

Monte Carlo Methods

- we have no knowledge of the environment.
- Monte Carlo methods require only experience.
- we solve problem based on averaging sample returns.
- Only for episodic tasks. only on completion of an episode values are updated.

[5.1] Monte Carlo Prediction:

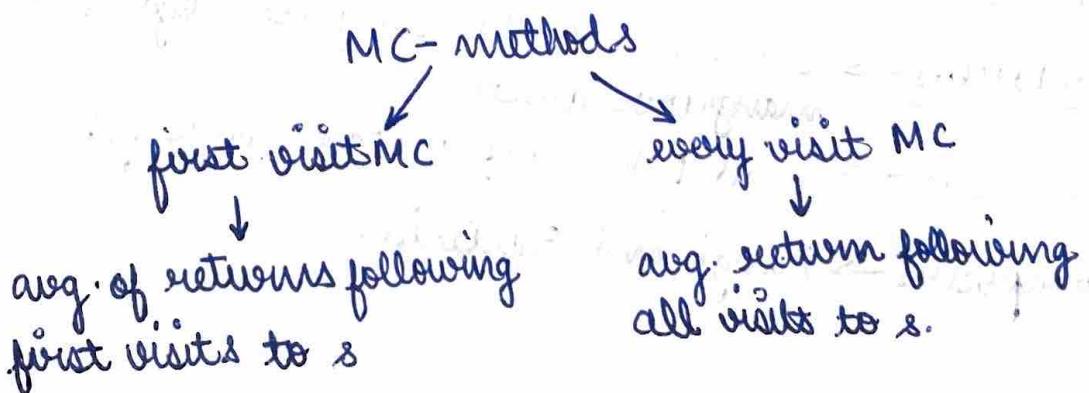
Learning state-value functions for a given policy.

we estimate $v_{\pi}(s)$ i.e. value of state s under policy π .

$$v_{\pi}(s) = E[G_t | s_t = s]$$

$$v_{\pi}(s) = E[R_{t+1} + \gamma v_{\pi}(s') | s_t = s]$$

Each occurrence of s in an episode is called a visit to state s .



First - visit MC prediction, for estimating $V \approx V_\pi$.

Input: a policy π to be evaluated

Initialise:

$V(s) \in \mathbb{R}$, arbitrarily $\forall s \in S$

Returns (s) \leftarrow an empty list, $\forall s \in S$.

Loop forever (for each episode):

generate an episode following $\pi: s_0, a_0, r_1, s_1, a_1, r_2, \dots$

$G_t \leftarrow 0$

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

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

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

Append G_t to Returns (s_t)

$V(s_t) \leftarrow \text{avg.}(\text{Returns}(s_t))$.

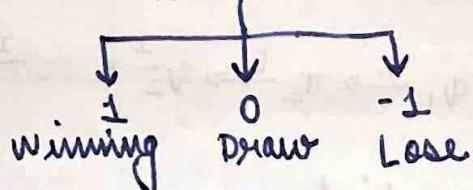
Example 5.1 \rightarrow Blackjack:

Object of the popular casino game is to obtain cards with sum close to 21 without exceeding it.

Face Cards $\rightarrow 10$, Ace $\rightarrow 1$ or 11 ,

Sticks on any sum greater than 17 or equal to it.

Rewards



$T=1$

[5.2] Monte Carlo estimation of action values:

If a model is not available, particularly useful to estimate action values, rather than state values.

similar to state value estimation but now we don't just look at state s but we look at (s, a) pairs.

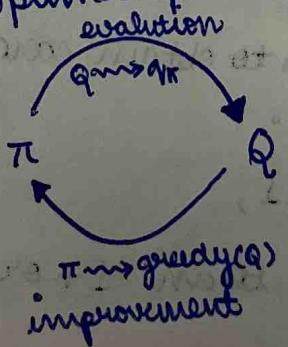
A visit to (s, a) means we are in state s and action a was taken in it.

The only complication is many state action pairs will never be visited.

To deal with this we mention in the start of the episodes state-action pair, and that every pair has a non-zero probability of being picked as the start.

[5.3] Monte Carlo Control:

We approximate Optimal policies



$$\pi_0 \xrightarrow{E} q_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} q_{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} q_{\pi_2} \xrightarrow{I} \dots \xrightarrow{I} \pi_* \xrightarrow{E} q_*$$

Policy evaluation is already studied.

$$\pi(s) \doteq \underset{a}{\operatorname{argmax}} \ q(s, a)$$

Policy improvement can be done by making each π_{k+1} greedy w.r.t q_{π_k} .

$$\begin{aligned} q_{\pi_k}(s, \pi_{k+1}(s)) &= q_{\pi_k}(s, \underset{a}{\operatorname{argmax}} q_{\pi_k}(s, a)) \\ &= \max_a q_{\pi_k}(s, a) \\ &\geq q_{\pi_k}(s, \pi_k(s)) \\ &\geq v_{\pi_k}(s) \end{aligned}$$

We make 2 assumptions

- The episodes have an exploring start
- Policy evaluation can be done with infinite number of episodes.

