

CMPUT 365: Monte Carlo Methods

Rupam Mahmood

Feb 7, 2022

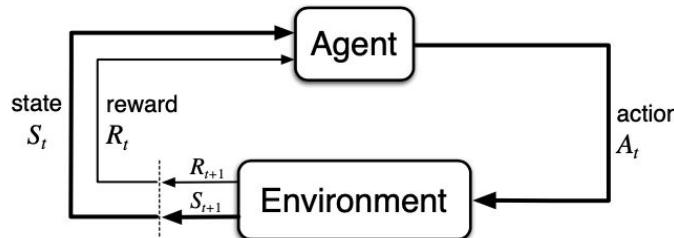
Admin: Coursera course 2 starts from this week

Due dates C2M1

Midterm:

- Exam time: 45 minutes + 5 minutes for uploading
- Syllabus: Coursera course 1 / Chapter 2-4 / up to Dynamic Programming
- The practice midterm is for you to understand the kind of questions to expect
- The answers should be as detailed as the sample answers
- Practicing and making a cheat sheet are way more helpful than hoping that you can figure out the answers by opening the book

Monte Carlo is a sample-based learning method



Benefits:



Vova Selin 对所有人说

下午 1:12

Since we usually don't know the environment dynamics accurately



David Wang 对所有人说

下午 1:14

It saves times. You don't have to update the entire model on every iteration.

...

- Learning from actual experience is striking because it requires no prior knowledge of the environment's dynamics, yet can still attain optimal behavior
- We saw a [colab](#) example before of averaging returns
- Monte Carlo (MC) methods are more easily understood in the episodic setting
- Updates can be made only after an episode is done
- Updates can be made for the first visit of a state or every visit of a state

MC prediction

First-visit MC prediction, for estimating $V \approx v_{\pi} \rightarrow$ Target policy, we want to estimate the value of policy π .

Input: a policy π to be evaluated

Initialize:

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

$Returns(s) \leftarrow$ an empty list, for all $s \in \mathcal{S}$

Loop forever (for each episode):

Generate an episode following $\pi \rightarrow S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

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

$G \leftarrow \gamma G + R_{t+1} \rightarrow$ keep calculating the return. each step's return are unbiased estimates.

Unless S_t appears in S_0, S_1, \dots, S_{t-1} : \rightarrow a sequence (S_t is not here)

Append G to $Returns(S_t)$

$\leftarrow V(S_t) \leftarrow \text{average}(Returns(S_t))$
↳ $V(S_t)$ is unbiased (指不带偏的无偏的). Bias: 带偏的. BT: estimate value 估计值. True value 真值

↳ average of all returns G from state S_0 . \rightarrow $E[X] = y$, y is an unbiased estimate of y .

$\therefore E[V_t(S_t)] = v_{\pi}(S_t)$

$\therefore E[G | S_t] = v_{\pi}(S_t)$

- After each update here, the average V remains an unbiased estimate of the value function
- How do you convert this into an every-visit one?

MC estimation of Action Values

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 $s \in \mathcal{S}$

$Returns(s) \leftarrow$ an empty list, for all $s \in \mathcal{S}$

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 $Returns(S_t)$

$V(S_t) \leftarrow \text{average}(Returns(S_t))$

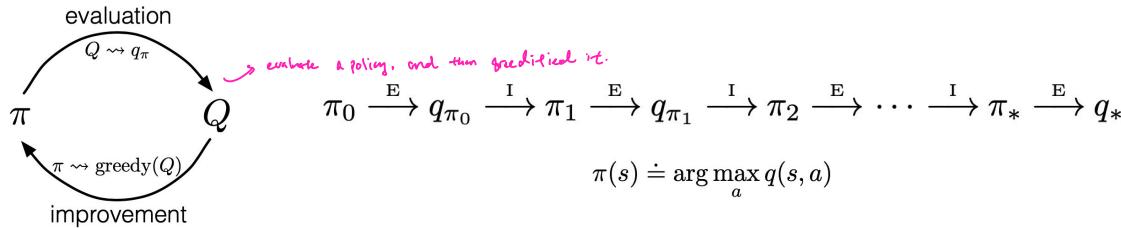
- How can we convert this pseudocode for action values?
- Sample-based methods fundamentally rely on exploration and randomness
- Two ways to maintain exploration: exploring start and using stochastic policies

random start state: ensure every state will be chosen as a start state equally (if have a start state randomly).

