

1. (20)

For the average reward setting, the differential state-value function is defined as:

$$v_{\pi}(s) \stackrel{\text{def}}{=} E_{\pi} \left[\sum_{k=t+1}^{\infty} (R_k - r(\pi)) \mid S_t = s \right],$$

where the average reward is

$$r(\pi) = \lim_{h \rightarrow \infty} E_{\pi} \left[\frac{1}{h} \sum_{t=0}^{h-1} R_{t+1} \right].$$

Derive the following Bellman equation for v_{π} :

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

Show steps.

$$\begin{aligned}
V_{\pi}(s) &= \mathbb{E}_{\pi} \left[\sum_{k=t+1}^{\infty} (R_k - r(s_k)) \mid S_t = s \right] \\
&= \mathbb{E}_{\pi} \left[R_{t+1} - r(s_t) + \sum_{k=t+2}^{\infty} (R_{k+1} - r(s_k)) \mid S_t = s \right] \\
&= \mathbb{E}_{\pi} \left[R_{t+1} - r(s_t) + \mathbb{E} \left[\sum_{k=t+2}^{\infty} (R_{k+1} - r(s_k)) \mid S_{t+1} = s \right] \mid S_t = s \right] \Rightarrow \text{Law of Total Expectation: } \mathbb{E}[\mathbb{E}[X|Y, Z]|Z] \\
&= \mathbb{E}_{\pi} \left[R_{t+1} - r(s_t) + \mathbb{E} \left[\sum_{k=t+2}^{\infty} (R_{k+1} - r(s_k)) \mid S_{t+1} = s \right] \mid S_t = s \right] \Rightarrow \text{Markov Property: } \mathbb{E}[X|Y, Z] = \mathbb{E}[X|Z] \\
&= \mathbb{E}_{\pi} \left[R_{t+1} - r(s_t) + V_{\pi}(s_{t+1}) \mid S_t = s \right] \Rightarrow V_{\pi}(s) = \mathbb{E}_{\pi} \left[\sum_{k=t+1}^{\infty} (R_k - r(s_k)) \mid S_t = s \right] \\
&= \sum_{a \in A, s' \in S} \mathbb{P}_{\pi}(A_t=a, S_{t+1}=s') \cdot R_{t+1} - r(s_t) + V_{\pi}(s') \Rightarrow \text{Bellman Equation: } V_{\pi}(s) = \sum_{x \in X} g(x) \cdot \mathbb{P}(X=x) \\
&= \sum_a \mathbb{P}_{\pi}(a|s) \cdot \sum_{s'} \mathbb{P}(s'; r|s, a) \cdot [r - r(s)] + V_{\pi}(s')
\end{aligned}$$

6. 用 $g(x)$ 表示 bellman equation 裡的 $V_{\pi}(s)$ 進

(20)

2. Modify the Tabular TD(0) algorithm for estimating v_π , to estimate q_π .

Tabular TD(0) for estimating v_π

1. Input: the policy π to be evaluated
2. Algorithm parameter: step size $\alpha \in (0, 1]$
3. Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$
4. Loop for each episode:
5. Initialize S
6. Loop for each step of episode:
7. $A \leftarrow$ action given by π for $S \Rightarrow$ estimate $q_{\pi, A}$, \therefore the target action. (根据 π 选择行动)
8. Take action A , observe R, S'
9. $V(S) \leftarrow V(S) + \alpha [R + \gamma V(S') - V(S)]$
10. $S \leftarrow S'$
11. until S is terminal

基于上一节的 Sarsa learning TD control, 只是将 \leftarrow action given by π for S 换为 q_π 。
即该表格: (由 π 演变而来):

① Tabular TD(0): estimate v_π , 通过假设 π 不变, 为 prediction.

$$\text{update } V \text{ 为 } V \leftarrow V(S) + \alpha [R + \gamma V(S') - V(S)]$$

② Sarsa: on policy TD control, estimate $Q \rightarrow q^*$, 通过 π 选择 optimal action value $q_{\pi, A}$, \therefore 为 update policy.

