

Monte Carlo Methods





Rupam Mahmood

February 10, 2020



Estimating expectation by sample averaging

- Sample average is an unbiased estimator of expectation \checkmark Lemma 1 (<u>Mahmood 2017</u>)
- Sample average is a consistent estimator of expectation \checkmark Lemma 3 (<u>Mahmood 2017</u>)



Unbiased and consistent estimation

The sample average $Z_n =$

Then we have $E[Z_n]$

And we have
$$P\left(\lim_{i\to\infty} Z_n = E[X_i]\right) = 1$$
; consistency of Z_n

Say X_i is an iid random variable

$$\frac{\sum_{i=}^{n} X_{i}}{n}$$
 is an estimate of $E[X_{i}]$

=
$$E[X_i]$$
; unbiasedness of Z_n

Monte Carlo prediction

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

Input: a policy π to be evaluated Initialize: $V(s) \in \mathbb{R}$, arbitrarily, for all $s \in S$ $Returns(s) \leftarrow an empty list, for all <math>s \in S$ Loop forever (for each episode): Generate an episode following $\pi: S_0, A_0, R_1, S_1, A_1, R_2, ..., S_{T-1}, A_{T-1}, R_T$ $G \leftarrow 0$ Loop for each step of episode, t = T - 1, T - 2, ..., 0: $G \leftarrow \gamma G + R_{t+1}$ Unless S_t appears in $S_0, S_1, \ldots, S_{t-1}$: Append G to $Returns(S_t)$ $V(S_t) \leftarrow \operatorname{average}(Returns(S_t))$

Friday's Monte Carlo code: Is it every-visit or first-visit?

```
1 import numpy as np
 2
 3 \text{ neps} = 1000; \text{ ns} = 1; \text{ g} = 0.9; \text{ probl} = 0.1
 4
 6
 8
 9 \text{ rets} = []
10 for ep in range(neps):
      s = 0; ret = 0; t = 0
11
12
       while s != 1:
13
           rnd = np.random.random()
           if rnd < pol[s, 0]: a = 0
14
15
          else: a = 1
16
17
           sp = s + a
18
           r = rewards[s, a]
19
            ret += r *(g**t)
20
            s = sp; t = t + 1
21
       rets.append(ret)
22 print("Average return:", np.mean(rets))
```

5 pol = np.zeros((ns, 2)); pol[0, 0] = probl; pol[0, 1] = 1 - probl

7 rewards = np.zeros((ns, 2)); rewards[0, 0] = -1; rewards[0, 1] = +1

Friday's Monte Carlo code: what's the difference between this and MC prediction algorithm from the book?

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

Input: a policy π to be evaluated Initialize: $V(s) \in \mathbb{R}$, arbitrarily, for all $s \in S$ $Returns(s) \leftarrow$ an empty list, for all $s \in S$ Loop forever (for each episode): Generate an episode following π : $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$ $G \leftarrow 0$ Loop for each step of episode, $t = T - 1, T - 2, \dots, 0$: $G \leftarrow \gamma G + R_{t+1}$ Unless S_t appears in S_0, S_1, \dots, S_{t-1} : Append G to $Returns(S_t)$ $V(S_t) \leftarrow$ average($Returns(S_t)$)

```
1 import numpy as np
 2
 3 \text{ neps} = 1000; \text{ ns} = 1; \text{ g} = 0.9; \text{ probl} = 0.1
 5 pol = np.zeros((ns, 2)); pol[0, 0] = probl; pol[0, 1] = 1 - probl
 6
 7 rewards = np.zeros((ns, 2)); rewards[0, 0] = -1; rewards[0, 1] = +1
 8
 9 rets = []
10 for ep in range(neps):
       s = 0; ret = 0; t = 0
11
12
       while s != 1:
           rnd = np.random.random()
13
           if rnd < pol[s, 0]: a = 0
14
           else: a = 1
15
16
17
            sp = s + a
18
           r = rewards[s, a]
19
           ret += r *(g**t)
           s = sp; t = t + 1
20
       rets.append(ret)
21
22 print("Average return:", np.mean(rets))
```

Monte Carlo version of classical policy iteration (with construction of greedy policies)

$$\pi_0 \xrightarrow{\mathrm{E}} q_{\pi_0} \xrightarrow{\mathrm{I}} \pi_1 \xrightarrow{\mathrm{E}} q_{\pi_1}$$

Could we use state value estimates?

Why is exploring starts necessary?

Why is exploring starts impractical?

What is the other impractical assumption here?



Monte Carlo control with generalized policy iteration

Monte Carlo ES (Exploring Starts), for estimating $\pi \approx \pi_*$

Initialize:

 $\pi(s) \in \mathcal{A}(s)$ (arbitrarily), for all $s \in S$ $Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in S$, $a \in \mathcal{A}(s)$ $Returns(s, a) \leftarrow \text{empty list, for all } s \in S, a \in \mathcal{A}(s)$

Loop forever (for each episode): Choose $S_0 \in S$, $A_0 \in \mathcal{A}(S_0)$ randomly such that all pairs have probability > 0 Generate an episode from S_0, A_0 , following $\pi: S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T$ $G \leftarrow 0$ Loop for each step of episode, $t = T - 1, T - 2, \ldots, 0$: $G \leftarrow \gamma G + R_{t+1}$ Unless the pair S_t, A_t appears in S_0, A_0 , Append G to $Returns(S_t, A_t)$ $Q(S_t, A_t) \leftarrow \operatorname{average}(Returns(S_t, A_t))$ $\pi(S_t) \leftarrow \operatorname{arg\,max}_a Q(S_t, a)$

$$S_1, A_1 \dots, S_{t-1}, A_{t-1}$$
:



Monte Carlo control without exploring start

On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$ Initialize: $\pi \leftarrow$ an arbitrary ε -soft policy $Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in S, a \in A$ $Returns(s, a) \leftarrow \text{empty list, for all } s \in S, a \in$ Repeat forever (for each episode): Generate an episode following π : S_0, A_0, R_1 , $G \leftarrow 0$ Loop for each step of episode, $t = T - 1, T - 2, \ldots, 0$: $G \leftarrow \gamma G + R_{t+1}$ Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \ldots, S_{t-1}, A_{t-1}$: Append G to $Returns(S_t, A_t)$ $Q(S_t, A_t) \leftarrow \operatorname{average}(Returns(S_t, A_t))$ $A^* \leftarrow \operatorname{arg\,max}_a Q(S_t, a)$ For all $a \in \mathcal{A}(S_t)$: $\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon / |\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon / |\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$

$$(s) \in \mathcal{A}(s)$$

$$\ldots, S_{T-1}, A_{T-1}, R_T$$

(with ties broken arbitrarily)