Monte Carlo ES, for estimating $\pi \approx \pi^*$

Initialise:

$$\pi(s) \in A(s) \text{ (arbitrarily)} \quad \forall s \in S$$

$$Q(s, a) \in R \text{ (arbitrarily)} \quad \forall s \in S, a \in A(s)$$

$$\text{Returns } (s, a) \leftarrow \text{empty list} \quad \forall s \in S, a \in A(s)$$

Loop forever (for each episode):

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

generate an episode from s_0, a_0 following $\pi: s_0, a_0, r_1, \dots, s_{T-1}, a_{T-1}, r_T$

$$G_t \leftarrow 0$$

unless the pair

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

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

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

Append G_t to Returns (s_t, a_t)

$$Q(s_t, a_t) \leftarrow \text{avg}(\text{Returns}(s_t, a_t))$$

$$\pi(s_t) \leftarrow \underset{a}{\operatorname{argmax}} Q(s_t, a)$$

[S.4] Monte Carlo Control without exploring starts:

Two approaches to solve this

on-policy methods

off-policy methods

On policy methods attempt to evaluate or improve a policy that is used to make decisions, whereas off-policy methods evaluate or improve the policy different from that used to generate data.

Monte Carlo ES in 5.3 is example of on-policy.

In on policy, we generally have a soft policy i.e. $\pi(a|s) > 0 \forall s \in S, a \in A(s)$, but gradually shifted closer to a deterministic policy.

We use ϵ -greedy i.e. with most of the time we choose maximal estimated action value but with probability ϵ , we instead select an action at random.

All non greedy actions are given minimal policy of selection

$\frac{\epsilon}{|A(s)|}$, and the remaining bulk $\frac{1-\epsilon + \frac{\epsilon}{|A(s)|}}{|A(s)|}$

ϵ -greedy policies are example of ϵ -soft policies for which $\pi(a|s) \geq \frac{\epsilon}{|A(s)|} \forall s \text{ and actions for some } \epsilon > 0$.

On-policy first-list MC control (for ϵ -soft) estimates $\pi \approx \pi^*$:

algorithm parameters: small $\epsilon > 0$

Initialise:

$\pi \leftarrow$ an arbitrary ϵ -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily) $\forall s \in S, a \in A(s)$

Returns(s, a) \leftarrow empty list, $\forall s \in S, \forall a \in A(s)$

Repeat forever (for each episode):

generate an episode following $\pi: s_0, a_0, r_1, s_1, a_1, r_2, \dots$

$G_t \leftarrow 0$

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

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

Unless the pair s_t, a_t appears in $s_0, a_0, s_1, a_1, \dots$

Append G_t to 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 $a \in A(s)$:

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

Any ϵ -greedy policy w.r.t q_π is an improver over any ϵ -soft policy π is assured by policy improvements.

Let π' be an ϵ -greedy policy:

$$q_\pi(s, \pi'(s)) = \sum_a \pi'(a|s) q_\pi(s|a)$$

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

$$\geq \frac{\varepsilon}{|A(s)|} \sum_a q_{\pi}(a|s) + (1-\varepsilon) \sum_a \frac{\pi(a|s) - \frac{\varepsilon}{|A(s)|} q_{\pi}(s|a)}{1-\varepsilon}$$

$$= \frac{\varepsilon}{|A(s)|} \sum_a q_{\pi}(s|a) - \frac{\varepsilon}{|A(s)|} \sum_a q_{\pi}(s|a) + \sum_a \pi(a|s) q_{\pi}(s|a)$$

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

If we consider a new environment that overrides the agent's choice of action with probability ε with a randomly picked action.

Let \tilde{v}_* and \tilde{q}_* denote new environments optimal value functions.

Then a policy π is optimal among policies iff $v_{\pi} = \tilde{v}_*$

$$\tilde{v}_*(s) = (1-\varepsilon) \max_a \tilde{q}_*(s|a) + \frac{\varepsilon}{|A(s)|} \sum_a \tilde{q}_*(s|a)$$

$$= (1-\varepsilon) \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma \tilde{v}_*(s')]$$

$$+ \frac{\varepsilon}{|A(s)|} \sum_a \sum_{s', r} p(s', r | s, a) [r + \gamma \tilde{v}_*(s')]$$

$$v_{\pi}(s) = (1-\varepsilon) \max_a q_{\pi}(s|a) \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')]$$

$$+ \frac{\varepsilon}{|A(s)|} \sum_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')]$$

If we compare both equations we get

$$v_{\pi} = \tilde{v}_{*}$$

[5.5] Off-Policy Prediction via Importance Sampling:

We use 2 policies \rightarrow 1. that is learned about and becomes optimal

2. Exploratory policy used to generate behaviour.

The policy being learned about is called target policy, policy used to generate behaviour, is called behaviour policy.

off-policy methods \rightarrow greater variance

off-policy methods \rightarrow slower convergence

\hookrightarrow more powerful & general

We take a prediction problem, both target and behaviour policies are fixed. We need to estimate v_{π} or q_{π} , we have episodes following another policy b , where $b \neq \pi$, in this case $b \rightarrow$ behaviour policy, $\pi \rightarrow$ target policy.

