# off-policy Evaluation

The goal is to evaluate a target policy using samples from another policy, called a behavior policy.

(policy of interest)

(policy generating behavior)

$$\pi^t : \text{target policy}$$
$$\pi^b : \text{behavior policy}$$

**Benefits**

Enables using samples from
- old policies
- exploratory policies
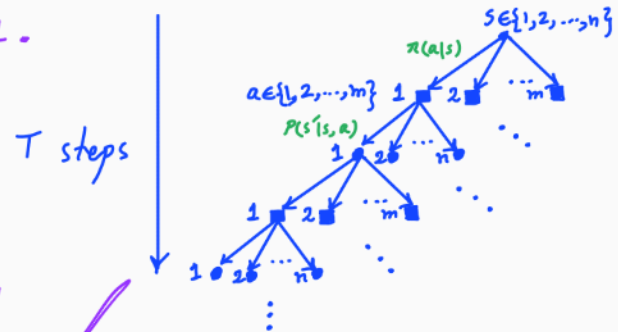- other agents/demonstrations

**Note:** Consider a sample finite trajectory from an MDP

rewards conform to the deterministic reward function

$$r_0 = R(s_0, a_0) \quad r_1 = R(s_1, a_1) \quad r_{T-1} = R(s_{T-1}, a_{T-1})$$

$$\tau_T = (s_0, a_0, r_0, s_1, a_1, r_1, \ldots, s_{T-1}, a_{T-1}, r_{T-1}, s_T)$$

generated based on an initial distribution $\mu_0$ and following a stationary, stochastic policy $\pi$.

T steps

$s \in \{1,2,\ldots,n\}$
$\pi(a|s)$
$a \in \{1,2,\ldots,m\}$
$P(s'|s,a)$

* The probability of $\tau_T$ being realized is

$$\tau_T$$

$$\mathbb{P}(\tau_T | \mu_0, \pi, P) = \mathbb{P}(s_0, a_0, s_1, a_1, \ldots, s_{T-1}, a_{T-1}, s_T | \mu_0, \pi, P)$$

$$= \mathbb{P}(s_T | s_0, a_0, s_1, a_1, \ldots, s_{T-1}, a_{T-1}, \mu_0, \pi, P)$$

$P(s_T | s_{T-1}, a_{T-1})$

$$\mathbb{P}(a_{T-1} | s_0, a_0, s_1, a_1, \ldots, s_{T-1}, \mu_0, \pi, P)$$

$\pi(a_{T-1} | s_{T-1})$

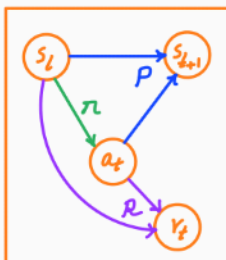$$\mathbb{P}(s_0, a_0, s_1, a_1, \ldots, s_{T-1} | \mu_0, \pi, P)$$

$$\tau_{T-1}$$

repeat the factorization

$$= \vdots$$

$$= \mu_0(s_0) \prod_{t=0}^{T-1} \pi(a_t | s_t) P(s_{t+1} | s_t, a_t)$$

$$\tau_T \longrightarrow \tau_{T-1} \longrightarrow \tau_{T-2} \longrightarrow \cdots \longrightarrow \tau_1 \longrightarrow \tau_0$$

unroll trajectory over time

$s_t$ $\xrightarrow{P}$ $s_{t+1}$
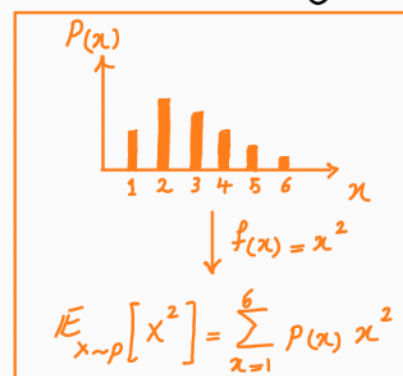$\pi$
$a_t$
$R$
$r_t$

**Importance sampling:** A general technique for evaluating properties of a target distribution, e.g., the expectation of a function over that distribution based on samples from another distribution.