→ this random uniform rating.

15%	20%
25%	25%

MC control



- MC control must use action-value estimates. Why?
- Let's consider a classical policy iteration version of MC (no model but full evaluation and improvement in turn)
- Policy evaluation is done as before, for example, using exploring start
- Policy improvement is done by making the policy greedy with respect to the current value
- According to the policy improvement theorem, the new policy usually becomes larger or better than the old policy? *⇒ Better, not larger.*
- Why this method is impractical?

MC control

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

Initialize:

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

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

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Loop forever (for each episode):

Choose $S_0 \in \mathcal{S}, A_0 \in \mathcal{A}(S_0)$ randomly *explore start* such that all pairs have probability > 0

Generate an episode from S_0, A_0 , following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

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

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

Unless 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 $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$ *gradually the value estimated.*

$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a)$

- We can give up on completing the evaluation step to avoid sampling infinite number of episodes at each evaluation step
- Stability is achieved only when both the policy and the value function are optimal
- But convergence is yet to be proved

MC control without exploring starts

- We can get rid of exploring start by using a stochastic policy
- We move the policy only to an epsilon-greedy policy. Still works as a GPI
- What's the difference between epsilon-soft and epsilon-greedy policy?
- any epsilon-greedy policy with respect to $q \pi$ is an improvement over any epsilon-soft policy π

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

Initialize:

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

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

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Loop forever (for each episode):

Choose $S_0 \in \mathcal{S}$, $A_0 \in \mathcal{A}(S_0)$ randomly such that all pairs have probability > 0

Generate an episode from S_0, A_0 , following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

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

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

Unless 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 $Returns(S_t, A_t)$

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

$\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$

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 $s \in \mathcal{S}, a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

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

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

Unless 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 $Returns(S_t, A_t)$

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

$A^* \leftarrow \text{argmax}_a Q(S_t, a)$ (with ties broken arbitrarily)

For all $a \in \mathcal{A}(S_t)$:

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

better than ε -greedy
 ε -greedy \Rightarrow ε -greedy
 π is ε -greedy

greedy \Rightarrow ε -greedy

greedy \Rightarrow ε -greedy

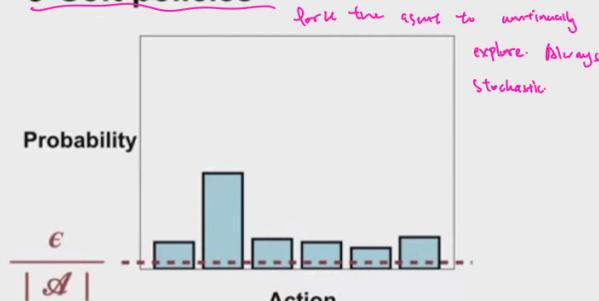
ϵ -Greedy Exploration

⇒ how can we learn all action values without exploring starts?

ϵ -Greedy policies



ϵ -Soft policies



~~Exploring Starts~~

Epsilon soft policies are always stochastic. Deterministic policy specify a single action to take in each state, stochastic policies instead specify the probability of taking action in each state in epsilon. Soft policies: all actions have a probability of at least Epsilon over the number of actions. They will eventually try all the actions.

History of RL

- Rich S. Sutton
- Andy Barto
- A. Harry Klopf
- Philosophy of science: empiricism, logical positivism, Karl Popper
- Phenomenology: Hubert Dreyfus

Last slido question

If return G is obtained from state s and following policy , is G an unbiased estimate of $v_{\pi}(s)$?

