

## Ex 4: Monte-Carlo Methods.

(Q.1)

a) Incremental first-visit MC policy evaluation (estimating value function of a policy in MDP)

→ First-visit MC prediction, for estimating  $V \approx V_{\pi}$

Input: policy  $\pi$  to be evaluated.

Initialize:

$v(s) \in \mathbb{R}$ , arbitrarily, for all  $s \in S$

$N(s) \leftarrow 0 \quad \forall s \in S$

Loop forever (for each episode):

Generate an episode following:  $s_0, a_0, r^1, s_1, a_1, r_2, \dots, s_{T-1}, a_{T-1}, r_T$

$G \leftarrow 0$

Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :

$G \leftarrow G + R_{t+1}$

Unless  $s_t$  appears in  $s_0, s_1, \dots, s_{t-1}$ :

Append  $G$  to  $\text{Returns}(s_t)$ .

$V(s_t) = \text{incremental-avg}(\text{Returns}(s_t), G)$ .

def incremental-avg(returns( $s_t$ ),  $G$ ):

$N(s_t) += 1$

$V(s_t) = V(s_t) + (G - V(s_t)) / N(s_t)$

return  $V(s_t)$ .

this fn → returns avg for each state,  $V(s_t)$ .

- updates it after each new return is observed.

$N(s_t) =$  keeps track of number of times state  $s_t$  has been visited.

b) To alter the pseudocode for Monte Carlo Es to maintain just the mean & a count (for each state-action pair) & update them incrementally,

Monte Carlo ES (Exploring starts), for estimating  $\pi \approx \pi^*$

Initialize:

$\pi(s) \in A(s)$  (arbitrarily), for all  $s \in S$

$q(s, a) \in \mathbb{R}$ . (arbitrarily), for all  $s \in S, a \in A(s)$

$N(s, a) \leftarrow 0 \quad \forall s \in S, a \in A(s)$ .

Loop forever (for each episode):

choose  $s_0 \in S, a_0 \in A(s_0)$  randomly such that all pairs have probability  $> 0$

Generate an episode from  $s_0, a_0, \text{ following } \pi: s_0, a_0, r_1, \dots$

$G \leftarrow 0$

$\dots, s_{T-1}, a_{T-1}, r_T$ .

Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :

$G \leftarrow \gamma G + r_{t+1}$

Unless the pair  $s_t, a_t$  appears in  $s_0, a_0, s_1, a_1, \dots, s_{t-1}, a_{t-1}$ :

Append  $G$  to Returns( $s_t, a_t$ )

$N(s_t, a_t) += 1$

$q(s_t, a_t) = q(s_t, a_t) + (G - q(s_t, a_t)) / N(s_t, a_t)$

$\pi(s_t) \leftarrow \operatorname{argmax}_a q(s_t, a)$

## Q.2) first-visit vs. every-visit

a) We know (in general)  $\rightarrow$

every-visit MC - expected to produce more accurate estimates  
of state values than first-visit MC.

$\downarrow$   
it considers all returns following a visit to a state.

but first-visit MC only considers return following the  
first visit to a state.

if we consider Black Jack problem,

there is constantly changing state in each episode

- The state only appears once, even if every-visit MC is used.
- Hence, it gets same result as using first-visit MC.

b) first-visit estimator:  $\gamma = 1$

- only considers the return following first-visit to a state.
- return following the 1<sup>st</sup> visit to the non-terminal state is 10.

$$\therefore V(s) = 10$$

$$G_t = R_{t+1} + \gamma G_{t+1}$$

$$G_0 = R_1 + G_1 = 10 \quad G_1 = R_2 + G_2 = 9.$$

$$G_{10} = 0.$$

• Every-visit estimator:

- considers all returns following visits to a state.
- Here, only state that is visited is the nonterminal state.

$$\begin{aligned}\therefore V(s) &= \frac{(1+2+3+\dots+10)}{10} = \frac{n(n+1)}{2} \cdot \frac{1}{10} \\ &= \cancel{\frac{10 \times 11}{2}} \times \frac{1}{10} = 5.5\end{aligned}$$

→ first-visit estimator produces a higher estimate as → it only considers return following first visit to the state. i.e. highest possible return.

→ But every-visit produces a lower estimate as it considers all returns following visits to the state, including returns from where episodes ends early

### Q. 5) OFF-policy Method -

a) Suppose we have sequence of returns  $G_1, G_2, \dots, G_{n-1}$ , all starting in the same state & each with a corresponding random weight  $w_i$

$$V_n = \frac{\sum_{k=1}^{n-1} w_k G_k}{\sum_{k=1}^{n-1} w_k}, \quad n \geq 2 \quad \dots \quad 5.7.$$

$$\text{To derive} \rightarrow V_{n+1} = V_n + \frac{w_n}{C_n} [G_n - V_n], \quad n \geq 1$$

derivation:

$$\begin{aligned} \rightarrow V_{n+1} &= \frac{\sum_{k=1}^n w_k g_k}{\sum_{k=1}^n w_k}, \quad n \geq 1 \\ &= \frac{\sum_{k=1}^{n-1} w_k g_k}{\sum_{k=1}^n w_k} + \frac{w_n g_n}{\sum_{k=1}^n w_k} \\ &= \frac{\sum_{k=1}^{n-1} w_k}{\sum_{k=1}^n w_k} \times \frac{\sum_{k=1}^{n-1} w_k g_k}{\sum_{k=1}^n w_k} + \frac{w_n g_n}{\sum_{k=1}^n w_k} \\ &= V_n - \frac{V_n \times w_n}{\sum_{k=1}^n w_k} + \frac{w_n g_n}{\sum_{k=1}^n w_k} \\ &= V_n + \left( G_n - V_n \right) \times \frac{w_n}{\sum_{k=1}^n w_k} \rightarrow c_n \\ &= V_n + \frac{w_n}{c_n} \times (G_n - V_n) \dots (5.8) \end{aligned}$$

- b) off-policy MC works by correcting bias in unweighted MC algorithm due to difference between behavior policy  $b$ , target policy  $\pi$ .
- can be done by weighting the unweighted returns by the importance sampling ratio,  $\pi(A_t|s_t) / b(A_t|s_t)$
  - But in the boxed algorithm for off-policy MC control,  $w_{\text{update}}$  only involves  $\frac{1}{b(A_t|s_t)}$
  - it's only used to update the weights  $c$ ,  $\rightarrow$  used to normalize the importance sampling ratio.
  - $c_{n+1} = c_n + \alpha_n \cdot \frac{1}{b(A_n|s_n)}$   $\leftarrow$  shows how  $w_{\text{update}}$  is used to normalize the importance sampling ratio.  
 $\downarrow$   
 weights  $c \rightarrow$  sum to 1.
  - $w_{\text{update}}$  in boxed algorithm for off-policy MC control is only used to normalize the importance sampling ratios.
  - As the <sup>importance</sup> sampling ratios have been normalized, they can be used to weight the unweighted returns
  - it will produce weighted returns, this can be used to update the value estimates using MC update rule.