Let $X$ be a discrete RV distributed according to probability measure $P$.

$$\mathbb{E}_{X \sim P}[f(X)] = \sum_{x \in X} P(x) f(x)$$

$$\boxed{\mathbb{E}_{X \sim P}[X] = \sum_{x \in X} P(x) x}$$

$$\boxed{\begin{array}{c} P(x) \\ \\ \downarrow f(x) = x^2 \\ \mathbb{E}_{X \sim P}[X^2] = \sum_{x=1}^{6} P(x) x^2 \end{array}}$$

Having $n$ i.i.d. samples $x_1, x_2, \ldots, x_n$ of $X$, drawn from $P(x)$, a simple Monte Carlo method can estimate this expected value by the empirical mean

$$\hat{\mathbb{E}}_{X \sim P}[f(X)] = \frac{1}{n} \sum_{i=1}^{n} f(x_i)$$

Now, consider another distribution over $X$ captured by probability measure $Q$.

$$\mathbb{E}_{X \sim P}[f(X)] = \sum_{x \in X} P(x) f(x) = \sum_{x \in X} Q(x) \underbrace{\frac{P(x)}{Q(x)}} f(x)$$

Likelihood ratio weight $w(x)$

$$= \mathbb{E}_{X \sim Q}\left[ \underbrace{\frac{P(x)}{Q(x)} f(x)}_{g(x)} \right]$$

requirement:
$Q(x) > 0$ if $P(x) f(x) \neq 0$

Having $n$ i.i.d. samples $x_1, x_2, ..., x_n$ of $X$, drawn from $Q(X)$, a simple Monte Carlo method can estimate this expected value by the empirical mean

$$\hat{\mathbb{E}}_{X \sim P}[f(x)] = \hat{\mathbb{E}}_{X \sim Q}\left[\frac{P(X)}{Q(X)}f(x)\right] = \frac{1}{n}\sum_{i=1}^{n}\frac{P(x_i)}{Q(x_i)}f(x_i)$$

$w(x_i)$ should be known

---

Importance sampling for off-policy evaluation: Let $G_T$ be the sample return corresponding to sample trajectory $\tau_T$ generated by $\pi^b$.

$$w(\tau_T) = \frac{P(\tau_T \mid \mathcal{M}_0, \pi^t, P)}{P(\tau_T \mid \mathcal{M}_0, \pi^b, P)} = \frac{\mathcal{M}_0(s_0)\prod_{t=0}^{T-1}\pi^t(a_t \mid s_t)P(s_{t+1}\mid s_t, a_t)}{\mathcal{M}_0(s_0)\prod_{t=0}^{T-1}\pi^b(a_t \mid s_t)P(s_{t+1}\mid s_t, a_t)}$$

$P$ is unknown but it cancels out $\longleftarrow$

$$= \prod_{t=0}^{T-1}\frac{\pi^t(a_t \mid s_t)}{\pi^b(a_t \mid s_t)}$$

Off-policy MC evaluation:

the return of the trajectory $\downarrow$

$w$

$$\mathbb{E}_{\pi^t, P}[G_T] \approx \hat{\mathbb{E}}_{\pi^b, P}\left[\overbrace{\prod_{t=0}^{T-1}\frac{\pi^t(a_t \mid s_t)}{\pi^b(a_t \mid s_t)}}^{} G_T\right] = \frac{1}{n}\sum_{i=1}^{n}w_i G_i$$

$$\hat{v}^{\pi_t}(s) \longleftarrow \hat{v}^{\pi_t}(s) + \alpha\left(\prod_{t=0}^{T-1}\frac{\pi^t(a_t \mid s_t)}{\pi^b(a_t \mid s_t)}G_T - \hat{v}^{\pi_t}(s)\right)$$