$$\text{update } Q \text{ 为 } Q \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$$

③ Q-learning: off-policy TD control, estimate $\pi_0 \rightarrow \pi_{\text{opt}}$, 通过选择 optimal policy π^* , \therefore 为 update policy.
behavior policy 为选择 max value 的 action

target policy (要更新) 为 π : greedy.

$$\text{更新: } Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', A') - Q(S, A)]$$

④ Expected Sarsa: 通过 predict, control, update, off-policy.

$$\text{更新: } Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \sum_a \pi_a(S) Q(S', a) - Q(S, A)]$$

⑤ Model based Q-planning: 基于 Q-learning 的基础上加上基于 model learning, 即:
从 S, A 为 a sample model, 为 model learning, 为 obtain.

从 sample next reward R , 为 sample next state S'

⑥ Tabular Dyna Q: 结合 Q learning + Q planning.

- 前向 RL direct RL (直接与环境交互) environmental interaction (experience) \Rightarrow 更新 $Q(S, A)$, 通过方式 1
- Q learning - 方式 2.
- 后向用之前得到的 (S, A) 为 S , 为来更新 model, 更新公式为 Q-learning - 方式 3.

⑦ Tabular Dyna Q+: 结合 Q learning + Q planning, 大体都是 Dyna Q - 方式.

但直接 model 更新中的 R 被为 $R + knf$.

并让直接在前面 direct RL 为 initialize $tau(S, a) = 0$, 更新 $tau(S, a) += 1$.

(20)

2. Modify the Tabular TD(0) algorithm for estimating v_π , to estimate q_π .

Tabular TD(0) for estimating v_π

1. Input: the policy π to be evaluated
2. Algorithm parameter: step size $\alpha \in (0, 1]$
3. Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$
4. Loop for each episode:
5. Initialize S
6. Loop for each step of episode:
7. $A \leftarrow$ action given by π for S
8. Take action A , observe R, S'
9. $V(S) \leftarrow V(S) + \alpha [R + \gamma V(S') - V(S)]$
10. $S \leftarrow S'$
11. until S is terminal

其他方法: ① Monte-Carlo: 一直到 episode 结束才更新, ② 使用 G-1 (将所有得到的奖励全部加起来)

• MC prediction, estimate V_π

```
MC prediction, for estimating  $V \approx v_\pi$ 
Input: a policy  $\pi$  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  $\pi$ :  $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$ 
   $G \leftarrow 0$ 
  Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :
     $G \leftarrow \gamma G + R_{t+1}$ 
    Append  $G$  to  $Returns(S_t)$ 
   $V(S_0) \leftarrow \text{average}(Returns(S_0))$ 
```

• MC exploring starts, estimate Q_π (control)

```
MC exploring starts, for estimating  $Q \approx q_\pi$ 
Input: a policy  $\pi$  to be evaluated
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  $\pi$ :  $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$ 
   $G \leftarrow 0$ 
  Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :
     $G \leftarrow \gamma G + R_{t+1}$ 
    Append  $G$  to  $Returns(S_t, A_t)$ 
   $Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$ 
   $Q(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$ 
```

• MC control for ϵ -greedy, estimate Q_π (predict)

```
MC control (for  $\epsilon$ -soft policies), estimates  $\pi \approx \pi_*$ 
Input: a policy  $\pi$  to be evaluated
Algorithm parameter: small  $\varepsilon > 0$ 
Initialize:
   $\pi \leftarrow$  an arbitrary  $\epsilon$ -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  $\pi$ :  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ 
   $G \leftarrow 0$ 
  Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :
     $G \leftarrow \gamma G + R_{t+1}$ 
    Append  $G$  to  $Returns(S_t, A_t)$ 
     $Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$ 
     $\pi(a | S_t) \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken arbitrarily)
    For all  $a \in \mathcal{A}(S_t)$ 
       $\pi(a | S_t) \leftarrow \begin{cases} 1 - \varepsilon / |\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon / |\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$ 
```

• off-policy MC, estimate Q_π :

Off policy MC

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

Difference between prediction and control in pseudocode

tabular TD(0) for q_π :
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

target policy determine what this estimate you hope that would converge to.
 here V_π : you need the model.
 here q_π : you can only have state and action pair.
 ... the target policy has to be π_θ !!

