{\Delta\theta}$$

2. "Vanilla PG" =  $\lim_{\Delta\theta \rightarrow 0}$

$$\frac{\text{step-IS}(\pi_{\theta+\Delta\theta}) - \text{step-IS}(\pi_{\theta})}{\Delta\theta}$$

$$\rightarrow = \sum_{h=1}^H \nabla \log \pi(a_h|s_h) \left( \sum_{h'=h}^H \gamma^{h'-1} \cdot r_{h'} \right)$$

|| in expectation.

3.

$$\frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim d^{\pi}} \left[ \underbrace{Q^{\pi}(s,a)}_{\text{"PG thm"}} \cdot \nabla \log \pi(a|s) \right]$$

4.  $\nabla J(\pi) = \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim d^{\pi}} \left[ (Q^{\pi}(s,a) - f(s)) \nabla \log \pi(a|s) \right]$

$f$  is arbitrary. e.g.  $f(s) = V^\pi(s), A^\pi(s, a)$ .

"baseline"  $\uparrow$

$$\begin{aligned} & \mathbb{E}_{(s, a) \sim d^\pi} [f(s) \nabla \log \pi(a|s)] \\ &= \mathbb{E}_{s \sim d^\pi} \left( \mathbb{E}_{a \sim \pi(\cdot|s)} [f(s) \nabla \log \pi(a|s)] \right) \\ &= \mathbb{E}_{s \sim d^\pi} \left[ f(s) \cdot \mathbb{E}_{a \sim \pi(\cdot|s)} [\nabla \log \pi(a|s)] \right]. \end{aligned}$$

$$\begin{aligned} \underbrace{\mathbb{E}_{a \sim \pi} \left[ \frac{\nabla \pi(a|s)}{\pi(a|s)} \right]}_{\Delta} &= \sum_a \pi(a|s) \cdot \frac{\nabla \pi(a|s)}{\pi(a|s)} \\ &= \sum_a \nabla \pi(a|s) \\ &= \nabla \sum_a \pi(a|s) \\ &= \nabla \cdot 1 = \vec{0} \end{aligned}$$

5.  $J(\pi) = \frac{1}{1-\gamma} \mathbb{E}_{(s, a) \sim d^\pi} [Q^\pi(s, a) \log \pi(a|s)]$

can be estimated by e.g. TD.

"actor-critic"  
 $\pi \leftarrow \downarrow Q^\pi$