

Importance Sampling "target policy/eval policy"

Motivating Task: estimate $J(\pi) := \mathbb{E}_{s \sim d_0} [V^\pi(s_0)]$.

- w/ on-policy data: roll-out w/ π , $\sum_{t=1}^H \gamma^{t-1} r_t$.
- How about off-policy data? "Off-policy Evaluation" (OPE)

→ Solution: can use Expected SARSA to learn Q^π .

$$\rightarrow J(\pi) = \mathbb{E}_{s \sim d_0} [Q^\pi(s, \pi)]$$

Problem: tabular TD-alg does not scale; must use func. approx. → BIAS!!!

Question: unbiased estimator for OPE??

Answer: YES! Importance Sampling.

Intro to IS / Importance Weighting / Inverse Propensity Scores (IPS)

Problem: want to estimate $\mathbb{E}_{x \sim p} [f(x)]$.

MC: $x_i \stackrel{iid}{\sim} p$, $\frac{1}{n} \sum_{i=1}^n f(x_i)$.

Setup: can't sample from p , but can sample from q .

IS: $x_i \stackrel{iid}{\sim} q$. $\frac{1}{n} \sum_{i=1}^n \frac{p(x_i)}{q(x_i)} f(x_i) \xrightarrow{\text{unbiased}} \mathbb{E}_{x \sim p} [f(x)]$.

$$x \sim q, \quad \frac{p(x)}{q(x)} f(x)$$

Proof of unbiasedness: want to show. $\mathbb{E}[\text{estimator}] = \text{what you want}$.
(Assume the domain of x is finite)

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

Remarks:

- MC is a special case: $p = q$. $\frac{p(x)}{q(x)} = 1$.
- $\frac{p(x)}{q(x)}$: "importance weight/ratio" converts $q \rightarrow p$.
- $\mathbb{E}_q \left[\frac{p(x)}{q(x)} \right] = 1$.

Application to Contextual Bandits (CB).

Episodic RL w/ $H=1$: in every "episode".

- $x \sim d_0$. (x is often called context).
- agent takes action $a \sim \pi_b$ (behavior policy).
- receives reward $r \sim R(\cdot | x, a)$.

OPE: what is $J(\pi) := \mathbb{E}[r | \pi]$ for $\pi \neq \pi_b$

Idea: each data point (x, a, r) .

- func of interest $(x, a, r) \mapsto r$.
- distribution of interest: $(x, a, r) \sim \pi$.

more formally: $x \sim d_0, a \sim \pi(\cdot | x), r \sim R(\cdot | x, a)$

estimand $J(\pi) = \mathbb{E}_{x \sim d_0, a \sim \pi(\cdot | x), r \sim R(\cdot | x, a)} [r]$.

" $(x, a, r) \sim p$ "

Actual data is sampled from a different dist. " q ":

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

\Rightarrow Directly apply IS:

$$\begin{aligned} \frac{p(x, a, r)}{q(x, a, r)} \cdot r &= \frac{p(x) \cdot p(a|x) \cdot p(r|x, a)}{q(x) \cdot q(a|x) \cdot q(r|x, a)} \cdot r \\ &= \frac{d_0(x) \cdot \pi(a|x) \cdot R(r|x, a)}{d_0(x) \cdot \pi_b(a|x) \cdot R(r|x, a)} \cdot r = \frac{\pi(a|x)}{\pi_b(a|x)} \cdot r \end{aligned}$$

Remarks:

- π_b must be stochastic in general, π can be stoch or deterministic
- Do not need full knowledge of π_b . "logging probability".
all you need is to record $(x, a, r, \pi_b(a|x))$.
- Useful to consider special case when π_b is uniform
 $\&$ π is deterministic

importance weight $\rho := \frac{\pi(a|x)}{\pi_b(a|x)} = \frac{\mathbb{I}[a = \pi(x)]}{1/|A|}$

$$= \begin{cases} |A| & \text{if } a = \pi(x). \\ 0 & \text{o.w.} \end{cases}$$

therefore, roughly $1/|A|$ of the data is useful.

- IS works the best when $\pi \approx \pi_b$.

Application to TD(0):

• Standard TD(0) learns V^π from on-policy data

e.g. (s, a, r, s') require $a \sim \pi$
 form 1-step target: $r + \gamma \cdot V(s') \xleftarrow{\text{update}} V(s)$

• Question: can we do TD(0) using off-policy data.

e.g. (s, a, r, s') , $a \sim \pi_b \neq \pi$.
 $\frac{\pi(a|s)}{\pi_b(a|s)} \cdot (r + \gamma \cdot V(s')) \xleftarrow{\text{update}} V(s)$

$$\begin{aligned} \mathbb{E} \left[\frac{\pi(a|s)}{\pi_b(a|s)} \cdot (r + \gamma V(s')) \right] &= \sum_a \pi_b(a|s) \cdot \frac{\pi(a|s)}{\pi_b(a|s)} \cdot \mathbb{E}[r + \gamma V(s') | a] \\ &= \sum_a \pi(a|s) \mathbb{E}[r + \gamma V(s') | a] = \mathbb{E}_{a \sim \pi} [r + \gamma V(s')] \end{aligned}$$

Application to DPE in MDPs

• Data traj. $\tau := (s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_H, a_H, r_H)$
 where $\forall t, a_t \sim \pi_b$.

