

Importance Sampling

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

~~$$\{x_i\} \text{ iid } p.$$~~

$$\frac{1}{n} \sum_i f(x_i).$$

$$\{x_i\} \text{ iid } q.$$

$$\text{IS: } \frac{1}{n} \sum_{i=1}^n \frac{p(x_i)}{q(x_i)} f(x_i).$$

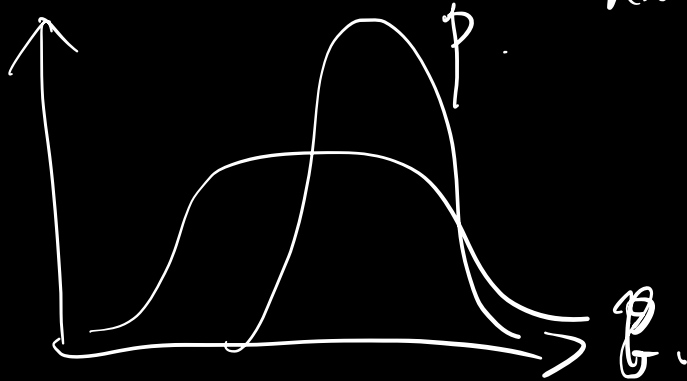
$$\mathbb{E}_{x \sim q} \left[\frac{p(x)}{q(x)} f(x) \right] = \int p(x) \cdot \frac{p(x)}{q(x)} f(x)$$

$$= \mathbb{E}_p[f].$$

Assumption:

$\forall x \text{ s.t. } p(x) > 0.$

we must have $q(x) > 0$.



$$\frac{p(x)}{q(x)} \leq C \quad \forall x.$$

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

(CR)

App. Contextual bandits

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

$$J(\pi) = \mathbb{E}_{\pi} [r]$$

$$\text{Data: } x = (s, a, r)$$

$$s \sim d_0, a \sim \pi_0 \\ r \sim R(\cdot | s, a) \\ x \sim q$$

$$J(\pi) = \mathbb{E}_q [f(x)]$$

$$(f(s, a, r) = r)$$

$$\text{IS: } \frac{p(x)}{q(x)} f(x) = \frac{p(s, a, r)}{q(s, a, r)} r$$

$$= \frac{d_0(s) \pi(a|s) R(r|s, a)}{d_0(s) \pi_0(a|s) R(r|s, a)} r$$

$$= \frac{\pi(a|s)}{\pi_0(a|s)} r$$

$$= \rho \cdot r$$

Case Study: $|A|=K$. $\pi_b(a|s) = \frac{1}{K}$.

π is deterministic ($\pi: S \rightarrow A$).

$$\frac{\mathbb{I}[\pi(s)=a]}{1/K} \cdot r = \begin{cases} K \cdot r & \text{if } \pi(s)=a \\ 0 & \text{if } \pi(s) \neq a. \end{cases}$$

\uparrow prob. $\quad \quad \quad \uparrow$ value.

$\text{w.p. } 1-1/K$

$$\begin{aligned} \text{Var}_{\mathcal{Q}}[p \cdot r] &\leq \mathbb{E}_{\mathcal{Q}}[p^2 r^2] \\ &\leq R_{\max}^2 \mathbb{E}_{\mathcal{Q}}[p^2] \\ &= R_{\max}^2 \mathbb{E}_{\mathcal{Q}}[p \cdot K] \\ &= K R_{\max}^2. \end{aligned}$$

$$p \cdot r \in [0, K R_{\max}].$$

π_b unif. π deterministic. $\Rightarrow r$ is const.

$$\text{Var}[p r] = r^2 \text{Var}[p].$$

$$= (\mathbb{E}[p^2] - \mathbb{E}[p]^2)$$

$$\begin{aligned}
&= v^2 (\mathbb{E}(K) - \mathbb{E}(K)) \\
&= v^2 \left(\mathbb{E}_{\pi_0} \left[\frac{\mathbb{I}[\pi(x) = a]}{1/k} \right] - 1 \right) \\
&= v^2 (k - 1).
\end{aligned}$$

$$\frac{1}{n} \sum_{i=1}^n \frac{\mathbb{I}[a_i = \pi(x_i)]}{1/k} \cdot r_i$$

$$= \frac{1}{n/k} \sum_{i: a_i = \pi(x_i)} r_i.$$

WIS: $\# \{ i: a_i = \pi(x_i) \}$

$$\frac{1}{n} \sum_{i=1}^n p_i r_i$$

↓

$$\sum_i p_i.$$

Doubly Robust: $\rho \cdot r$ ($\rho = \frac{\pi(a|s)}{\pi_b(a|s)}$)
 $\hat{R}(s, a)$.