0 1 2

Yes

8%

No

92%

Participants can vote at slido.com with code #X978

Same as asking $E_b[G_t | S_t = s] \stackrel{?}{=} v_{\pi}(s)$.

$$E_b[G_t | S_t = s]$$

$$E_{A_k \sim b}[G_t | S_t = s]$$

A random variable G_t doesn't tell much about the probability distribution. (we don't know how the action is chosen by the agent). \therefore G_t is not an unbiased estimate of $v_{\pi}(s)$.

One action is drawn from π in this distribution.

[Let's review unbiasedness]

$$Z_n = \frac{\sum_{i=1}^n x_i}{n}$$

Z_n = sample average estimator.
Based on the value of *n*

x_i → R.V.

x_i → R.V.

of $E_{x \sim p}[x]$ where $x_i \stackrel{iid}{\sim} p$?

slido question ...

Same as asking $E_{x_i \sim p}[z_n] \stackrel{?}{=} E_{x \sim p}[x]$

x_i → each of the action is drawn from policy *T_i*

x → each of the action is drawn from policy *p*.

[what if data comes from a different distribution?]

Is $Z_n = \frac{\sum_{i=1}^n x_i}{n}$ an unbiased estimate

of $E_{X \sim p}[X]$ where $x_i \sim d$?

slide question ...

Same as asking $E_{X_i \sim d}[Z_n] \stackrel{?}{=} E_{X \sim p}[X]$.

[Unbiasedness of sample average is connected to the unbiasedness of the summand]

$$Z_n = \frac{x_1 + x_2 + \dots + x_n}{n}$$

↑ summand
each of this summand are unbiased estimates.

say $\mu \doteq E_{X \sim p}[X]$.

Then $E_{X_i \sim p}[X_i] = \mu$ as well

because x, x_i are all i.i.d.

And $E_{X_i \sim p}[x_i] = \mu$

$\Rightarrow E_{X_i \sim p}[Z_n] = \mu$.

$$E_{X_i \sim p}[Z_n] = E_{X_i \sim p}\left[\frac{\sum_{i=1}^n X_i}{n}\right]$$

$$= \frac{1}{n} \sum_{i=1}^n E_{X_i \sim p}[X_i]$$

operation of expectation of each of this X_i .

$$= \frac{1}{n} \sum_{i=1}^n \mu$$

$$= \frac{1}{n} \times n \mu = \mu.$$

[Therefore $E_{X_i \sim d} [Z_n] \neq \mu$]

$$Z_n = \frac{\sum_{i=1}^n X_i}{n} ; \quad X_i \sim d$$

each of the summand are drawn from the different distribution d (not p).

$$E_{X_i \sim d} [X_i] \neq \mu = E_{X \sim p} [X]$$

mean of distribution $d \neq$ mean of distribution p .
 $\therefore E_b [G_t | S_t = s] \neq V_\pi(s)$ policy π .

Estimate return under policy b .

$b \neq \pi$. \therefore 值估计错误!!

$$V(s) = \frac{\sum_{t \in T(s)} G_t}{|T(s)|} ; G_t \sim b$$

$$E_b [G_t | S_t = s] \neq V_\pi(s)$$

$$E_b [V(s)] \neq V_\pi(s).$$

$$\mu = V_\pi(s)$$

How about $E_{X_i \sim d} [Z_n'] = ?$

$$Z_n' = \frac{\sum_{i=1}^n Y_i}{n}, \text{ where}$$

$$Y_i = \frac{p(X_i)}{d(X_i)} X_i \quad \begin{array}{l} \text{summand} \\ \text{importance sampling ratio} \end{array}$$

probability (density) of X_i under a distribution.