• want to estimate $J(\pi) := \mathbb{E} \left[\sum_{t=1}^H \gamma^{t-1} r_t \mid \forall t, a_t \sim \pi \right]$.

Apply IS: let p be the dist. of τ if $\forall t, a_t \sim \pi$.
 let q be the dist. of τ if $\forall t, a_t \sim \pi_b$.

IS estimator: $\frac{p(\tau)}{q(\tau)} \cdot \sum_{t=1}^H \gamma^{t-1} r_t$.

$$p(\tau) = d_0(s_1) \cdot \pi(a_1|s_1) \cdot \cancel{P(s_2|s_1, a_1)} \cdot \pi(a_2|s_2) \cdot \cancel{R(r_2|s_2, a_2)} \\ \uparrow \\ P(s_3|s_2, a_2) \cdot \dots \cdot P(s_H|s_{H-1}, a_{H-1}) \cdot \pi(a_H|s_H) \cdot \cancel{R(r_H|s_H, a_H)}$$

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

$$\Rightarrow \frac{p(\tau)}{q(\tau)} = \frac{\pi(a_1|s_1) \cdot \pi(a_2|s_2) \cdot \dots \cdot \pi(a_H|s_H)}{\pi_b(a_1|s_1) \cdot \pi_b(a_2|s_2) \cdot \dots \cdot \pi_b(a_H|s_H)} \\ = p_1 \cdot p_2 \cdot p_3 \cdot \dots \cdot p_H = \underline{p_{1:H}}$$

Remarks:

• Special case: π_b is unif. random. & π is deterministic.

$$p_t = \frac{\pi(a_t|s_t)}{\pi_b(a_t|s_t)} = \frac{\mathbb{I}[a_t = \pi(s_t)]}{1/|A|} = \begin{cases} 0, & \text{o.w.} \\ |A|, & \text{if } a_t = \pi(s_t) \end{cases}$$

$$\Rightarrow p_{1:H} \cdot \sum_{t=1}^H \gamma^{t-1} r_t = \begin{cases} 0, & \text{if } \exists t, s_t, a_t \neq \pi(s_t). \\ |A|^H \cdot \sum_{t=1}^H \gamma^{t-1} r_t, & \text{if } \forall t, a_t = \pi(s_t). \end{cases}$$

occurs w/ prob. $1/|A|^H$

Improvement to (per-trajectory) ZS.

Idea: want to estimate $\mathbb{E}[\sum_{t=1}^H \gamma^{t-1} r_t | \pi]$.

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

Estimate $\mathbb{E}[r_t | \pi]$: states & actions after timestep t .

"PDIS"
"step-wise IS"

↓
is irrelevant.

$$p_{1:t} \cdot r_t \Rightarrow \sum_{t=1}^H \gamma^{t-1} p_{1:t} \cdot r_t$$

Application to multi-step TD. $s_t, a_t, r_t, s_{t+1}, a_{t+1}, V_{t+1}, s_{t+2}$.

• on-policy 2-step TD. $(s, [a, r, s', a', [r', s'']])$
 $a, a' \sim \pi$.

2-step TD target: $r + \gamma \cdot r' + \gamma^2 V(s'')$ ← update $V(s)$.

• off-policy? $(s, a, r, s', a', r', s'')$ where $a, a' \sim \pi_b$

importance weighted target: $\frac{\pi(a|s) \cdot \pi(a'|s')}{\pi_b(a|s) \cdot \pi_b(a'|s')} \cdot (r + \gamma \cdot r' + \gamma^2 V(s''))$

Improve? $\frac{\pi(a|s)}{\pi_b(a|s)} \cdot r + \frac{\pi \cdot \pi}{\pi_b \cdot \pi_b} \cdot \gamma (r' + \gamma V(s'))$.

$$\mathbb{E}_{\substack{a \sim \pi \\ a' \sim \pi}} [r + \gamma r' + \gamma^2 V(s'')]$$

$$= \mathbb{E}_{\substack{a \sim \pi \\ a' \sim \pi}} [r] + \gamma \mathbb{E}_{\substack{a \sim \pi \\ a' \sim \pi}} [r' + \gamma V(s'')]$$

$$= \mathbb{E}_{a \sim \pi} [r] + \dots$$