

Monte Carlo Methods

Whereas Dynamic Programming learns the task in a iterative bottom-up approach, Monte Carlo method learns via episodic top-down approach. This is useful when you have to learn $v_*(s)$ or $q_*(s, a)$ via sampling because you have no model $p(s', r' | s, a)$. So, it is less complex because you can focus on the solving the task via episodic experience, and you don't need to know the whole model of the environment. You use policy π to tackle the task for an entire episode:

$$S_0, A_0, R_1, S_1, A_1, \dots, S_{T-1}, A_{T-1}, R_T.$$

At the end of episode, you compute the return for every state visited:

$$v_\pi(s) = E_\pi[G_t | S_t = s], \text{ and } q_\pi(s, a) = E[G_t | S_t = s, A_t = a].$$

where $G_t = \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1}$.

Therefore, the estimated value for a state s or state action (s, a) is the average of all the returns that the agent has collected in that state or action of the state:

$$V_\pi(s) = \frac{1}{N} \sum_{k=1}^N G_{s_k}, \text{ and } Q_\pi = \frac{1}{N} \sum_{k=1}^N G_{s,a_k}.$$

Calculating Values

For a generated trajectory $(S_0, A_0, R_1, S_1, A_1, \dots, S_{T-1}, A_{T-1}, R_T)$, we calculate the returns for each moment t :

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T.$$

We can't use V anymore because we don't have the model $\pi(s', a' | s, a)$ and calculating v requires the evaluation of **each individual action** in a state, so we can only learn for Q which already implicitly learn the model:

$$q_\pi(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma v(s')].$$

where as if you want to choose an action from V , you still need to pick the action a with the highest value, and to do that, you need know the next state for every action a leads to and you need to know model p to get probability distributions for the next state s' from (s, a) .

Importance of Exploration

Since $Q(s, a)$ is an estimate that is improved by the agent collecting experience following unoptimal policies. Therefore, the estimates may not be accurate especially in the beginning of the learning and there is a chance that the a bad estimate prevents the agent from ever choosing (s, a) that might become optimal in the future. To prevent this, we make sure that all actions are chosen from time to time using:

1. Exploring starts with random state $S_0 \sim S$ and random action $A_0 \sim A(S_0)$ and,

2. Stochastic policies: $\pi(a|s) > 0, \forall a \in A(s)$.

Stochastic Policies

You can either generate the experience with the same policy you're trying to optimize (On-Policy), or generate experience with an exploratory policy b different from the one we're going to optimize.

ϵ -greedy Policy

You select a random action at probability ϵ , and select highest $Q(s, a)$ at probability $1 - \epsilon$:

$$\pi(a|s) = \begin{cases} 1 - \epsilon + \epsilon_r & \text{for } a = a^* \\ \epsilon_r & \text{for } a \neq a^* \end{cases}$$

where $\epsilon_r = \frac{\epsilon}{|A|}$.