We require that every action taken under π is also at least occasionally taken under b . i.e $\pi(a|s) > 0 \Rightarrow b(a|s) > 0$. This is called assumption of coverage.

b must be stochastic in states where it is not identical to π , but π may be deterministic.

In control \rightarrow target policy \rightarrow deterministically greed w.r.t action value functions.

\hookrightarrow behaviour policy \rightarrow stochastic and exploratory.

Almost all off-policy methods use importance sampling

a general technique
for estimating expected values
under one distribution given
samples from other.

We weight returns acc. to the relative probability of their
trajectories under the target and behaviour policies, called
importance sampling ratio.

Given starting state $s_t \rightarrow A_t, s_{t+1}, A_{t+1}, s_{t+2}, \dots, s_T$

$$\Pr \{ A_t, s_{t+1}, A_{t+1}, \dots, s_T \}$$

$$= \pi(A_t | s_t) p(s_{t+1} | s_t, A_t) \pi(A_{t+1} | s_{t+1}) \dots p(s_T | s_{T-1}, A_{T-1})$$

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

$k=t$

Thus the relative probability under target & behaviour
is

$$\frac{\prod_{k=t}^{T-1} \pi(A_k | s_k) p(s_{k+1} | s_k, A_k)}{\prod_{k=t}^{T-1} b(A_k | s_k) p(s_{k+1} | s_k, A_k)}$$

$k=t$

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

$$\frac{\prod_{k=t}^{T-1} \pi(A_k | s_k)}{\prod_{k=t}^{T-1} b'(A_k | s_k)}$$

We have $G_{1:t}$ from the π policy

$$E[G_{1:t} | S_t = s] = v_\pi(s)$$

This cannot be averaged in terms of $v_\pi(s)$

But

$$E[\sum_{t=1:T-1} G_{1:t} | S_t = s] = v_\pi(s)$$

We define set of all timestamps in which state s is visited, denoted $\mathcal{I}(s)$. This is for an every-visit method, for a first-visit method $\mathcal{I}(s)$ would only include that were the first time stamps to visit s within their episodes. Let $T(t)$ denote first time of termination. Then $\{G_{1:t}\}_{t \in \mathcal{I}(s)}$ are the returns that pertain to state s and mean $\{\sum_{t=1:T(t)-1} G_{1:t}\}_{t \in \mathcal{I}(s)}$

To estimate $v_\pi(s)$

$$V(s) \doteq \frac{\sum_{t \in \mathcal{I}(s)} \sum_{t=1:T(t)-1} G_{1:t}}{|\mathcal{I}(s)|}$$

↓
Ordinary importance sampling

$$V(s) \doteq \frac{\sum_{t \in \mathcal{I}(s)} \sum_{t=1:T(t)-1} G_{1:t}}{\sum_{t \in \mathcal{I}(s)} \sum_{t=1:T(t)-1} 1}$$

↓
Weighted importance sampling.

suppose we obtain $K(s)$ returns for class s .

$$V_{\pi}(s) \approx \frac{1}{K(s)} \sum_{k=1}^{K(s)} \frac{\rho_{t(k)} : T(k)-1}{G_t^{(k)}}$$

↓
Ordinary

$$V_{\pi}(s) \approx \frac{\sum_{k=1}^{K(s)} \rho_{t(k)} : T(k)-1}{\sum_{k=1}^{K(s)} G_t^{(k)}}$$

↓
Weighted

[5.6] Incremental Implementation:

In ordinary importance sampling, returns are scaled by the importance sampling ratio $\rho_{t:T(t)-1}$, then simply averaged; we can use incremental updates of averaging but use scaled returns instead.

For weighted importance sampling:

suppose we have G_1, G_2, \dots, G_{m-1} all starting from the same state with a corresponding random weight w .

$$V_m = \frac{\sum_{k=1}^{m-1} w_k G_k}{\sum_{k=1}^{m-1} w_k} \quad m \geq 2$$

To update

$$V_{m+1} = V_m + \frac{w_m}{G_m} [G_m - V_m], \quad m \geq 1$$

$$G_{m+1} = G_m + w_{m+1}$$