

- Expectation over actions :- known
- Expectation over states : not known.

$$\pi(a_1|s) ; \pi(a_2|s), \dots$$

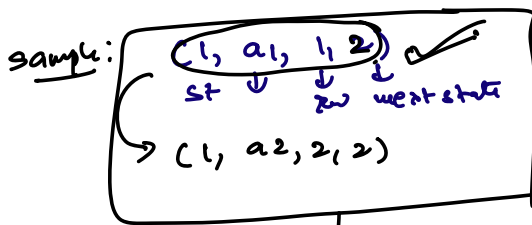
$$\sum_a \frac{Q(s', a) \cdot \pi(a|s)}{\pi(a|s)}$$

Q4:

MDP with 2 states 1 & 2.

$\pi$

$$V^\pi(2) = 2 + V^\pi(1).$$



$$\begin{aligned} V_0 &= 0 \\ \alpha &= \frac{1}{n+1} \\ \gamma &= 1. \end{aligned}$$

$$V_2(1) \neq V_2(2)$$

timestep

$$V_{n+1}(i) = V_n(i) + \alpha ( \dots )$$

$$V_{n+1}(i_n) = (1 - \alpha_n) V_n(i_n) + \alpha_n [r_n + \gamma V_n(i_{n+1})].$$

$$(i_n, \pi_n, i_{n+1})$$

First sample:

$$(1, a_1, 1, 2) ; (1, a_2, 2, 2)$$

$$V_0(1) = 0 ; V_0(2) = 0$$

$$V_1(1) = (1 - 1) V_0(1) + 1 (1 + 1(0)) = 1$$

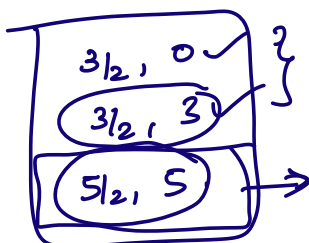
$$\rightarrow V_1(2) = 2 + 1 = 2$$

Second sample:-

$$V_2(1) = (1 - \frac{1}{2}) (1) + \frac{1}{2} (2 + 1 \cdot 2) \quad \left| \quad \frac{1}{2} + \frac{1}{2} (2) = 3/2 \right.$$

$$= \frac{1}{2} + 2 = 5/2$$

$$V_2(2) = 2 \cdot 5/2 = 5$$



⑤.  $C(1, a_1, 1, 2, a_2, 2, 1, a_3, 4, 2)$ .  $\gamma = 0.5$

$V(i)$

1-step TD  $\downarrow$

$\boxed{1} + \gamma(0)$   $\rightarrow$  ①

$1 + 0.5(2) + \gamma^2(0)$   $\rightarrow$  ②

$1 + 0.5(2) + (0.5)^2(4) + (0.5)^3(0)$   $\rightarrow$  ③

Off-policy Prediction / control :-

Prediction :-  $\pi$  : given obs:  $V^\pi$

Model-information is not known.

• samples are obtained from different policy  $\mu$ .

Monte-Carlo estimation :-

$\tau \rightarrow \{ \underline{s_1}, \underline{a_1}, \underline{s_2}, \underline{a_2}, \underline{s_3}, \underline{a_3}, \dots \}$  trajectory given.

$\begin{cases} \underline{\underline{Pr_\mu(\tau)}} :- Pr\{s_1, a_1, s_2, \dots\} \\ \underline{\underline{Pr_\pi(\tau)}} :- \end{cases}$

$= \underbrace{Pr\{s_1\}}_{\mu(a_1|s_1)} \times \underbrace{P(s_2|s_1, a_1)}_{\mu(a_2|s_2)} \times \dots$

$Pr(s_1) \times \pi(a_1|s_1) \times P(s_2|s_1, a_1) \times \pi(a_2|s_2) \times \dots$

$J_\tau = \frac{Pr_\pi(\tau)}{Pr_\mu(\tau)}$

$\pi(a_1|s_1) \times \pi(a_2|s_2) \times \dots$

$\mu(a_1|s_1) \times \mu(a_2|s_2) \times \dots$

$(0.001)^{1000}$

$Pr\{x_1, x_2, x_3\} = Pr\{x_1\} \cdot$

$Pr\{x_2|x_1\} \cdot$

$Pr\{x_3|x_1, x_2\}$

$(0.19)^{1000}$

$(0.001)^{1000}$

$$\begin{aligned}
 \cdot \quad E_{\mu}[\bar{x}] &= \underbrace{\mu(x_1)} \cdot x_1 + \underbrace{\mu(x_2)} \cdot x_2 \\
 \cdot \quad E_{\pi}[\bar{x}] &= \pi(x_1) \cdot x_1 + \pi(x_2) \cdot x_2 \\
 &= \mu(x_1) \cdot \underbrace{\frac{\pi(x_1)}{\mu(x_1)}} \cdot x_1 + \mu(x_2) \cdot \underbrace{\frac{\pi(x_2)}{\mu(x_2)}} \cdot x_2 \\
 &= E_{\mu}[\underbrace{\rho(x)} \cdot x]
 \end{aligned}$$