$$\hat{R}(s, \pi) + \rho (r - \hat{R}(s, a))$$

$$\mathbb{E}_{a \sim \pi_b(\cdot|s)} [\hat{R}(s, \pi) - \rho \hat{R}(s, a)] = 0.$$

MDP | H-step finite horizon
 (S, A, P, R, d_0, H)

Φ

$s_1 \sim d_0, a_1 \sim \pi, r_1 = R(s_1, a_1), s_2 \sim P(\cdot | s_1, a_1), a_2 \sim \pi.$

$s_H, a_H, r_H.$

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=1}^H r_t \right]$$

OPE: data $(s_1, a_1, r_1, \dots, s_H, a_H, r_H)$

$a_{1:H} \sim \pi_b.$

\mathcal{D}

$$P(S_1, a_1, r_1, \dots, S_H, a_H, r_H)$$

$$\sum_{t=1}^H \gamma_t$$

$$Q(S_1, a_1, r_1, \dots, S_H, a_H, r_H)$$

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

$$\cancel{d_0(s_1)} \cdot \pi_b(a_1|s_1) \cdot \cancel{P(s_2|s_1, a_1)} \cdot \pi_b(a_2|s_2) \dots$$

$$= \prod_{t=1}^H \frac{\pi(a_t|s_t)}{\pi_b(a_t|s_t)} =: \prod_{t=1}^H \rho_t =: \rho_{1:H}$$

$$IS: \rho_{1:H} \cdot \sum_{t=1}^H \gamma_t$$

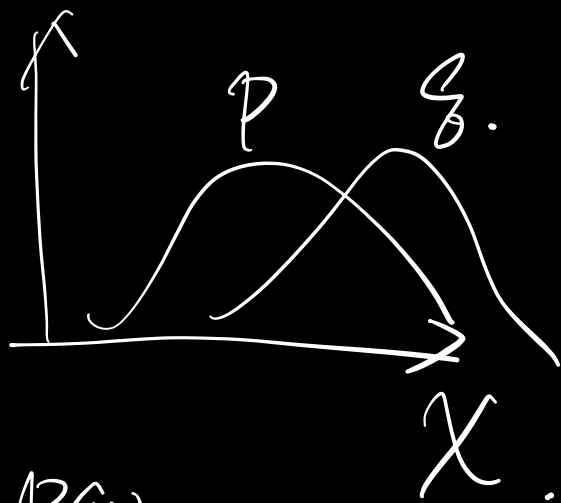
π_b unif. random. $\pi_b(a|s) = 1/K$.

deterministic π .

$$IS: \left(\frac{\prod_{t=1}^H \mathbb{I}[\pi(s_t) = a_t]}{1/K} \right) \left(\sum_{t=1}^H \gamma_t \right)$$

$$\prod_{t=1}^H \mathbb{I}[\pi(s_t) = a_t] \cdot \frac{1}{K}$$

$$K \cdot \Pi [\pi(s_t) = a_t, \forall t]$$



$$P_t = \frac{\pi(a_t | s_t)}{\pi_b(a_t | s_t)} \in C$$

$$P_{1:H} \in C^H$$

$$\max_x \frac{P(x)}{q(x)} \geq 1$$

$$C \leq 1 + \frac{1}{H}$$

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=1}^H r_t \right]$$

$$= \sum_{t=1}^H \mathbb{E}_{\pi} [r_t]$$

$$\mathbb{E}_{\pi} [r_1] = \mathbb{E}_{\pi_b} [P_{1:H} \cdot r_1]$$

$$= \mathbb{E}_{\pi_b} [P_1 \cdot r_1]$$

$$\mathbb{E}_{\pi} [r_t] = \mathbb{E}_{\pi_b} [P_{1:t} \cdot r_t]$$

Step-wise IS

$$\sum_{t=1}^H P_{1:t} r_t$$

(compare to $\sum_{t=1}^H P_{1:H} r_t$).

