

# Application of IS in RL $M = (S, A, P, R, \gamma, do)$ .

Assume all trajectories terminate in  $H$  steps.

Data:  $(s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_H, a_H, r_H)$  where  $s_1 \sim do, a_t \sim \pi_b$

Goal: estimate  $J(\pi) := \mathbb{E} \left[ \sum_{h=1}^H \gamma^{h-1} r_h \mid \pi \right]$ .

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

$$\tau \sim \phi : \tau \sim \pi.$$

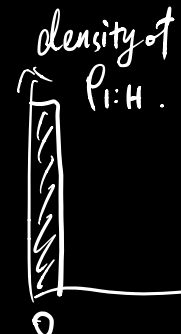
$$\tau \sim \phi : \tau \sim \pi_b.$$

$$f: \tau \mapsto \sum_{h=1}^H \gamma^{h-1} r_h \quad \mathbb{E}_{\tau \sim \phi} [f(\tau)]$$

$$J(\pi) = \mathbb{E}_{\tau \sim \phi} [f(\tau)] = \mathbb{E}_{\tau \sim \phi} \left[ \frac{p(\tau)}{q(\tau)} f(\tau) \right]$$

$$\frac{p(\tau)}{q(\tau)} = \frac{do(s_1) \cdot \pi(a_1|s_1) \cdot R(r_1|s_1, a_1) \cdot P(s_2|s_1, a_1) \cdot \pi(a_2|s_2) \dots R(r_H|s_H, a_H)}{do(s_1) \cdot \pi_b(a_1|s_1) \cdot R(r_1|s_1, a_1) \cdot P(s_2|s_1, a_1) \cdot \pi_b(a_2|s_2) \dots R(r_H|s_H, a_H)}$$

$$= \prod_{h=1}^H \frac{\pi(a_h|s_h)}{\pi_b(a_h|s_h)} = \prod_{h=1}^H \rho_h =: \rho_{1:H}$$



$$\boxed{\text{Traj-IS}} \left( \sum_{h=1}^H \gamma^{h-1} r_h \right) \cdot \rho_{1:H} \quad \leftarrow \text{never use!}$$

Special case:  $\pi_b$  is unif rand.  $\pi$  deterministic.  
 only traj w/ all actions matching  $\pi$  gets non-zero weight (weight =  $|A|^H$ , w/ prob.  $1/|A|^H$ ).

$$\boxed{\text{Step IS.}} \quad J(\pi) = \mathbb{E} \left[ \sum_{h=1}^H \gamma^{h-1} r_h \mid \pi \right] = \sum_{h=1}^H \gamma^{h-1} \mathbb{E} [r_h \mid \pi]$$

Idea: use IS to estimate  $\mathbb{E} [r_h \mid \pi]$  for  $h=1, 2, \dots, H$ .

How to estimate  $\mathbb{E} [r_1 \mid \pi] \leftarrow \rho_1 \cdot r_1$  ( $\rho_1 = \frac{\pi(a_1|s_1)}{\pi_b(a_1|s_1)}$ )

Final:  $\sum_{h=1}^H \gamma^{h-1} \rho_{1:h} \cdot r_h \leftarrow$  (compare: traj:  $\sum_{h=1}^H \gamma^{h-1} \rho_{1:H} \cdot r_h$ )

# Alt. interpretation of per-step IS

Let  $V_0 := 0$ .

Claim:  $V_H$  is precisely the step-IS estimator.

$$V_{H-h+1} := \underbrace{P_h}_{\text{unbiased}} \cdot \underbrace{(r_h + \gamma V_{H-h})}_{\text{unbiased for } Q^\pi(s_h, a_h)} \leftarrow \underbrace{\pi}_{\text{target policy}} \downarrow \text{unbiased for } \underbrace{V^\pi(s_{h+1})}_{\text{unbiased for } V^\pi(s_0)}$$

$$Q^\pi(s_h, \pi) \Rightarrow V_H \text{ unbiased estim of } V^\pi(s_0)$$

$$V^\pi(s_h) \Rightarrow V_H \text{ unbiased estim of } V^\pi(s_0)$$

Bandit DR:  $\hat{R}(s, \pi) + \rho \cdot (r - \hat{R}(s, a))$ . ideally good estimate of  $r$ .

DR for RL [Jiang & Lu '16]:  $\hat{Q}(s, a)$ : approximate  $Q^\pi$

$$V_{H-h+1}^{\text{DR}} = \hat{Q}(s_h, \pi) + \rho_h (r_h + \gamma V_{H-h}^{\text{DR}} - \hat{Q}(s_h, a_h))$$

## Policy Gradient (on-policy policy optimization).

$\{\pi_\theta : \theta \in \Theta\}$  policy class.

Goal: optimize  $J(\pi_\theta)$ . Alg: (S)GD on  $J(\pi_\theta)$ .

What we need: calculate stochastic gradient.  
i.e. expectation =  $\nabla_\theta J(\pi_\theta)$ .

on-policy:  $s_1, a_1, r_1, \dots, s_H, a_H, r_H \sim \pi_\theta$ .

Nontrivial:  $J(\pi_\theta) = V^{\pi_\theta}(s_0) = d_0^\top (I - \gamma P^\pi)^{-1} R^\pi$ . depends on MDP dynamics.

