

Notes on Reinforcement Learning

Randy Sun

March 11, 2018

1 Markov Decision Process

This section introduces several basic elements of Markov Decision Process. Also, some important equations would be presented.

1.1 Fundamental elements

Markov Decision Process (MDP) is a tuple, $(S, A, P, \pi, R, \gamma)$.

- S is a finite set which contains a set of states, $\{s_1, s_2, \dots, s_N\}$. $S(s, a) \subseteq S$ is a set of states which are available after state s and action a .
- A is an action set which contains a set of actions, $\{a_1, a_2, \dots, a_M\}$. $A(s) \subseteq A$ is a set of actions which are available after action s .
- P are transition probability measures.

$$P_{s,s'}^a = \Pr(s_{t+1} = s' | s_t = s, a_t = a). \quad (1.1)$$

- π is policy. $\pi(a_t = a | s_t = s)$ is the probability of taking an action, a , under state s and policy, π . Moreover,

$$\pi(s_{t+1} = s' | s_t = s) = \sum_{a \in A(s)} \pi(a_t = a | s_t = s) P_{s,s'}^a. \quad (1.2)$$

- R are rewards.

$$R_{t+1} = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+1+k} \quad (1.3)$$

where

$$r_{t+1} = r(s_t, a_t, s_{t+1}) \quad (1.4)$$

depend on the current state, s_{t+1} , last state, s_t , and last action, a_t .

$$R_s^a = E_{\pi}[r_{t+1} | s_t = s, a_t = a] = \sum_{s' \in S(s,a)} P_{s,s'}^a r_{t+1}(s_t = s, a_t = a, s_{t+1} = s'). \quad (1.5)$$

- γ is the discount rate, which is smaller than 1.

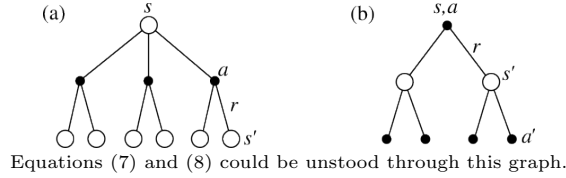
There are two important functions, which play the crucial roles in reinforcement learning, *state-value function for policy π* , V^π and *action-value function for policy π* , Q^π .

$$\begin{aligned} V^\pi(s) &= E_\pi [R_{t+1} | s_t = s] \\ &= \sum_{a \in A(s)} \pi(a_t = a | s_t = s) Q^\pi(s, a) \end{aligned} \quad (1.6)$$

means the expected reward given the current state is s under policy π , and

$$\begin{aligned} Q^\pi(s, a) &= E_\pi [R_{t+1} | s_t = s, a_t = a] \\ &= \sum_{s' \in S(s, a)} P_{s, s'}^a (r(s_t = s, a_t = a, s_{t+1} = s') + \gamma V^\pi(s')) \\ &= R_s^a + \gamma \sum_{s' \in S(s, a)} P_{s, s'}^a V^\pi(s') \end{aligned} \quad (1.7)$$

means the expected reward given the current state is s and current action a under policy π . As the state-value function, action-value function can be written as the linear combination of state-value function, either.



Reinforce Learning cares about looking for the *best policy*, π^* . It aims to estimate

$$\begin{aligned} \max_{\pi} V^\pi(s) &= V^{\pi^*}(s) \\ &= \max_{a \in A(s)} Q^{\pi^*}(s, a) \end{aligned} \quad (1.8)$$

and

$$\max_{\pi} Q^\pi(s, a) = Q^{\pi^*}(s, a). \quad (1.9)$$

It's intuitively to write down the last equation in (9), since the best policy under a state, s , would be taking the action with the maximum expected reward.