For prediction. \wedge Which is different than the following

Sarsa (on-policy TD control) for estimating $Q \approx q_*$

From Policy Evaluation to policy improvement.
 estimating the optimal policy.

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

 Initialize S

 Choose A from S using policy derived from Q (e.g., ε -greedy)

 Loop for each step of episode:

 Take action A , observe R, S'

 Choose A' from S' using policy derived from Q (e.g., ε -greedy)

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$$

$S \leftarrow S'; A \leftarrow A'$;

 until S is terminal

control part ε -greedy. π_θ is policy derived from Q !! π_θ is behavior policy. π_θ is control.

behavior policy: policy derived from Q .

Exactly the same update.

approximating qpi

can also called SARSA for prediction!!

$\Rightarrow \pi_\theta$ policy π_θ to π_θ .

What would you modify in the above to get the pseudo code for TD(0)?

→ Given, ... 3 用被忽略的先验概率.

$$E_{\text{Final}}[P(A|S) (R + \delta V(S')) - V(S) | S=s]$$

$$= \sum_{a, s', r} P_b(A=s, C=r | S=s) P(a|s) P(r+s | S=s) \Rightarrow \text{Lotus Figure} = \sum_{x \in S} \delta(x) \cdot P(X=x)$$

由 $b(a|s)$ 中的 $\delta(s')$ 为 s' 得到了一个 a ，即 a 由先验概率被忽略到 s' ， r 。

$$= \sum_{a, s', r} b(a|s) \underbrace{P(s', r | s, a)}_{b(a|s)} (r + \delta V(S'))$$

$$= \sum_{a, s', r} P(S', r | s, a) \delta(a|s) (r + \delta V(S'))$$

$$= E_{\text{Final}}[R + \delta V(S') - V(S) | S=s]$$

3. (20)

Show that the following off-policy TD(0) update for v_π is correct for transition $s, a \sim b, s', r$:

$$v(s) \leftarrow v(s) + \alpha \left[\hat{p}(a|s)(r + \gamma v(s')) - v(s) \right]$$

by showing the following:

$$E_{A \sim b} \left[\hat{p}(a|s)(r + \gamma v(s')) - v(s) \mid s = s \right]$$

$$= E_{A \sim \pi} \left[r + \gamma v(s') - v(s) \mid s = s \right].$$

Here, $\hat{p}(a|s) = \frac{\pi(a|s)}{b(a|s)}$, π and b are

policy distributions with assumption

$$\pi(a|s) > 0 \Rightarrow b(a|s) > 0.$$

4. (20)

Give the specification of the off-policy
Expected Sarsa control method.

Ans. For the transition S, A, R, S' , where the action A_{t+1} is drawn from policy distribution b , update the action value the following way:

$$Q(S, A) \leftarrow Q(S, A) + \alpha \left[R + \gamma \sum_{a'} \pi(a'|S') Q(S', a') - Q(S, A) \right],$$

where the target policy $\pi \neq b$ and
uses greedification such as ϵ -greedy;
 $\alpha > 0$. ^{control}

(20)

5. An agent is in a 3-state MDP, $S = \{1, 2, 3\}$, where each state has two actions $A = \{1, 2\}$. Assume the agent observed the following trajectory:

$$S_0 = 1, A_0 = 1, R_1 = 1, S_1 = 2, A_1 = 2, R_2 = -1, \\ S_2 = 3, A_2 = 1, R_3 = 2, S_3 = 1, A_3 = 1, R_4 = 2, S_4 = 1.$$

The agent uses Tabular Dyna-Q.

which of the following are possible (or not possible) simulated transition $\{S, A, R, S'\}$ given the above observed trajectory with a deterministic model and random search control.