$$v_{H-t+1} := \boxed{P_t} (r_t + v_{H-t}).$$

$$\frac{\pi(a_t | s_t)}{\pi_b(a_t | s_t)}$$

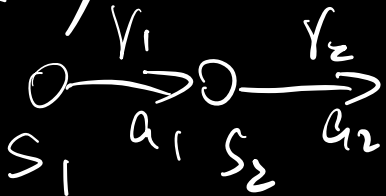
unbiased estimate of

$$Q_t^\pi(s_t, \pi) = V_t^\pi(s_t)$$

$$Q_t^\pi(s_t, a_t)$$

$$V_{t+1}^\pi(s_{t+1})$$

where $a_{t+1} \sim \pi_b$.



$$v_1 = v_0 = 0 = P_H \cdot r_H$$

$$\hat{Q} : S \times A \rightarrow \mathbb{R}$$

$$v_{H-t+1} := \hat{Q}(s_t, \pi) + P_t (r_t + v_{H-t} - \hat{Q}(s_t, a_t))$$

Policy Gradient

$$\max_{\pi_{\theta} \in \Pi} J(\pi_{\theta}).$$

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\pi_{\theta}} \left[\begin{matrix} \downarrow \\ \dots \\ s_1, a_1, r_1, \dots, s_H, a_H, r_H \end{matrix} \right]$$

$$\nabla_{\theta} J(\pi_{\theta}) \Big|_{\theta=\theta_0} = \mathbb{E}_{\pi_{\theta_0}} [\dots]$$

$$\nabla_{\theta} J(\pi_{\theta}).$$

$$\nabla J(\pi) = \nabla \mathbb{E}_{\pi} \left[\sum_{t=1}^H r_t \right].$$

$$= \nabla \mathbb{E}_{\tau \sim \pi} [R(\tau)].$$

$$= \nabla \sum_{\tau} \underbrace{P^{\pi}(\tau)}_{\downarrow} R(\tau).$$

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

$$= \sum_{\tau} \underbrace{P^{\pi}(\tau)}_{\downarrow} R(\tau) \nabla \log P^{\pi}(\tau).$$

$$= \mathbb{E}_{\pi} \left[\underline{R(\tau)} \nabla \log P^{\pi}(\tau) \right].$$

$(s_1, a_1, r_1, \dots, s_H, a_H, r_H)$
||
 τ .

$$\begin{aligned} \nabla \log P^{\pi}(z) &= \nabla \log \left(d_0(s_1) \cdot \pi(a_1|s_1) P(s_2|s_1, a_1) \right. \\ &\quad \left. \pi(a_2|s_2) \dots \pi(a_H|s_H) \right) \\ &= \cancel{\nabla \log d_0(s_1)} + \nabla \log \pi(a_1|s_1) + \cancel{\nabla \log P(s_2|s_1, a_1)} \\ &\quad + \dots + \nabla \log \pi(a_H|s_H) \\ &= \sum_{t=1}^H \nabla \log \pi(a_t|s_t) \end{aligned}$$

$$\nabla_{\theta} J(\pi_{\theta}) = \lim_{\Delta \theta \rightarrow 0} \frac{J(\pi_{\theta + \Delta \theta}) - J(\pi_{\theta})}{\Delta \theta}$$

$$\mathbb{E}_{\pi_{\theta}} [\dots]$$

"REINFORCE"

$$\nabla J(\pi) = \mathbb{E}_{\pi} \left[\left(\sum_{t=1}^H \nabla \log \pi(a_t|s_t) \right) \left(\sum_{t=1}^H \gamma_t \right) \right]$$

"Vanilla PG"

$$\nabla J(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=1}^H \nabla \log \pi(a_t|s_t) \sum_{t'=1}^H \gamma_{t'} \right]$$

$$\begin{aligned} & \overbrace{t'=t} \\ & \left(Q^{\overline{r}}(s_+, a_+) \right. \\ & \quad \left. - f(s_+) \right). \end{aligned}$$