1.2 Bellman Equation

To estimate, we rewrite V^π , Q^π , V^{π^*} , Q^{π^*} by the form of Bellman equation.

$$\begin{aligned}
V^\pi(s) &= E_\pi[R_{t+1}|s_t = s] \\
&= E_\pi[r_{t+1} + \gamma R_{t+2}|s_t = s] \\
&= \sum_{a, s'} \Pr_\pi(a_t = a, s_{t+1} = s' | s_t = s) (E_\pi[r_{t+1} + \gamma R_{t+2} | s_t = s, a_t = a, s_{t+1} = s']) \\
&= \sum_{a, s'} \Pr_\pi(a_t = a, s_{t+1} = s' | s_t = s) (r_{t+1}(s_t = s, a_t = a, s_{t+1} = s') + E_\pi[\gamma R_{t+2} | s_t = s, a_t = a, s_{t+1} = s']) \\
&= \sum_{a, s'} \Pr_\pi(a_t = a, s_{t+1} = s' | s_t = s) (r_{t+1}(s_t = s, a_t = a, s_{t+1} = s') + E_\pi[\gamma R_{t+2} | s_{t+1} = s']) \\
&= \sum_{a, s'} \Pr_\pi(a_t = a, s_{t+1} = s' | s_t = s) (r_{t+1}(s_t = s, a_t = a, s_{t+1} = s') + \gamma V^\pi(s')) \\
&= \sum_{a \in A(s)} \pi(a_t = a | s_t = s) \sum_{s' \in S(s, a)} P_{s, s'}^a (r_{t+1}(s_t = s, a_t = a, s_{t+1} = s') + \gamma V^\pi(s'))
\end{aligned} \tag{1.10}$$

and

$$\begin{aligned}
Q^\pi(s, a) &= E_\pi[R_{t+1} | s_t = s, a_t = a] \\
&= E_\pi[r_{t+1} + \gamma R_{t+2} | s_t = s, a_t = a] \\
&= \sum_{s', a'} \Pr_\pi(a_{t+1} = a', s_{t+1} = s' | s_t = s, a_t = a) (E_\pi[r_{t+1} + \gamma R_{t+2} | s_t = s, a_t = a, s_{t+1} = s', a_{t+1} = a']) \\
&= \sum_{s', a'} \Pr_\pi(a_{t+1} = a', s_{t+1} = s' | s_t = s, a_t = a) (R_{s, s'}^a + \gamma Q^\pi(s', a')) \\
&= R_s^a + \sum_{s' \in S(s, a)} P_{s, s'}^a \gamma V^\pi(s').
\end{aligned} \tag{1.11}$$

Actually, the equations (12) and (13) could be obtained directly from equations (7) and (8).

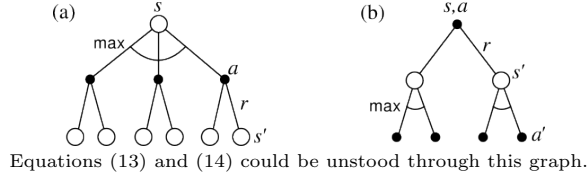
To have the Bellman equation form of V^{π^*} , we can move on from equation (9),

$$\begin{aligned}
V^{\pi^*}(s) &= \max_a Q^{\pi^*}(s, a) \\
&= \max_a \left(R_s^a + \gamma \sum_{s' \in S(s, a)} P_{s, s'}^a V^{\pi^*}(s') \right).
\end{aligned} \tag{1.12}$$

As for Q^{π^*} , since equation (12) holds for any policy, including best policy π^* ,

combing equations (12) and (13) could obtain the result.

$$\begin{aligned}
Q^{\pi^*}(s, a) &= R_s^a + \gamma \sum_{s' \in S(s, a)} P_{s, s'}^a V^{\pi^*}(s') \\
&= R_s^a + \gamma \sum_{s' \in S(s, a)} P_{s, s'}^a \max_a Q^{\pi^*}(s', a).
\end{aligned} \tag{1.13}$$



2 Model-based Dynamic Programming

In this section, we will introduce policy iteration and value iteration. Moreover, we will introduce policy improvement theorem.

2.1 Policy iteration

Policy iteration is an model-based algorithm to find the best policy and its corresponding state-value function. It needs $P_{s, s'}^a$ as the input for the algorithm, so it is classified to the model-based category. Policy iteration algorithm consists of two parts: *Policy evaluation* and *Policy improvement*.

Given a policy π , the algorithm of Policy evaluation aims to estimate $V^\pi(s) \forall s \in S$. It estimates the value iteratively by model-based dynamic programming.

Algorithm 1 Policy evaluation

Input: A policy π , transition probabilities $P_{s, s'}^a$, expected rewards R_s^a , and discount rate γ .

Output: state-value function $V^\pi(s) \forall s \in S$.

