# Importance Sampling (IS)

Problem: estimate $\mathbb{E}_{x \sim p}[f(x)]$, where $p \in \Delta(X)$ and $f: X \to \mathbb{R}$.

MC: draw $x_1, \dots, x_n \overset{iid}{\sim} p$, estm: $\frac{1}{n} \sum_{i=1}^{n} f(x_i)$. | *if f is bounded. Hoeffding's ineq applies.*

(in short, will write: $x \sim p$. $f(x)$ as the estimator).

→ Example: MC policy eval. $J(\pi) := \mathbb{E}_\pi \left[ \sum_{t=1}^{H} \gamma^{t-1} r_t \right]$.
$x \longleftrightarrow \tau = (s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_H, a_H, r_H)$
$p \longleftrightarrow \tau \sim \pi$. (from do).
$f \longleftrightarrow \tau \longmapsto \sum_{t=1}^{H} \gamma^{t-1} r_t$.

IS: no access to $x \sim p$. but can sample $x \sim g$. where $g \in \Delta(X)$.

if: $\forall x$ where $p(x) > 0$, we have $g(x) > 0$. then: ←

$x \sim g$, $\boxed{\dfrac{p(x)}{g(x)}} f(x)$ ← | *IS estimator. IPS, IW.* | *IS: design g. IPS, IW: $\frac{p(x)}{g(x)}$.*

*imp. weights/ratio, density ratio, etc.*

$\boxed{\text{Claim: IS is unbiased}}$   assume $X$ is finite.

Proof: $\mathbb{E}_{x \sim g}\left[\dfrac{p(x)}{g(x)} f(x)\right] = \sum_{x \in X} g(x) \cdot \left(\dfrac{p(x)}{g(x)} f(x)\right) = \mathbb{E}_{x \sim p}[f(x)]$

ideally: $\dfrac{p(x)}{g(x)}$ should be small. | if $\max_x \dfrac{p(x)}{g(x)} \le C$. $\quad f(x) \in [-B, B]$. $\quad \downarrow \quad [-CB, CB]$.

$\boxed{\text{Fact}: \mathbb{E}_{x \sim g}\left[\dfrac{p(x)}{g(x)}\right] = 1}$   $\left\|\dfrac{p}{g}\right\|_\infty$   $\mathbb{E}_{x \sim g}\left[\dfrac{p(x)^2}{g(x)^2}\right]$.

Example: OPE in contextual bandit

$S$: contexts, $A$: actions, $R: S \times A \to \Delta([0,1])$.
let $d_0 \in \Delta(S)$ be the ctx dist.   *behavior/logging policy.*
Have data: $\{(s, a, r)\}$: $s \sim d_0$, $a \sim \pi_b(\cdot | s)$, $r \sim R(\cdot | s, a)$
  *target/eval policy.*
Goal: estimate $J(\pi) := \mathbb{E}[r | \pi]$.
IPS: $\dfrac{\pi(a|s)}{\pi_b(a|s)} \cdot r := \underline{\rho} \cdot r$.

Proof of unbiasedness: let $(s,a,r) \sim g \iff s \sim d_0, \boxed{a \sim \pi_b,} r \sim R$ (data).

let $(s,a,r) \sim p \iff s \sim d_0, \boxed{a \sim \pi}, r \sim R$

$$J(\pi) = \mathbb{E}_{(s,a,r) \sim p}[r] = \mathbb{E}_{(s,a,r) \sim g}^{\text{``data''}} \left[ \frac{p(s,a,r)}{g(s,a,r)} \cdot r \right].$$

$$= \mathbb{E}_{(s,a,r) \sim g} \left[ \frac{p(s) \cdot p(a|s) \cdot p(r|s,a)}{g(s) \cdot g(a|s) \cdot g(r|s,a)} \Big/ r \right].$$

$$= \mathbb{E}_{(s,a,r) \sim g} \left[ \frac{d_0(s) \cdot \pi(a|s) \cdot R(r|s,a)}{d_0(s) \cdot \pi_b(a|s) \cdot R(r|s,a)} \cdot r \right] = \mathbb{E}_{(s,a,r) \sim g} \left[ \frac{\pi(a|s)}{\pi_b(a|s)} r \right].$$

---

Variance of IS: Consider special case where

· $\pi$ is deterministic. $a = \pi(s)$      let $K = |A|$.

· $\pi_b$ is unif over $A$.    (or $\pi_b \sim U$).   $\impliedby$.

IS: $a \sim U$, $\rho \cdot r$, where $\rho = \dfrac{\pi(a|s)}{\pi_b(a|s)} = \dfrac{\mathbb{I}(a = \pi(s))}{1/K}$.

---

Let's further assume that $\underline{\underline{r}}$ is const (ind. of $s,a$, has no rndness).
$\boxed{r_0}$

$$\text{Var}[\rho \underline{r}] = r_0^2 \, \text{Var}[\rho].$$

$$= r_0^2 \cdot \left( E[\rho^2] - (E[\rho])^2 \right) = r_0^2 \left( E[\rho^2] - 1 \right).$$

$$= r_0^2 \left( \mathbb{E} \left[ \frac{\mathbb{I}(a = \pi(s))}{1/K^2} \right] - 1 \right).$$

$$= r_0^2 \left( K \, \mathbb{E} \left[ \frac{\mathbb{I}(a = \pi(s))}{1/K} \right] - 1 \right) = \boxed{r_0^2 (K-1)}.$$

---

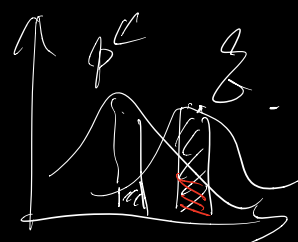IS: $\{(s_i, a_i, r_i)\}_{i=1}^n$

$$\frac{1}{n} \sum_{i=1}^n \frac{\mathbb{I}(a_i = \pi(s_i))}{1/K} \, r_i = \boxed{\frac{1}{n/K}} \sum_{i: a_i = \pi(s_i)} r_i.$$

Can we address this?

· Improvement 1: WIS. (self-normalized IPS).

In the above spec. case:

$$\frac{\sum_{i: a_i = \pi(s_i)} r_i}{\left| \{ i : a_i = \pi(s_i) \} \right|} .$$



General case:

$$\frac{\sum_{i=1}^{n} \rho_i \cdot r_i}{\underbrace{\sum_{i=1}^{n} \rho_i}_{\xrightarrow{\text{expe}} n}}$$

(biased but consistent),

(unbiased).

- Improvement 2: DR (doubly robust)

In the spec case: $\rho \cdot \boxed{r} \implies \hat{r}_0 + \rho(r - \boxed{r_0})$

$$\mathbb{E}[\rho \cdot r_0] = r_0 \cdot \mathbb{E}[\rho] = r_0$$

General case: $\boxed{\hat{R}}$ $S \times A \to \mathbb{R}$.
arbitrary function.

$$DR = \mathbb{E}_{a' \sim \pi}\left[ \hat{R}(s, a') \right] + \rho \cdot (r - \hat{R}(s, a)).$$

"control variate"

[Dudik, Li, Langford].

**Why DR if we have good $\hat{R}$?**

→ DR has low var. if $\hat{R} \approx R$.

→ DR is always unbiased even if $\hat{R}$ is poor.

→ IS is a special case of DR: $\hat{R} \equiv 0$.

regress
$(s, a) \to r$.
on separate data
to fit $\hat{R}$.

Applications to RL.