- i. $\{S=1, A=1, R=2, S'=1\}$
- ii. $\{S=3, A=1, R=2, S'=1\}$
- iii. $\{S=1, A=1, R=1, S'=2\}$
- iv. $\{S=2, A=2, R=-1, S'=3\}$
- v. $\{S=3, A=2, R_3=2, S'=1\}$.

Just mention possible or not possible for each.

Worksheet 9: Planning, Learning & Acting

CMPUT 397
March 19, 2021

2. An agent is in a 4-state MDP, $\mathcal{S} = \{1, 2, 3, 4\}$, where each state has two actions $\mathcal{A} = \{1, 2\}$. Assume the agent saw the following trajectory,

Time step 1: $S_0 = 1, A_0 = 2, R_1 = -1, \rightarrow S_1 = 1, A_1 = 1, R_2 = 1, \rightarrow S_2 = 2, A_2 = 2, R_3 = -1, \rightarrow S_3 = 2, A_3 = 1, R_4 = 1, \rightarrow S_4 = 3, A_4 = 1, R_5 = 100, \rightarrow S_5 = 4$
Final state.

and uses Tabular Dyna-Q with 5 planning steps for each interaction with the environment.

- (a) Once the agent sees S_5 , how many Q-learning updates has it done with **real experience**? How many updates has it done with **simulated experience**?
- (b) Which of the following are possible (or not possible) simulated transitions $\{S, A, R, S'\}$ given the above observed trajectory with a deterministic model and random search control?
- $\{S = 1, A = 1, R = 1, S' = 2\} \checkmark \Rightarrow S_1=1, A_1=1, R_2=1, S_2=2.$
 - $\{S = 2, A = 1, R = 1, S' = 3\} \times$
 - $\{S = 2, A = 2, R = -1, S' = 2\} \checkmark \Rightarrow S_2=2, A_2=2, R_3=-1, S_3=2.$
 - $\{S = 1, A = 2, R = -1, S' = 1\} \checkmark \Rightarrow S_0=1, A_0=2, R_1=-1, S_1=1.$
 - $\{S = 3, A = 1, R = 100, S' = 5\} \times$

a). Tabular Dyna-Q

Initialize $Q(s, a)$ and $Model(s, a)$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$
Loop forever:

- $S \leftarrow$ current (nonterminal) state
- $A \leftarrow \varepsilon$ -greedy(S, Q)
- Take action A ; observe resultant reward, R , and state, S'
- $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$ ← **real experience** (updated by real experience).
SARS' sequence.
- $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment)
- Loop repeat n times:
 - $S \leftarrow$ random previously observed state
 - $A \leftarrow$ random action previously taken in S
 - $R, S' \leftarrow Model(S, A)$ ← **simulate step**
 - $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$ ← **simulated update** (update 2 by simulated experience)

For each update we do based on the real experience, we do $n \times$ many updates using the simulated experience.

For each update we do based on the real experience, we do $n \times$ many updates using the simulated experience.

The question shows that there are 5 time steps of state-action pairs and their subsequent next state and reward that we observed from the real experience, and the question also says that we will do 5 planning steps, so $n=5$

So 5 updates based on real experience, and for each of these updates we do 5 planning steps

$\rightarrow 5 \times 5 = 25$

For each update we do based on the real experience, we do $n \times$ many updates using the simulated experience.

对于我们基于真实经验所做的每一次更新，我们使用模拟经验做5次更新。

问题显示，有5个时间步的 state-action 对及其随后的下一个状态和奖励是我们从真实经验中观察到的。问题还说我们将做5个计划步骤，所以 $n=5$

所以有5个基于真实经验的更新，对于每个更新，我们做5个规划步骤

$\rightarrow 5 \times 5 = 25$

Worksheet 12: Control with Approximation

CMPUT 397
April 7, 2022

4. (Exercise 10.6 S&B) Suppose there is an MDP that under any policy produces the deterministic sequence of rewards $+1, 0, +1, 0, +1, 0, \dots$ going on forever. Technically, this is not allowed because it violates ergodicity; there is no stationary limiting distribution μ_π and the limit

