

Apply IS to bandits

IS:  $f: X \rightarrow \mathbb{R}$ ,  $p, g \in \Delta(X)$ ,  $\mathbb{E}_p[f] = \mathbb{E}_g[\frac{p}{g} \cdot f]$

Contextual bandit: in each round,

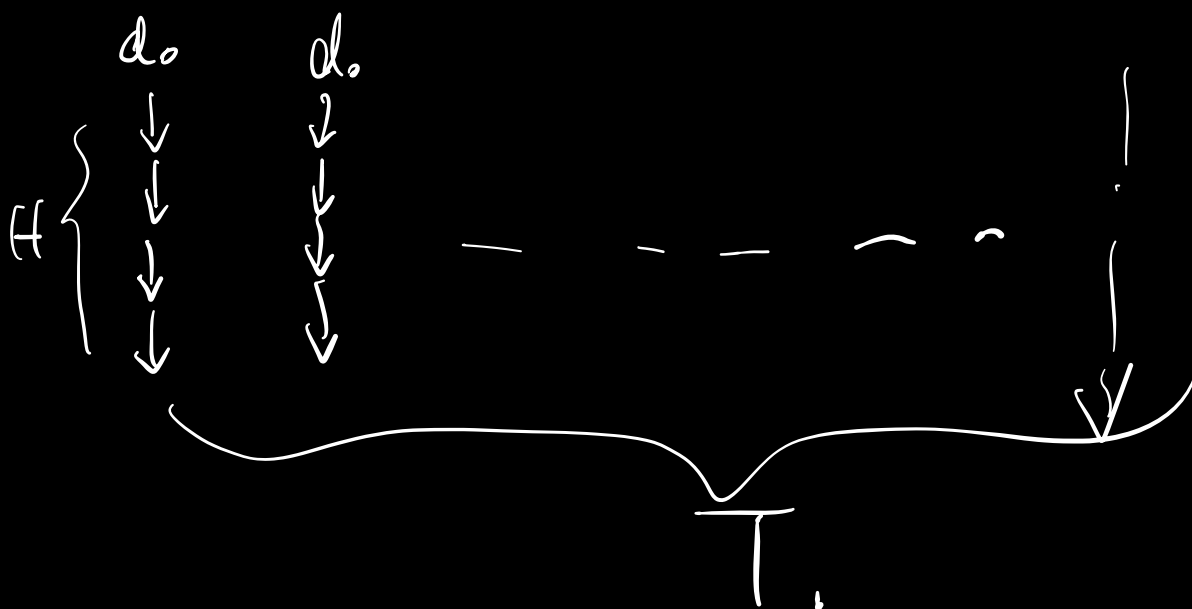
- (1) context  $x \sim d_0$
- (2) learner takes action  $a$

(3) reward  $r \sim R(\cdot | x, a)$   $\mathbb{E}_{x \sim d_0} [R(x, \pi(x))]$

Off-policy Evaluation (OPE)

Want to estimate  $J(\pi) = \mathbb{E}[r | \pi]$

Data:  $\{(x, a, r)\}$ ,  $a \sim \underline{\pi}_b \neq \pi$



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

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

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

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

$$p(x, a, r) = p(x) \cdot p(a|x) \cdot p(r|x, a)$$

$$= d_0(x) \cdot \pi(a|x) \cdot R(r|x, a)$$

$$q(x, a, r) = d_0(x) \cdot \pi_b(a|x) \cdot R(r|x, a)$$

$$\Rightarrow \rho = \frac{\pi(a|x)}{\pi_b(a|x)} \cdot \mathbb{E}_{\pi} [r] = \mathbb{E}_{\underline{\pi_b}} [\rho \cdot r]$$

IS.

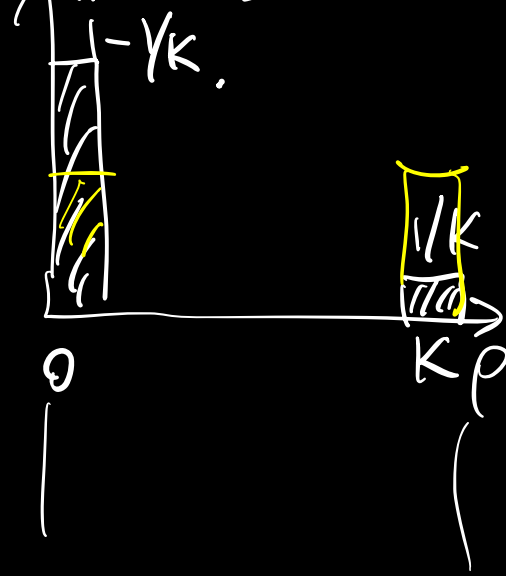
$$\mathbb{E}_{\underline{\pi_b}} [\rho] = 1 \quad |A| = K \quad \underline{\pi_b} = \text{Unif}$$

$\pi$  deterministic.

$$\frac{\pi(a|x)}{\pi_b(a|x)} = \begin{cases} K & \text{if } a = \pi(x) \\ 1 & \text{otherwise} \end{cases}$$

↑ prob. mass

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



Case study:  $R(\cdot|x,a)$  is deterministic and constant.

