

①. MC estimation \rightarrow Prediction \rightarrow Estimate value function
 \hookrightarrow Control: $\#$ how to opt. make use of trajectory
 \hookrightarrow find optimal policy.

②) off-policy estimation:-

Prediction :- (π) : V^π value fn for policy π .

Generate trajectory: 1). $S_0, a \sim \pi, S_1, a \sim \pi, \dots$

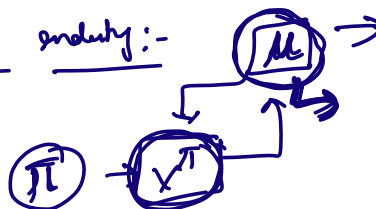
off-policy:-

Samples will

be generated from

different policy μ

1. gaming industry:-



Both are stochastic policies :-



off-policy estimation:-

① $(\Omega, \mathcal{F}, \mu)$ $X: \Omega \rightarrow \mathbb{R}$

② $(\Omega, \mathcal{F}, \pi)$

Q1: $E_{\mu} [X] = \sum_{x \in X} \mu(x) \cdot x$

$E_{\pi} \left[\frac{\pi(x)}{\mu(x)} \cdot X \right] \stackrel{!}{=} \sum_{x \in X} \mu(x) \cdot \frac{\pi(x)}{\mu(x)} \cdot x$

$$z = \sum_{x \in \mathcal{X}} \pi(x) \cdot x$$

$$= E_{\pi} [x]$$

$$p(x) = \frac{\pi(x)}{\mu(x)} ; \quad p(x) \cdot x$$

(*)

$$V^{\pi}(s) = E_{(s_0, s_1, s_2, \dots)} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t) \right]$$

$$= \sum_{s_0, s_1, \dots} P_{\pi}(s_0, s_1, s_2, \dots) \underbrace{\sum_{t=0}^{\infty} \gamma^t r(s_t)}_x$$

$$V^{\pi}(s) = E_{\mu} \left[f(s_0, s_1, s_2, \dots) \right]$$

$$\sum_{t=0}^{\infty} \gamma^t r(s_t)$$

$$V^{\mu}(s) = \sum_{s_0, s_1, \dots} P_{\mu}(s_0, s_1, s_2, \dots) \sum_{t=0}^{\infty} \gamma^t r(s_t)$$

Q:- If you get trajectories from μ , can you compute value function of π ?

$$p(s_0, s_1, \dots) = \frac{P_{\pi}(s_0, s_1, \dots)}{P_{\mu}(s_0, s_1, \dots)}$$

$$P_{\pi}(s_0, a_0, s_1, a_1, s_2) \cdot$$

$$P(s_0) \times \pi(a_0 | s_0) \times P(s_1 | s_0, a_0) \times \pi(a_1 | s_1) \times P(s_2 | s_1, a_1)$$

$$P(x_1, x_2, x_3) = P(x_1) \cdot P(x_2 | x_1) \cdot P(x_3 | x_1, x_2)$$

$$(\pi), (\mu)$$

$$P_{\mu}(s_0, a_0, s_1, a_1, s_2) :-$$

$$P(s_0) \times \mu(a_0 | s_0) \times P(s_1 | s_0, a_0) \times \mu(a_1 | s_1) \times P(s_2 | s_1, a_1)$$

$$\frac{P_{\pi}(s_0, a_0, s_1, a_1, s_2)}{P_{\mu}(s_0, a_0, s_1, a_1, s_2)} = \frac{\pi(a_0 | s_0) \times \pi(a_1 | s_1)}{\mu(a_0 | s_0) \times \mu(a_1 | s_1)} \quad \rightarrow \text{we know}$$

$$\begin{aligned} \mu: S &\rightarrow A \text{ deterministic} \\ \mu: S \times A &\rightarrow [0, 1] \rightarrow \text{stochastic} \end{aligned}$$

environment:

$$\pi: S, a \rightarrow s'$$

$$(\pi). \text{ I know policy } (\mu) \Rightarrow (s_0, a_0 \sim \mu, s_1, a_1 \sim \mu, s_2, \dots)$$

$$\pi \Rightarrow$$

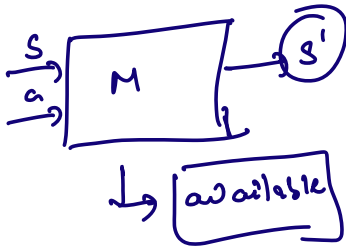
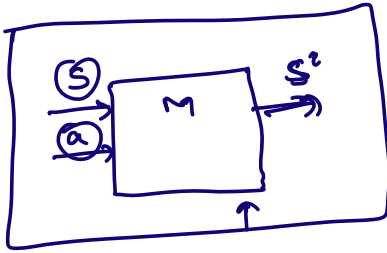
on policy:-

$$\begin{aligned} \pi :- & \begin{aligned} & 1. s_0, a_0, r_0, s_1, a_1, r_1, \dots \\ & \quad (R_1) = \sum_{t=0}^{\infty} \gamma^t r_t \\ & 2. s_0, a_1, r_0^* \dots \\ & \quad (R_2) \\ & 3. \quad (R_3) \\ & 4. \dots \end{aligned} \\ V^{\pi} = & \frac{\sum R_i}{n} = \left(\frac{R_1}{n} + \frac{R_2}{n} + \frac{R_3}{n} + \dots \right) \end{aligned}$$

Generating ✓

offspring:

μ :



1.

R_1^M $(s_0, a \sim \mu, s_0, s, a \sim \mu, \dots)$

2.

R_2^M

3.

R_3^M

$\sqrt{\pi}$

$i^?$

$$\frac{R_1^M + R_2^M + \dots}{n}$$

$$\sqrt{\pi} \leftarrow \frac{\sum_i s_i R_i^M}{n}$$