$$\lim_{t \rightarrow \infty} \mathbb{E}[R_t | S_0, A_{0:(t-1)} \sim \pi]$$

does not exist. Nevertheless, the average reward,

$$\lim_{h \rightarrow \infty} \frac{1}{h} \sum_{t=1}^h \mathbb{E}[R_t | S_0, A_{0:(t-1)} \sim \pi]$$

is well defined; What is it? Now consider two states in this MDP. From A, the reward sequence is exactly as described above, starting with a $+1$, whereas, from B, the reward sequence starts with a 0 and then continues with $+1, 0, +1, 0, \dots$ The differential return

$$G_t = R_{t+1} - r(\pi) + R_{t+2} - r(\pi) + R_{t+3} - r(\pi) + \dots$$

A: 1 0 1 0 ...
B: 0 1 0 1 0 ...

is not well defined for this case as the limit does not exist. To repair this, one could alternately define the value of a state as

$$v_\pi(s) = \lim_{\gamma \rightarrow 1} \lim_{h \rightarrow \infty} \sum_{t=0}^h \gamma^t (\mathbb{E}_\pi[R_{t+1} | S_0 = s] - r(\pi))$$

Under this definition, what are the values of states A and B?

Markov reward process:
 $C_1 = \frac{1}{1} = 1$
 $C_2 = \frac{1+0}{2} = \frac{1}{2}$
 $C_3 = \frac{1+0+1}{3} = \frac{2}{3} \rightarrow$ 2/3 is odd.
 $C_4 = \frac{1+0+1+0}{4} = \frac{1}{2} \rightarrow$ even number: all half
 $C_5 = \frac{1+0+1+0+1}{5} = \frac{3}{5} \rightarrow$ odd number: reducing. (converging to $\frac{1}{2}$, i.e. never going below 0.5!!).
 $C_6 = \frac{1+0+1+0+1+0}{6} = \frac{1}{2}$
 $C_7 = \frac{0}{7} \rightarrow$ 0/7 is odd.

$\lim_{h \rightarrow \infty} C_h = \frac{1}{2}$ is the average reward.

Here we have a period of 2: $\therefore \underline{V_{\pi}(A)} = \frac{1}{2} (2(R_1 - r(\pi)) + (R_2 - r(\pi))$
 $\qquad\qquad\qquad = \frac{1}{2} (2(1 - \frac{1}{2}) + (0 - \frac{1}{2}))$
 $\qquad\qquad\qquad = \frac{1}{2} (1 - \frac{1}{2})$
 $\qquad\qquad\qquad = \frac{1}{4}$
 $\qquad\qquad\qquad$ Starting from B, \rightarrow B LL 1/2, $\therefore R_1=0, R_2=1$.
 $\underline{V_{\pi}(B)} = \frac{1}{2} (0(R_1 - r(\pi)) + (R_2 - r(\pi))$
 $\qquad\qquad\qquad = \frac{1}{2} (0(0 - \frac{1}{2}) + (1 - \frac{1}{2}))$
 $\qquad\qquad\qquad = -\frac{1}{4}$



Confirming accuracy through Bellman equation:

$$\begin{aligned} V_{\pi}(A) &= \frac{1}{4} \\ V_{\pi}(B) &= -\frac{1}{4} \\ \therefore V_{\pi}(A) &= 1 - \frac{1}{2} + V_{\pi}(B) \\ &= 1 - \frac{1}{2} - \frac{1}{4} \\ &= \frac{1}{2} - \frac{1}{4} \\ &= \frac{1}{4} \rightarrow \text{correct!!} \\ V_{\pi}(B) &= 0 - \frac{1}{2} + V_{\pi}(A) \\ &= -\frac{1}{2} + \frac{1}{4} \\ &= -\frac{1}{4} \rightarrow \text{correct!!} \end{aligned}$$