- 1: Initialize $V^\pi(s) = 0 \forall s \in S$.
 - 2: **while** True **do**
 - 3: **for** $s \in S$ **do**
 - 4: $V^\pi(s) = \sum_{a \in A(s)} \pi(a_t = a | s_t = s) \left(R_s^a + \gamma \sum_{s' \in S(s, a)} P_{s, s'}^a V^\pi(s') \right)$
 - 5: **end for**
 - 6: **if** $V_{k+1}^\pi(s) = V_k^\pi(s) \forall s \in S$ **then**
 - 7: break
 - 8: **end if**
 - 9: **end while**
 - 10: Return $V^\pi(s) \forall s \in S$.
-

Now, we are able to evaluate a policy. We want to improve the policy. By the policy improvement theorem (we leave the proof in subsection 3), if

$$V^\pi(s) \leq Q^\pi(s, \pi'(s)) \quad (2.1)$$

for all s ,

$$V^\pi(s) \leq V^{\pi'}(s) \quad (2.2)$$

for all s . So, to have better policy, we use equation (1.7) to find the better policy π' . Let

$$\pi'(s) = \operatorname{argmax}_a \sum_{s' \in S(s,a)} P_{s,s'}^a (r(s_t = s, a_t = a, s_{t+1} = s') + \gamma V^\pi(s')) \quad (2.3)$$

for each s . The algorithm is as follows:

Algorithm 2 Policy improvement

```

1: while True do
2:   policy stable = True
3:   for each  $s$  in  $S$  do
4:     Let  $b = \pi(s)$ .
5:     Let  $\pi(s) = \operatorname{argmax}_a \sum_{s' \in S(s,a)} P_{s,s'}^a (r(s_t = s, a_t = a, s_{t+1} = s') + \gamma V^\pi(s'))$ 
6:     if  $b \neq \pi(s)$  then
7:       policy stable = False
8:     end if
9:   end for
10:  if policy stable then
11:    break
12:  end if
13: end while

```

Now, we are able to improve the policy. Then, we could combine the policy evaluation and policy improvement to find the best policy. The whole process of

Algorithm 3 Policy iteration

```

1: Initialize  $\pi(s) \in A(s)$  and  $V^\pi(s)$  arbitrarily.
2: while True do
3:   Policy evaluation on  $\pi$ .
4:   Policy improvement on  $\pi$  and get  $\pi'$ .
5:   if  $\pi = \pi'$  then
6:     break
7:   end if
8:    $\pi = \pi'$ 
9: end while

```

policy iteration could be unstood as the following graph:

$$\pi_0 \xrightarrow{E} V^{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V^{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} \dots \xrightarrow{I} \pi^* \xrightarrow{E} V^*,$$

2.2 Value iteration

Although there are some guarantees that the policy iteration would converge to the best policy, it sometimes suffers from protractive procedure. There are other ways to find the best policy. Instead of iteratively finding the best policy, one could try to find the state-value function directly in first, then use the state-value function to find the best policy.

Value iteration uses the following way to update the state-value function:

$$V_{n+1}(s) = \max_a \sum_{s' \in S(s,a)} P_{s,s'}^a (r(s_t = s, a_t = a, s_{t+1} = s') + \gamma V^\pi(s')) \quad (2.4)$$

Then, it would converge to the value-state function of the best policy. Lastly, use equation (1.6) to find the best policy. Note that the difference between the

Input: θ small enough.

```

1: Initialize  $V(s) = 0 \forall s \in S$ .
2: while True do
3:    $d = 0$ .
4:   for  $s \in S$  do
5:      $v = V(s)$ .
6:      $V(s) = \max_a \sum_{s' \in S(s,a)} P_{s,s'}^a (r(s_t = s, a_t = a, s_{t+1} = s') + \gamma V^\pi(s'))$ .
7:      $d = \max(d, |v - V(s)|)$ .
8:   end for
9:   if  $d \leq \theta$  then
10:    break
11:  end if
12: end while
13: Let  $\pi(s) = \operatorname{argmax}_a \sum_{s' \in S(s,a)} P_{s,s'}^a (r(s_t = s, a_t = a, s_{t+1} = s') + \gamma V^\pi(s'))$ 
     $\forall s \in S$ .
```

two update steps of value iteration and policy iteration could be understood through figure(1)(a) and figure(2)(a).