IS:  $\frac{\mathbb{I}[\pi(x)=a]}{1/K} \cdot C$

$\{(x^i, a^i, r^i)\}_{i=1}^n$

$$\begin{aligned}
 \frac{1}{n} \sum_{i=1}^n \frac{\mathbb{I}[\pi(x^i)=a^i]}{1/K} \cdot C &= \frac{K}{n} \cdot \sum_{i: \pi(x^i)=a^i} C \\
 &= \frac{\sum_{i: \pi(x^i)=a^i} C}{n/K}.
 \end{aligned}$$

Alt:

$$\frac{\sum_{i: \pi(x^i)=a^i} c}{|\{i: \pi(x^i)=a^i\}|}$$

$$\rho^i = \frac{\pi(a^i|x^i)}{\pi_b(a^i|x^i)}$$

$$IS: \frac{1}{n} \sum_i \rho^i \cdot r^i$$

$$WIS: \sum_i \frac{\rho^i r^i}{\sum_j \rho^j}$$

Doubly Robust:  $IS: \rho \cdot r$

DR: Input:  $\hat{R}(x, a)$

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

Unbiased:  $E_{\pi_b} [\hat{R}(x, \pi) - \rho \hat{R}(x, a)] = 0$

Remark:  $\rho \cdot r = \frac{\pi(a|x)}{\pi_b(a|x)} \cdot r$

$$(x, a, r, \pi_b(a|x)) \rightarrow$$

"logging prob" proposal

MDP: finite-horizon episodic MDP:

$$(S, A, P, R, d_0, H).$$

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

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

OPE: data  $\tau = (s_1, a_1, s_2, a_2, \dots, s_H, a_H) \sim \pi_b$ .

$$p: \tau \sim \pi.$$

$$q: \tau \sim \pi_b.$$

$$f: \tau \mapsto R(s_1, a_1) + R(s_2, a_2) + \dots + R(s_H, a_H) \\ =: R(\tau).$$

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

$$= \mathbb{E}_{\tau \sim q} \left[ \frac{p(\tau)}{q(\tau)} \cdot R(\tau) \right].$$

~~$$p(\tau) = d_0(s_1) \cdot \pi(a_1 | s_1) \cdot P(s_2 | s_1, a_1) \cdot \pi(a_2 | s_2) \dots$$~~

~~$$q(\tau) = d_0(s_1) \cdot \pi_b(a_1 | s_1) \cdot P(s_2 | s_1, a_1) \cdot \pi_b(a_2 | s_2) \dots$$~~

$$= \prod_{h=1}^H \frac{\pi(a_h | s_h)}{\pi_0(a_h | s_h)} =: \prod_{h=1}^H \rho_{h,} =: \rho_{1:H}$$

$\pi_s = \text{unif.}$       $\pi$ : deterministic.

$$\rho_{1:H} = \prod_{h=1}^H \frac{\mathbb{I}[a_h = \pi(x_h)]}{1/K} = K^H \prod_{h=1}^H \mathbb{I}[a_h = \pi(x_h)]$$

$$= K^H \cdot \mathbb{I}[\forall h, a_h = \pi(x_h)]$$

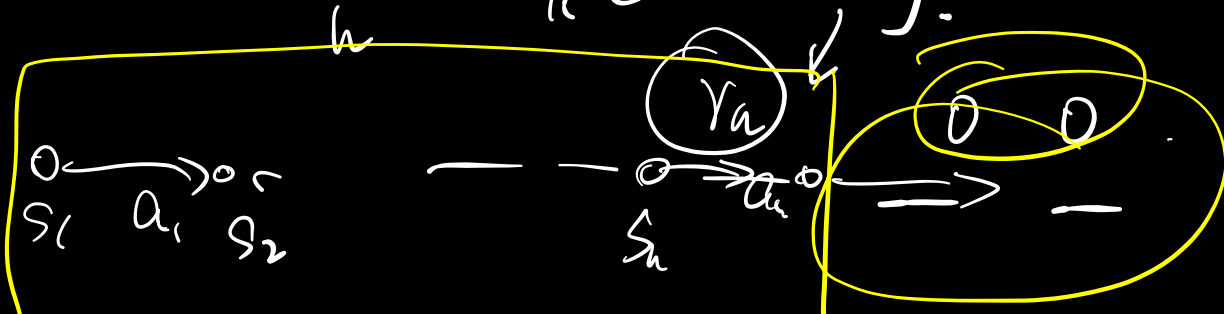
$$\prod_{h=1}^H \rho_{h,} =$$

$$\left(1 + O\left(\frac{1}{H}\right)\right)^H \approx e.$$



$$J(\pi) = \mathbb{E}_{\pi} \left[ \sum_h R(s_h, a_h) \right] \rightarrow R(\tau)$$

$$= \sum_h \mathbb{E}_{\pi} [R(s_h, a_h)]$$



$$P_{i:h} = \gamma_h \cdot V_h$$

$$P_{i:H} \sum_n \gamma_n \quad \mapsto \quad \sum_h P_{i:h} \cdot V_h$$

$$\mathbb{E}_{\pi} [R(s_n, a_n)] = \mathbb{E}_{\pi_h} \left[ \frac{p(s_n, a_n)}{z(s_n, a_n)} R(\cdot) \right]$$

$$\frac{d_{\pi}^{\pi}(s, a)}{d_{\pi}^{\pi_0}(s, a)}$$