using Bellman: $V_{\pi}(B) = \lim_{\gamma \rightarrow 1} (\gamma(1 - V_{\pi}(B)) + \gamma(0 - V_{\pi}(B)) + \gamma^2 (0 - V_{\pi}(B)) + \dots)$
 $\qquad\qquad\qquad = \lim_{\gamma \rightarrow 1} (1 - \frac{1}{2} + \gamma(0 - \frac{1}{2}) + \gamma^2 (0 - \frac{1}{2}) + \dots)$
 $\qquad\qquad\qquad = \lim_{\gamma \rightarrow 1} (1 - \frac{1}{2} - \frac{\gamma}{2} + \frac{\gamma^2}{2} \dots)$
 $\qquad\qquad\qquad = \frac{1}{2} \lim_{\gamma \rightarrow 1} (1 - \gamma + \gamma^2 - \gamma^3 \dots)$
 $\qquad\qquad\qquad = \frac{1}{2} \lim_{\gamma \rightarrow 1} \frac{1}{1-\gamma}$
 $\qquad\qquad\qquad = \frac{1}{2} \cdot \frac{1}{2} \cdot 4$
 $\qquad\qquad\qquad = \frac{1}{4}$

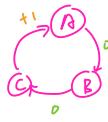
Worksheet 12: Control with Approximation

CMPUT 397
April 7, 2022

5. (Exercise 10.7 SEC) Consider a Markov reward process consisting of a ring of three states A, B, and C, with state transitions going deterministically around the ring. A reward of +1 is received upon arrival in A and otherwise the reward is 0. What are the differential values of the three states using

$$v_\pi = \lim_{\gamma \rightarrow 1} \lim_{h \rightarrow \infty} \sum_{t=0}^h \gamma^t (\mathbb{E}_\pi[R_{t+1}|S_0 = s] - r(\pi)) \rightarrow \text{If you have to find the answer equal this way, 不要直接用这个来求 answer!!}$$

Markov reward process:



$$\therefore \text{Average reward } r(\pi) = \frac{1+0+0}{3} = \frac{1}{3}$$

→ 3/3 得到 1 的 reward.

$$\left. \begin{aligned} v_\pi(C) &= (1 - \frac{1}{3}) + v_\pi(A) \\ v_\pi(B) &= (0 - \frac{1}{3}) + v_\pi(C) \end{aligned} \right\} \text{由 -1/3 得到此 state.}$$

$$v_\pi(A) = (0 - \frac{1}{3}) + v_\pi(B)$$

$$\begin{aligned} v_\pi(B) &= \frac{1}{3}(3(R_1 - r(\pi)) + 2(R_2 - r(\pi)) + R_3 - r(\pi)) \\ &= \frac{1}{3}(3(0 - \frac{1}{3}) + 2(0 - \frac{1}{3}) + 1(1 - \frac{1}{3})) \end{aligned}$$

由 1/3 得到, reward factor: 0 → 0 → 1

$$= \frac{1}{3}(-1 - \frac{2}{3} + \frac{2}{3})$$

$$= -\frac{1}{3}$$

∴ can find $v_\pi(B)$ and $v_\pi(C)$ by Bellman equation. 但因为用了 Bellman equation 找 answer, 所以你用这种方法就错.

$$v_\pi(C) = 1 - \frac{1}{3} + v_\pi(A) = \frac{2}{3} - \frac{1}{3} = \frac{1}{3}$$

$$v_\pi(A) = 0 - \frac{1}{3} + v_\pi(C) = -\frac{1}{3} + \frac{1}{3} = 0$$

由 1/3 得到 0 的 reward, 由 1 得到 reward 为 0, ∴ 0 是 valuable 的奖励, ∴ $v_\pi(A)$ 为 0.

$$B \quad \text{2 reward } (2/3) \times 0.1, \therefore B$$

由 1/3, ∴ $v_\pi(B) = 0$.

$$C \quad 1 - \text{reward } (1/3) \times 0.1, \therefore C$$

最高, ∴ $v_\pi(C) = 1$.

∴ Q.