$t = t_s$ if the trajectory does not start from $s$

# Off-policy TD evaluation:

$$\mathbb{E}_{\pi^t, p}\left[ R(S_t, A_t) + \gamma V(S_{t+1}) \mid S_t = s \right]$$

$$= \mathbb{E}_{\substack{A_t \sim \pi^t(s) \\ S_{t+1} \sim P(\cdot \mid S_t, A_t)}}\left[ R(S_t, A_t) + \gamma V(S_{t+1}) \mid S_t = s \right]$$

$$\approx \hat{\mathbb{E}}_{\substack{A_t \sim \pi^b(s) \\ S_{t+1} \sim P(\cdot \mid S_t, A_t)}}\left[ \frac{\pi^t(A_t \mid S_t)}{\pi^b(A_t \mid S_t)}(R(S_t, A_t) + \gamma V(S_{t+1})) \mid S_t = s \right]$$

$$\hat{V}^{\pi_t}(s) \leftarrow \hat{V}^{\pi_t}(s) + \alpha \left( \frac{\pi^t(A_t \mid S_t)}{\pi^b(A_t \mid S_t)}(R(S_t, A_t) + \gamma \hat{V}^{\pi_t}(S_{t+1})) - \hat{V}^{\pi_t}(S_t) \right)$$

---

## off-policy Optimization

The goal is to optimize a target policy using samples from another policy, called a behavior policy.

We can use a generalized policy iteration approach, where the policy evaluation step is performed using importance sampling.

# off-policy MC method with soft exploration

Iterative method:

- Initialize $\hat{Q}$

- Initialize $\pi^t$ as a greedy policy with respect to $\hat{Q}$

- Initialize $G(s,a) = \emptyset$ for all $s \in S, a \in A$

- Repeat

  - Build $\pi^b$ as a soft policy with respect to $\hat{Q}$

  - Generate a sample trajectory $\tau$ using $\pi^b$

  *policy evaluation*
  - For $(s,a) \in S \times A$ s.t. $(s,a) \in \tau$:

    - $t_{s,a} = $ first time $(s,a)$ is visited

    - $G = \prod\limits_{t=t_{s,a}}^{T} \dfrac{\pi^t(a_t|s_t)}{\pi^b(a_t|s_t)} \left( \sum\limits_{t=t_{s,a}}^{T} \gamma^{t-t_{s,a}} R(s_t, a_t) \right)$

    - $G(s,a) \leftarrow G(s,a) \cup \{G\}$

    - $N(s,a) \leftarrow N(s,a) + 1$

    - $\hat{Q}(s,a) = \dfrac{1}{N(s,a)} \sum\limits_{G \in G(s,a)} G$

  *policy improvement*
  - Build $\pi^t$ as a greedy policy with respect to $\hat{Q}$

# off-policy TD method (Q learning)

Iterative method:

- Initialize $\hat{Q}$ s.t. $\hat{Q}(s,a) = 0$ for terminal states

- Initialize $\pi^t$ as a greedy policy with respect to $\hat{Q}$

- Repeat

  - Sample $s_0 \in S$

  - Repeat until the episode terminates

    - Build $\pi^b$ as a soft policy with respect to $\hat{Q}$
      (e.g., $\mathcal{E}$-greedy)

    - Take action $a$ according to $\pi^b(a|s)$ at $s_0$

    - Observe reward $R(s,a)$ and next state $s'$

    - policy evaluation and improvement
      $\hat{Q}(s,a) \longleftarrow \hat{Q}(s,a) + \alpha [R(s,a) + \gamma \max_{a' \in A} \hat{Q}(s',a') - \hat{Q}(s,a)]$

    - policy improvement · Build $\pi^t$ as a greedy policy with respect to $\hat{Q}$

    - Update the current state