2.3 Policy Improvement Theorem

Given that $\forall s \in S, V^\pi(s) \leq Q^\pi(s, \pi'(s))$, then,

$$\begin{aligned}
V^\pi(s) &\leq Q^\pi(s, \pi'(s)) \\
&= \sum_{s' \in S(s, \pi'(s))} P_{s, s'}^{\pi'(s)} \left(R_{s, s'}^{\pi'(s)} + \gamma V^\pi(s') \right) \\
&\leq \sum_{s' \in S(s, \pi'(s))} P_{s, s'}^{\pi'(s)} \left(R_{s, s'}^{\pi'(s)} + \gamma Q^\pi(s', \pi'(s')) \right) \\
&= \sum_{s' \in S(s, \pi'(s))} P_{s, s'}^{\pi'(s)} \left(R_{s, s'}^{\pi'(s)} + \gamma \sum_{s'' \in S(s', \pi'(s'))} P_{s', s''}^{\pi'(s')} \left(R_{s', s''}^{\pi'(s')} + \gamma V^\pi(s'') \right) \right) \\
&= \sum_{s' \in S(s, \pi'(s))} P_{s, s'}^{\pi'(s)} r_{t+1}(s, \pi'(s), s') + \gamma \sum_{s'' \in S(s', \pi'(s'))} P_{s', s''}^{\pi'(s')} r_{t+2}(s', \pi'(s'), s'') + \dots \\
&\leq \dots\dots \\
&\leq \dots\dots \\
&= \sum_{a, s', a', s'', \dots} \Pr_{\pi'}(a_t = a, s_{t+1} = s', a_{t+1} = a', \dots | s_t = s) (r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} \dots) \\
&= V^{\pi'}(s).
\end{aligned} \tag{2.5}$$

That is to say, if the expected reward under a state s and an action $a = \pi'(s)$ is greater than expected reward under state s , $\forall s \in S$, the state value function of the new policy π' would be greater than the old one.

3 Monte Carlo estimate

Monte Carlo control is an approach to estimate the value-state function or action-state function. It is based on sampling many episodes and try to estimate it by averaging all the returns. In general, there are two ways to estimate, *first-visit MC estimate* and *every-visit MC estimate*. They share some crucial properties but have slightly different theoretical properties. Here, we consider the former one.

Algorithm 4 First-visit MC control for estimate state-value function

```
1: Initialize policy  $\pi$ , state-value function  $V^\pi$  and empty list  $L(s) \forall s \in S$ .
2: while True do
3:   Generate an episode with  $\pi$ .
4:   for each state  $s$  in the episode do
5:     Let  $G$  be the return of the state  $s$ .
6:     Append  $G$  to  $L(s)$ .
7:      $V(s) = \text{average}(L(s))$ .
8:   end for
9: end while
```

3.1 On-policy control

Before we move on to on-policy approach, we should discuss some thoughts about policy. We introduce several common stochastic strategies below,

1. *Greedy strategy*

$$\pi(a_t = a | s_t = s) = \begin{cases} 1 & \text{if } a = \underset{a}{\operatorname{argmax}} Q^\pi(s, a) \\ 0 & \text{if } a \neq \underset{a}{\operatorname{argmax}} Q^\pi(s, a) \end{cases} \quad (3.1)$$

Greedy strategy selects the action with the largest reward, it is *deterministic*.

2. *ϵ -greedy strategy*

$$\pi(a_t = a | s_t = s) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{|A(s)|} & \text{if } a = \underset{a}{\operatorname{argmax}} Q^\pi(s, a) \\ \frac{\epsilon}{|A(s)|} & \text{if } a \neq \underset{a}{\operatorname{argmax}} Q^\pi(s, a) \end{cases} \quad (3.2)$$

ϵ -greedy strategy strikes a balance between *exploitation* and *exploration*. It makes sure that all actions would have the chances to be executed.

Here, *exploitation* means always selecting the action with the largest expected return (greedy strategy). And the *exploration* means we should select other possible actions so that we may find out a better action than the action we see it as the best one.

Basically, there are two types of MC control: *with exploring starts* and *without exploring starts*. We first introduce the former algorithm.