Think in terms of summands of Z_n' .

$$E_{X_i \sim d} [Y_i] = \mu \Rightarrow$$

$$E_{X_i \sim d} [Z_n'] = \mu.$$

$$\begin{aligned} E_{X_i \sim d} [Y_i] &= E_{X_i \sim d} \left[\frac{p(X_i)}{d(X_i)} X_i \right] \\ &= \sum \frac{p(x)}{d(x)} x \underbrace{d(x)}_{\substack{\text{probability under the base R.V.}}} \\ &= \sum p(x) x = E_{x \sim p} [x] = \mu \end{aligned}$$

assumption of coverage

$$\underbrace{p(x) > 0}_{\substack{\text{when our p.w. is} \\ \text{non-zero}}} \Rightarrow \underbrace{d(x) > 0}_{\substack{\text{then d.w. is non-zero.}}}$$

$$a_t \frac{p(a_t | \pi)}{p(a_t | \theta)} \xrightarrow{\text{wants to predict } \pi.}$$

drawn from π

Unbiased and consistent estimation

Say $X_i \sim p$ is an iid random variable

The sample average $Z_n = \frac{\sum_{i=1}^n X_i}{n}$ is an estimate of $E_{X \sim p}[X] = \sum_x xp(x)$

So is X_i

Then we have $E_{\underbrace{X_i \sim p}_{\text{Sample average}}}[Z_n] = E_{X \sim p}[X]$; unbiasedness of Z_n

And we have $P\left(\lim_{n \rightarrow \infty} Z_n = E_{X \sim p}[X]\right) = 1 \iff \underbrace{Z_n}_{\text{as } n \rightarrow \infty, \text{ this estimate is still random.}} \xrightarrow{a.s.} E_{X \sim p}[X]$; consistency of Z_n

going closer and closer to the thing we want.

On the other hand, we have $E_{X_i \sim p}[X_i] = E_{X \sim p}[X]$, but not $X_i \xrightarrow{a.s.} E_{X \sim p}[X]$

When samples are from a different distribution ...

Say $X_i \sim d$ is an iid random variable (note the difference in distribution)

Let's call d the data distribution, and p the target distribution

The sample average $Z_n = \frac{\sum_{i=1}^n X_i}{n}$ is a *bad* estimate of $E_{X \sim p}[X] = \sum_x xp(x)$

Because now we have $E_{X_i \sim d}[Z_n] = E_{X \sim d}[X] \neq E_{X \sim p}[X]$

And $Z_n \xrightarrow{a.s.} E_{X \sim d}[X] \neq E_{X \sim p}[X]$

When samples are from a different distribution ...

Obviously, $X_i \sim d$ is a worse estimate of $E_{X \sim p}[X]$

How about $Y_i = \frac{p(X_i)}{d(X_i)} X_i$, where $X_i \sim d$?

If d provides adequate coverage of p : $p(x) > 0$ implies $d(x) > 0$,

$$E_{X_i \sim d} [Y_i] = E_{X_i \sim d} \left[\frac{p(X_i)}{d(X_i)} X_i \right] = \sum_x \frac{p(x)}{d(x)} x d(x)$$

$$= \sum_x x p(x) = E_{X \sim p}[X]$$

When samples are from a different distribution, we can use importance sampling correction

$\frac{p(X_i)}{d(X_i)}$ is known as the **importance sampling ratio**

It can be used to correct the discrepancy between target and data distributions

The following importance sampling estimator is an unbiased and consistent

estimator of $E_{X \sim p}[X]$

$$Z'_n = \frac{\sum_{i=1}^n Y_i}{n}, \text{ where } Y_i = \frac{p(X_i)}{d(X_i)} X_i \text{ and } X_i \sim d$$

[Last slido question]

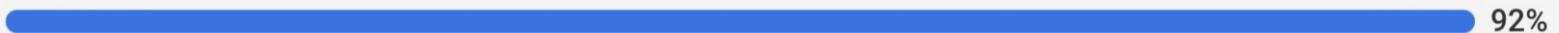
If return G is obtained from state s and following policy , is G an unbiased estimate of $v_{\pi}(s)$?