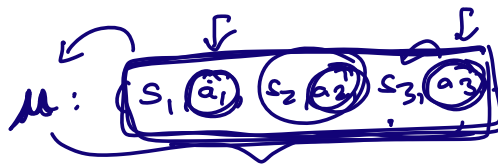
$\rho(x) = \frac{\pi(x)}{\mu(x)}$

$$V^{\pi}(s) = E[G_{\pi}] \rightarrow \sum_{t=0}^{\infty} \gamma^t r_t \rightarrow r(s_t, a_t \sim \pi).$$

$$V^{\mu}(s) = E[G_{\mu}]$$

$$\rightarrow \pi: (s_1, a_1, s_2, a_2, s_3, a_3, \dots)$$

$r_1 \quad r_2 \quad r_3$



$$r_1 + \gamma r_2 + \gamma^2 r_3$$

$$V^{\pi}(s_1) \leftarrow$$

$$(r_1 + \gamma r_2 + \gamma^2 r_3) \cdot \rho_{\pi}$$

$$\rho_{\pi}(r_1 + \gamma r_2 + \gamma^2 r_3)$$

Ordinary  
I.S.

$$V(s) = \frac{\sum_{\tau} G_{\tau} \rho_{\tau}}{|\tau|}$$

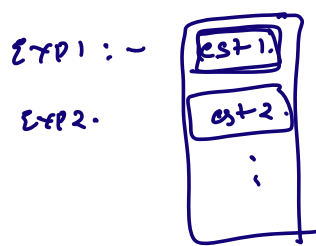
$$\text{contin: } \frac{\sum G_{\tau}}{|\tau|}$$

$\rightarrow$  infinite variance.

weighted  
I.S.

$$V(s) = \frac{\sum_{\tau} \rho_{\tau} G_{\tau}}{\sum_{\tau} \rho_{\tau}}$$

$$\frac{\rho_{\tau}}{\sum_{\tau} \rho_{\tau}}$$



$$E[x] = x^*$$

$$\rightarrow x_0, x_1, x_2, \dots$$

$$x_1, x_2, x_3, \dots$$

$$E[x_i] = x^*$$

Incremental implementation :-

$$v_n = \frac{\sum_{k=1}^{n-1} w_k q_k}{\sum_{k=1}^{n-1} w_k}$$

$$\left\{ \begin{array}{l} v_{n+1} = v_n + \frac{w_n}{c_n} [q_n - v_n] \\ c_{n+1} = c_n + w_{n+1} \end{array} \right\}$$

$$\frac{x_1 + x_2 + \dots + x_n}{n}$$

$$x_{n+1} = (1-\alpha) x_n + \alpha r_n$$

$$\downarrow$$

$$\frac{1}{n+1}$$

Prediction

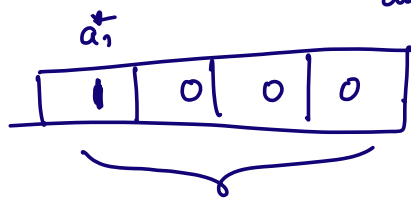
Control :-

Prob :- compute the optimal policy, with samples derived from a diff. behavior policy 'u'.

target policy :- fixed  $\pi$   
 $\rightarrow$  move it.

$Q(s, a)$  estimates

at each time-steps :-



$$\left\{ \begin{array}{l} \textcircled{1} \pi(s) \leftarrow \arg \max_a Q(s, a) \\ \textcircled{2} \pi(s) \leftarrow \begin{cases} \arg \max Q(s, a) & \text{w.p } \epsilon \\ \text{random} & \text{w.p } 1-\epsilon \end{cases} \\ \textcircled{3} \text{UCB.} \end{array} \right\}$$

Summarize:-

- (1). Importance Sampling :- (weights tractable appropriately).
- (2).  $\hookrightarrow$  eliminates unknown & only on known quantities.
- (3). weighted i.i.s : ordinary i.i.s has infinite variance.

$$\pi : \underline{V}^{\pi} : E[\underline{G}^{\pi}]$$

$\mu$  : Some samples are used

$$\text{if } \pi(a|s) > 0 \quad \text{then} \quad \mu(a|s) > 0. \quad \checkmark$$

$$\boxed{\mu(a^*|s) > 0 \text{ but } \pi(a^*|s) = 0}$$

Example 4.3 :- Gamblers problem