

Monte Carlo Methods

- estimate value functions
- discover optimal policies

Monte Carlo : experience sample episodes

sample sequences of
states, actions, and rewards
from actual or simulated interaction
with an environment

- actual experience
- simulated experience

"Average sample returns"

Experiences \rightarrow episodes

- DP: compute value functions
- MC: learn value functions

"Simply average the returns observed after

visits to that state"



"As more returns are observed, the average should converge to the expected value."

A set of episodes



visit

state s

first-visit MC method

estimate $V_{\pi}(s)$ as the average of the returns

every-visit MC method

average all the returns following first visit to s

$S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$

$$G \leftarrow 0$$

$$t = T-1, T-2, \dots, 0$$

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

Unless S_t in S_0, S_1, \dots, S_{t-1}

$$\begin{cases} \text{Append } G \text{ to } \text{Return}(S_t) \\ V(S_t) \leftarrow \text{average}(\text{Returns}(S_t)) \end{cases}$$

$$n \uparrow \rightarrow \text{error} \downarrow \left(\frac{1}{\sqrt{n}} \right)$$

MC: estimates for each state are independent.

DO NOT bootstrap

① Learn from actual experience

② Learn from simulated experience

③ independent.

Maintaining exploration (Sufficient exploration)

exploring starts: episodes begin with state-action pairs randomly selected to cover all possibilities

action-value function Q

$$\pi(s) \doteq \arg \max_a Q(s, a)$$

$$Q_{\pi_k}(s, \pi_{k+1}(s)) = Q_{\pi_k}(s, \arg \max_a Q_{\pi_k}(s, a))$$

every pair has $= \max_a Q_{\pi_k}(s, a)$
a non-zero probability
of being selected as the start $\geq Q_{\pi_k}(s, \pi_k(s))$
 $\geq V_{\pi_k}(s)$

MC Control with Exploring Start:

episode $t = T-1, T-2, \dots, 0$

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

if S_t, A_t in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to Returns (S_t, A_t)

$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$

$\pi(S_{t+1}) \leftarrow \underset{a}{\operatorname{argmax}} Q(S_t, a)$

MC Control without Exploring Starts

- On-policy: evaluate or improve the policy that is used to make decisions
- off-policy: evaluate or improve the policy different from that used to generate data (may be unrelated to the policy followed)

On-policy Control:

policy "soft" $\pi(a|s) > 0$

ϵ -greedy policies ϵ probability
select an action at random

$\left\{ \begin{array}{ll} \frac{\epsilon}{|A(s)|} & \text{probability non-greedy action} \\ 1 - \epsilon + \frac{\epsilon}{|A(s)|} & \text{greedy action} \end{array} \right.$

Episode: $t = T-1, T-2, \dots, 0$

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

if S_t, A_t in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$

Append G to $\text{Returns}(S_t, A_t)$

$$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$$

$$A^* \leftarrow \underset{a}{\operatorname{argmax}} Q(S_t, a)$$

For all $a \in A(S_t)$:

$$\pi(a | S_t) \leftarrow \begin{cases} 1 - \epsilon + \frac{\epsilon}{|A(S_t)|} & \text{if } a = A^* \\ \frac{\epsilon}{|A(S_t)|} & \text{if } a \neq A^* \end{cases}$$

Problem: sufficient exploration
(maintaining exploration)

with
exploring start

without
exploring start

{ on-policy $\pi(a|s) > 0$ \rightarrow deterministic optimal policy
off-policy

Monte Carlo Control without Exploring Starts

$$Q_{\pi}(s, \pi'(s)) = \sum_a \pi'(a|s) Q_{\pi}(s, a)$$

$$= \frac{\epsilon}{|A(s)|} \sum_a Q_{\pi}(s, a) + (1-\epsilon) \max_a Q_{\pi}(s, a)$$

$$\geq \frac{\epsilon}{|A(s)|} \sum_a Q_{\pi}(s, a) + (1-\epsilon) \sum_a \frac{\pi(a|s) - \frac{\epsilon}{|A(s)|}}{1-\epsilon} Q_{\pi}(s, a)$$

$$= \frac{\epsilon}{|A(s)|} \sum_a Q_{\pi}(s, a) - \frac{\epsilon}{|A(s)|} \sum_a Q_{\pi}(s, a) + \sum_a \pi(a|s) Q_{\pi}(s, a)$$

$$= V_{\pi}(s)$$

$$\epsilon\text{-soft} : \pi(a|s) \geq \frac{\epsilon}{|A(s)|}$$

On-policy method: use ϵ -greedy policy

Select an action at random ϵ probability

$$\begin{cases} \frac{\epsilon}{|A(s)|} & \text{for non-greedy actions} \\ 1 - \epsilon + \frac{\epsilon}{|A(s)|} & \text{for greedy actions} \end{cases}$$

Off-policy Prediction

On-policy method: Learn action values not for the optimal policy, but for a near-optimal policy that still explores

{ target policy: learned policy, \rightarrow optimal policy
behavior policy: generate policy \rightarrow more exploratory

Estimate V_{π} , Q_{π} , \textcircled{b} ($b \neq \pi$)
 \uparrow
 $\textcircled{\pi}$ behavior policy
target policy

Convergence $\pi(a|s) > 0$, $b(a|s) > 0$

Importance Sampling (ratio)

weight returns according to the relative probability of their trajectories occurring under the target

and behavior policies.

State-action trajectory:

$$A_t, S_{t+1}, A_{t+1}, \dots, S_T$$

$$P_r \{A_t, S_{t+1}, A_{t+1}, \dots, S_T | S_t, A_{t:T-1} \sim \pi\}$$

$$= \pi(A_t | S_t) p(S_{t+1} | S_t, A_t) \pi(A_{t+1} | S_{t+1}) \dots$$

$$p(S_T | S_{T-1}, A_{T-1})$$

$$= \prod_{k=t}^{T-1} \pi(A_k | S_k) p(S_{k+1} | S_k, A_k)$$

$$P_{t:T-1} = \frac{\prod_{k=t}^{T-1} \pi(A_k | S_k) p(S_{k+1} | S_k, A_k)}{\prod_{k=t}^{T-1} b(A_k | S_k) p(S_{k+1} | S_k, A_k)}$$

$$= \prod_{k=t}^{T-1} \frac{\pi(A_k | S_k)}{b(A_k | S_k)}$$

$$\mathbb{E}[P_{t:T-1} G_t | S_t = s] = V_\pi(s)$$

Scale the return by the ratios and average the returns:

$$V(s) = \frac{\sum_{t \in \mathcal{J}(s)} \rho_{t:T(t)-1} G_t}{|\mathcal{J}(s)|}$$

Ordinary importance
sampling

$$N(s) = \frac{\sum_{t \in \mathcal{J}(s)} \rho_{t:T(t)-1} G_t}{\sum_{t \in \mathcal{J}(s)} \rho_{t:T(t)-1}}$$

Weighted average

(preferred)

Variance of the importance sampling-scaled returns:

$$\text{Var}[x] = E[(x - \bar{x})^2]$$

$$= E[x^2 - 2x\bar{x} + \bar{x}^2]$$

$$= E[x^2] - \bar{x}^2$$

Expected square of the importance-sampling-scaled return:

$$\mathbb{E}_b \left[\left(\frac{1}{T} \sum_{t=0}^{T-1} \frac{\pi(A_t | S_t)}{b(A_t | S_t)} G_t \right)^2 \right]$$