8. (15)

What is the difference between the following two updates? Which one is a stochastic-gradient and a semi-gradient update?

a. $\underline{w}_{t+1} \stackrel{?}{=} \underline{w}_t -$

$$\alpha \left[\nabla_{\underline{w}_t} \frac{1}{2} \left(R_{t+1} + \gamma \hat{v}(S_{t+1}, \underline{w}) - \hat{v}(S_t, \underline{w}_t) \right)^2 \right]$$

半梯度因为梯度中是中 $\nabla_{\underline{w}_t} \underline{w}_t$!

$$= \underline{w}_t + \frac{1}{2} \alpha [2(R_{t+1} + \gamma \hat{v}(S_{t+1}, \underline{w}) - \hat{v}(S_t, \underline{w}_t)) \nabla_{\underline{w}_t} \hat{v}(S_t, \underline{w}_t)]$$

\underline{w}_t

b. $\underline{w}_{t+1} \stackrel{?}{=} \underline{w}_t -$

$$\alpha \left[\nabla_{\underline{w}_t} \frac{1}{2} \left(R_{t+1} + \gamma \hat{v}(S_{t+1}, \underline{w}_t) - \hat{v}(S_t, \underline{w}_t) \right)^2 \right]$$

半梯度: 涵包梯上梯式中所有含有 \underline{w}_t 的项.

$$= \underline{w}_t + \frac{1}{2} \alpha \cdot 2(R_{t+1} + \gamma \hat{v}(S_{t+1}, \underline{w}_t) - \hat{v}(S_t, \underline{w}_t))$$

$$\cdot (\nabla_{\underline{w}_t} \hat{v}(S_t, \underline{w}_t) - \nabla_{\underline{w}_t} \hat{v}(S_{t+1}, \underline{w}_t))$$

↓
调速慢, 得以将之前的梯度为正.

TD(0) 中没有, 而 TD(1) 中含有中, 不能表示 $\nabla_{\underline{w}_t}$.

Ans.

(a) becomes

$$\underline{\omega}_{t+1} = \underline{\omega}_t - \alpha \left[\nabla_{\underline{\omega}_t} \frac{1}{2} (R_{t+1} + \gamma \hat{v}(S_{t+1}, \underline{\psi}) - \hat{v}(S_t, \underline{\omega}_t)) \right] \Big|_{\underline{\psi} = \underline{\omega}_t}$$

对这一项进行求导

因为对 $\underline{\omega}_t$ 进行求导，所以梯度矩阵上含有 $\underline{\omega}_t$ 的项。

$$= \underline{\omega}_t + \frac{1}{2} \alpha \left[2(R_{t+1} + \gamma \hat{v}(S_{t+1}, \underline{\psi}) - \hat{v}(S_t, \underline{\omega}_t)) \nabla_{\underline{\omega}_t} \hat{v}(S_t, \underline{\omega}_t) \right] \Big|_{\underline{\psi} = \underline{\omega}_t}$$

$$= \underline{\omega} + \alpha (R_{t+1} + \gamma \hat{v}(S_{t+1}, \underline{\omega}_t) - \hat{v}(S_t, \underline{\omega}_t)) \nabla_{\underline{\omega}_t} \hat{v}(S_t, \underline{\omega}_t),$$

which is the semi-gradient TD(0) update.

Therefore, this one is semi-gradient.

(b) becomes

$$\underline{\omega}_{t+1} = \underline{\omega}_t - \alpha \nabla_{\underline{\omega}_t} \frac{1}{2} (R_{t+1} + \gamma \hat{v}(S_{t+1}, \underline{\omega}_t) - \hat{v}(S_t, \underline{\omega}_t))^2$$

$$= \underline{\omega}_t + \frac{1}{2} \alpha 2(R_{t+1} + \gamma \hat{v}(S_{t+1}, \underline{\omega}_t) - \hat{v}(S_t, \underline{\omega}_t)) \times (\nabla_{\underline{\omega}_t} \hat{v}(S_t, \underline{\omega}_t) - \nabla_{\underline{\omega}_t} \hat{v}(S_{t+1}, \underline{\omega}_t)),$$

which is a gradient update and different from TD(0).

The main difference is the last term that does not exist in TD(0).