0 1 2

Yes

 8%

No

 92%

Participants can vote at slido.com with code #X978

Same as asking $E_b [q_t | S_t = s] \stackrel{?}{=} v_{\pi}(s)$.

$E_b[G_t | S_t = s] \neq v_{\pi}(s)$ because $E_b[G_t | S_t = s] = v_b(s)$
and $b \neq \pi$ in general.

Importance sampling for off-policy prediction

We want to estimate v_π whereas samples are from a different policy $b \neq \pi$

We call b the behavior policy, and π the target policy

Then the importance sampling ratio for a trajectory corresponding to return G_t is

$$\rho_{t:T-1} \doteq \frac{P(A_t, S_{t+1}, A_{t+1}, \dots, S_T | S_t, A_{t:T-1} \sim \pi)}{P(A_t, S_{t+1}, A_{t+1}, \dots, S_T | S_t, A_{t:T-1} \sim b)}$$

Annotations:

- A_t (red arrow)
- For action-state sequence (pink arrow)
- Action drawn from the probability of the same sequence (pink arrow)
- you want to estimate unknown for agent (green arrow)
- $\pi(A_k | S_k)$ (green circle)
- $b(A_k | S_k)$ (green circle)
- known for the agent (green arrow)

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

Importance sampling estimator for off-policy prediction: $V(s) \doteq \frac{\sum_{t \in \mathcal{T}(s)} \rho_{t:T(t)-1} G_t}{|\mathcal{T}(s)|}$

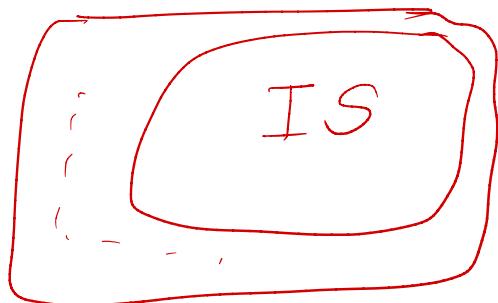
→ unbiased estimate.

Here $A_k \sim b$.

Hence, $E_b[V(s)] = v_{\pi}(s)$.

$$E_{A \sim b} \left[\rho_{t:T(t)-1} G_t \right] = v_{\pi}(s).$$

$V(s)$ is an importance sampling estimator and also an off-policy estimator.



Importance sampling for off-policy prediction

Sample average estimator for on-policy prediction: $\hat{V}(s) \doteq \underbrace{\frac{\sum_{t \in \mathcal{T}(s)} G_t}{|\mathcal{T}(s)|}}_{\text{sample average}}$

$\mathcal{T}(s)$ contains all time steps in which state s is visited

G_t denotes the return after t up through $T(t)$

$T(t)$ denotes the first time of termination after t

Importance sampling estimator for off-policy prediction: $\hat{V}(s) \doteq \frac{\sum_{t \in \mathcal{T}(s)} \rho_{t:T(t)-1} G_t}{|\mathcal{T}(s)|}$

HS

can a decision made based on the feedback loop sometime deviates away from your ultimate goal because here you are now focused on the immediate reward? \Rightarrow the algorithm is decided for focus on the 长期收益 (长期回报) (discount rate), 回报衰减 (discounting).

Worksheet 6: Monte Carlo

CMPUT 397
February 8, 2022

1. (Exercise 5.4 S&B) The pseudocode for *Monte Carlo ES* is inefficient because, for each state-action pair, it maintains a list of all returns and repeatedly calculates their mean. How can we modify the algorithm to have incremental updates for each state-action pair?

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

Initialize:

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

Loop forever (for each episode):

Choose $S_0 \in \mathcal{S}$, $A_0 \in \mathcal{A}(S_0)$ randomly such that all pairs have probability > 0
Generate an episode from S_0, A_0 , following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$
 $G \leftarrow 0$

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

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

Unless 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 $Returns(S_t, A_t)$
 $Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$
 $\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$