Abbrev:  $\pi = \pi_\theta$ .  $\nabla = \nabla_\theta$ .  $R: \tau \mapsto \sum_{h=1}^H \gamma^{h-1} r_h \leftarrow$

$$\nabla J(\pi) = \nabla_\theta \left( \sum_{\tau} R(\tau) P^\pi(\tau) \right)$$

$$\Delta P^\pi(\tau) = d_0(s_1) \pi(a_1 | s_1) R(r_1 | s_1, a_1) \dots$$

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

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

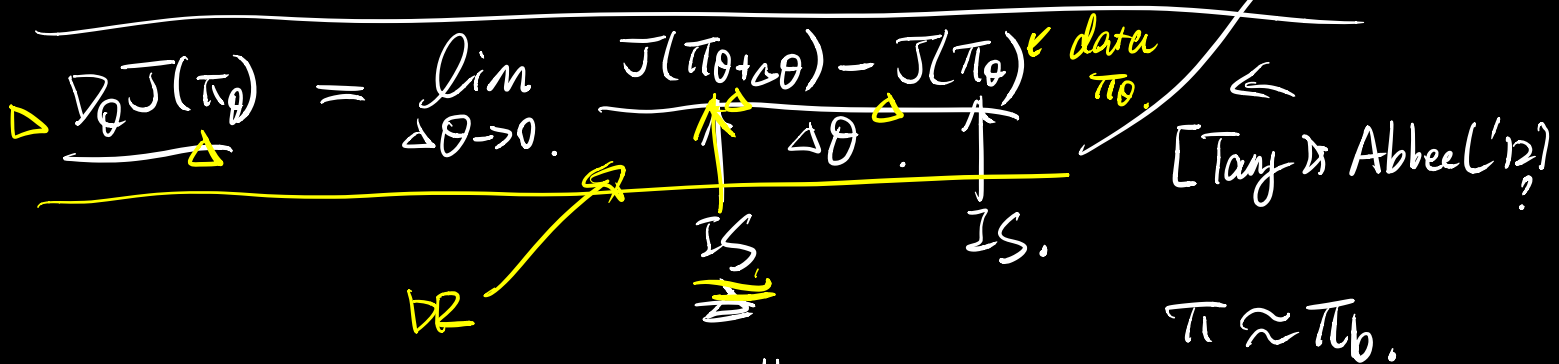
$\nabla y = \frac{y}{y} \nabla \log y$   
 b/c.  $\nabla \log y = \left[ \frac{1}{y} \right] \nabla y$

$$= \sum_{\tau} R(\tau) P^\pi(\tau) \nabla \left( \log(d_0(s_1) \cdot \pi(a_1 | s_1) \cdot \underbrace{R(r_1 | s_1, a_1)}_{\nabla=0} \cdot \underbrace{P(s_2 | s_1, a_1)}_{\nabla=0} \dots \underbrace{R(r_H | s_H, a_H)}_{\nabla=0} \right)$$

$$= \sum_{\tau} R(\tau) P^\pi(\tau) \nabla \left( \underbrace{\log d_0(s_1)}_{\nabla=0} + \log \pi(a_1 | s_1) + \log \dots \right)$$

$$= \sum_{\tau} R(\tau) P^\pi(\tau) \sum_{n=1}^H \nabla \log \pi(a_n | s_n)$$

$$= \mathbb{E} \left[ R(\tau) \cdot \sum_{n=1}^H \nabla \log \pi(a_n | s_n) \mid \pi \right] \leftarrow \text{"REINFORCE"}$$



REINFORCE:  $\left( \sum_{h=1}^H \gamma^{h-1} r_h \right) \sum_{n=1}^H \nabla \log \pi(a_n | s_n) \leftarrow \left( \sum_{h=1}^H \gamma^{h-1} r_h \right)$

$$\nabla J(\pi) = \nabla \left( \sum_{\tau} R(\tau) P^\pi(\tau) \right)$$

$$= \nabla \left( \sum_{h=1}^H \gamma^{h-1} \mathbb{E}[r_h | \tau] \right)$$

$\sum_{\tau_{pre}^h} r_h \cdot P^\pi(\tau_{pre}^h) \leftarrow P_{(i+1)}$

Standard PG:  $\sum_{h=1}^H \gamma^{h-1} \nabla \log(a_h | s_h) \cdot \left( \sum_{k=h}^H \gamma^{k-h} r_k \right)$  ← MC return

Actor-critic: estimate  $\hat{Q} \approx Q^\pi$ .  $Q^\pi(s_h, a_h)$

$$\sum_{h=1}^H \gamma^{h-1} \nabla \log(a_h | s_h) \hat{Q}(s_h, a_h)$$

"Baseline":  $\nabla \gamma^{h-1} \nabla \log(a_h | s_h) \left( \sum_{k=h}^H \gamma^{k-h} r_k - \underbrace{f(s_h)}_{\text{DR}}$

for any fixed  $f = \hat{V}^\pi(s_h)$ .

DR-PG [Huang & Jiang '20].

FQZ:  $\arg \min_f \mathbb{E} \left( \underbrace{f(s, a)}_f - r - \underbrace{V_{f_{k-1}}(s')}_{\pi^*} \right)^2$

Variant:  $(V, \pi)$

↑ State-function for  $V^*$

off-policy TD:  $\mathbb{E} \left[ \underbrace{V(s)}_{\substack{\text{argmin} \\ V \in \mathcal{V}}} - \frac{\pi(a|s)}{\pi_b(a|s)} \left( r + \underbrace{V_{k-1}(s')}_{V^\pi} \right) \right]^2$