

(*) Monte-carlo estimation :-

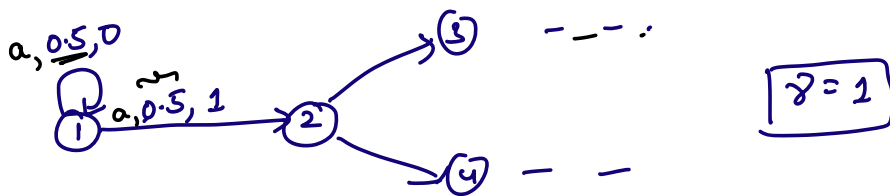
1). given samples $\begin{cases} \text{First visit} \\ \text{Every visit} \end{cases}$

2). prediction, control.
 \hookrightarrow optimal policy.

(2). off-policy :- target policy, behavior policy.

$\begin{cases} \text{prediction :} \\ \text{control} \rightarrow \text{optimal policy.} \end{cases}$

Can we still improve the sample efficiency :-



given policy π , estimate value functions $V^\pi(i)$?

$V^\pi(2) = 2$. ✓ $\rightarrow \pi$

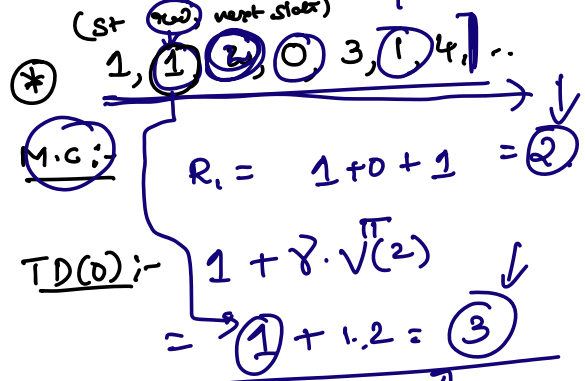
Actual estimate of $V^\pi(i)$?

$$(*) V^\pi(i) = \frac{1}{2} [0 + V^\pi(i)] + \frac{1}{2} [1 + 2]$$

$$\Rightarrow \frac{V^\pi(i)}{2} = \frac{3}{2}$$

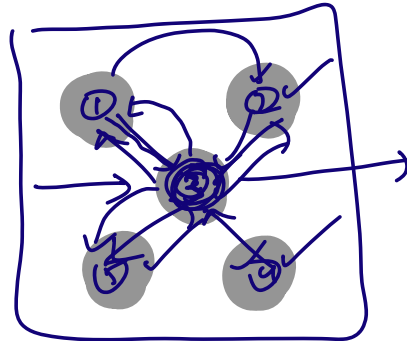
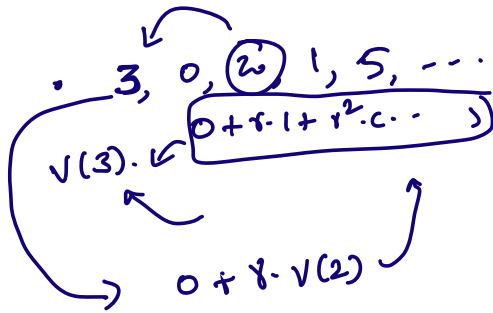
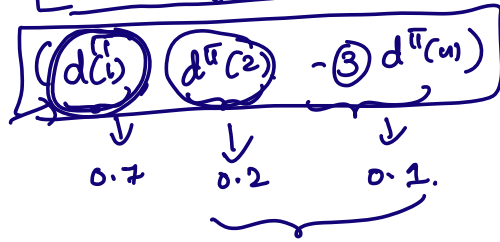
$$\Rightarrow \boxed{V^\pi(i) = 3}$$

one sample :-
 (st, act, next state)



MDP + policy \rightarrow Markov chain.

\rightarrow stationary distn.



$$V(3) \approx V^{\pi}(3)$$

1, $a_1, \gamma, 3, a_2, \gamma_2, 4, \dots$

$$* V(1) = \gamma + \gamma \cdot V(3) \rightarrow$$

[TD(0) update:]

$$x = E[c]$$

c_1, c_2, c_3, \dots samples

$$\rightarrow x_{n+1} = (1 - \alpha_n) x_n + \alpha_n (c_n)$$

$x_n \rightarrow x$

$(V_1(s_1), r_1(c_1), \dots)$

Prediction :-

$$\pi : V^{\pi} \quad (s, a \sim \pi, s')$$

$$V(s) = E_{s'} [\gamma + \gamma V^{\pi}(s')]$$

$$V_{n+1}(s) = (1 - \alpha_n) V_n(s) + \alpha_n [\gamma + \gamma V_n(s')]$$

$$\underline{V_n(s)} \rightarrow \underline{V^\pi(s)}, \forall s \in S.$$

$n \rightarrow \infty$

Q-learning :- control :- Optimal policy.

optimal Q-Bellman eqn:-

$$Q(s, a) = E_{s'} [r + \gamma \max_b Q(s', b)].$$

\downarrow
 x

\downarrow
 C

$$(s, a, r, s').$$

\downarrow
given.

$$Q_{n+1}(s, a) = (1 - \alpha_n) Q_n(s, a) + \alpha_n [r + \gamma \max_b Q_n(s', b)].$$

\downarrow
Step-Size

$$\sum \alpha_n = \infty ; \sum \alpha_n^2 < \infty$$

\downarrow

$$\left\{ \frac{1}{n} \right\}.$$

$$1, \frac{1}{2}, \frac{1}{3}, \left(\frac{1}{4} \right), \dots$$

⑦. please derive the same for SARSA & expected SARSA.

⑦ TD (0) prediction }
 ⑦ Q-learning. }

Mountain car :-

