

Given  $\pi$ . and data:  $(s_1, a_1, r_1, \dots, s_H, a_H, r_H)$   
 we want to evaluate.  
 $a_t: \# \sim \pi_b$   
 "target/eval policy". "behavior policy"  
 logging.

Importance Sampling / reweighting  
 Inverse Propensity Score (IPS).

$X$  large finite space.  
 $p, g \in \Delta(X)$ .  $f: X \rightarrow [0, 1]$ . known.  
 Want to estimate  $E_p[f] := E_{x \sim p}[f(x)]$ .  
 $X_1, X_2, \dots, X_n \sim p, \Rightarrow \frac{1}{n} \sum_{i=1}^n f(x_i)$   
 "Monte Carlo".

What if:  $X_1, X_2, \dots, X_n \sim g$ .

Estimator:  $\frac{1}{n} \sum_{i=1}^n \frac{p(x_i)}{g(x_i)} f(x_i)$ .  
 Long weight/

density ratio / IPS

"Estimator":  $\frac{P(x)}{g(x)} f(x)$ , where  $X \sim g$ .

Unbiased:  $E_g \left[ \frac{P(x)}{g(x)} f(x) \right] = E_p [f]$ .

$$= \sum_x g(x) \frac{P(x)}{g(x)} f(x) = E_p [f].$$

1.  $\left\| \frac{P}{g} \right\|_\infty < +\infty$ .

$$\frac{\left\| \frac{P}{g} \right\|_\infty}{\varepsilon^2}.$$

2: Hoeffding:  $\left\| \frac{P}{g} \right\|_\infty \sqrt{\frac{1}{2a} \ln \frac{2}{\delta}}$ .

3:  $E_g \left[ \frac{P(x)}{g(x)} \right] = \sum_x g(x) \frac{P(x)}{g(x)} = 1$ .

4:  $\frac{P(x)}{E_g(x)} f(x) \in [0, \left\| \frac{P}{g} \right\|_\infty]$ .

$\overline{\text{Var}}(\text{IPS}) = O\left(\left\| \frac{P}{g} \right\|_\infty X \in [0, a]\right)$ .  
 $\text{Var} X \leq O(a^2)$ .

$\uparrow$

$\frac{P(x)}{Q(x)}$  is large.