Algorithm 5 Monte Carlo control with exploring starts.

```

1: Initialize policy  $\pi$ , action-value function  $Q^\pi$  and empty list  $L(a, s) \forall s \in S$  and  $a \in A(s)$ .
2: while True do
3:   Choose  $s_0$  and  $a_0$  s.t all pairs have probabilities  $> 0$ .
4:   Generate all pairs from  $s_0$  and  $a_0$ .
5:   for each pair  $s$  and  $a$  in the episode do
6:     Let  $G$  be the return of first occurrence of the pair  $s$  and  $a$ .
7:     Append  $G$  to  $L(s, a)$ .
8:      $Q(s, a) = \text{average}(L(s, a))$ .
9:   end for
10:  for each state  $s$  in the episode do
11:     $\pi(s) = \underset{a}{\operatorname{argmax}} Q^\pi(s, a)$ .
12:  end for
13: end while

```

Note that step 3 makes sure all the pairs would appear as the beginnings of episodes.

Now, could we avoid the human efforts in step 3 and 4? Apparently, if we use greedy strategy, it won't work. So, we should make our policy non-greedy, i.e, we hope our policy to be $\epsilon - soft$.

Algorithm 6 Monte Carlo control without exploring starts.

```

1: Initialize policy  $\pi$  to be  $\epsilon - greedy$ , state-value function  $V^\pi$  and empty list  $L(a, s) \forall s \in S$  and  $a \in A(s)$ .
2: while True do
3:   Generate an episode with  $\pi$ .
4:   for each pair  $s$  and  $a$  in the episode do
5:     Let  $G$  be the return of first occurrence of the pair  $s$  and  $a$ .
6:     Append  $G$  to  $L(s, a)$ .
7:      $Q(s, a) = \text{average}(L(s, a))$ .
8:   end for
9:   for each state  $s$  in the episode do
10:     $a^* = \underset{a}{\operatorname{argmax}} Q^\pi(s, a)$ .
11:
12:    for  $a \in A$  do
13:       $\pi(a|s) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{|A|} & \text{if } a = a^* \\ \frac{\epsilon}{|A|} & \text{if } a \neq a^* \end{cases}$ 
14:    end for
15:  end for
16: end while

```

This seems to be a good strategy. However, could it better the original

policy? All we want to know is if $Q^\pi(s, \pi'(s)) \geq V^\pi(s) \forall s \in S$?

$$\begin{aligned}
Q^\pi(s, \pi'(s)) &= \sum_a \pi'(a|s) Q^\pi(s, a) \\
&= \frac{\epsilon}{|A|} \sum_a Q^\pi(s, a) + (1 - \epsilon) \max_a Q^\pi(s, a) \\
&\geq \frac{\epsilon}{|A|} \sum_a Q^\pi(s, a) + (1 - \epsilon) \sum_a \frac{\pi(a|s) - \frac{\epsilon}{|A|}}{1 - \epsilon} Q^\pi(s, a) \\
&= \frac{\epsilon}{|A|} \sum_a Q^\pi(s, a) - \frac{\epsilon}{|A|} \sum_a Q^\pi(s, a) + \sum_a \pi(a|s) Q^\pi(s, a) \\
&= V^\pi(s)
\end{aligned} \tag{3.3}$$

So, this algorithm works.

3.2 Off-policy control

As stated previously, we estimate the target policy, π , while following another policy, π' . This strategy could be successful under the idea of *important sampling*. In general, we may want to estimate $E(f(X))$ for any possible function f . Nevertheless, sometimes it's hard to sample, π , of the random variable X . Thus, we use another distribution π' to estimate it.

$$\begin{aligned}
E_\pi(f(X)) &= \int_{\Omega} f(x) \pi(x) dx \\
&= \int_{\Omega} f(x) \frac{\pi(x)}{\pi'(x)} \pi'(x) dx \\
&= E_{\pi'} \left(f(X) \frac{\pi(X)}{\pi'(X)} \right) \\
&\underset{N \rightarrow \infty}{\approx} \sum_{i=1}^N \frac{f(x^i) \frac{\pi(x^i)}{\pi'(x^i)}}{N} \text{ where } x^i \sim \pi'.
\end{aligned} \tag{3.4}$$

Define

$$w_i = \frac{\pi(x^i)}{\pi'(x^i)} \tag{3.5}$$

then equation (3.1) could be written as

$$\sum_{i=1}^N \frac{f(x^i) w_i}{N} \tag{3.6}$$

Obviously, this is an unbiased estimator of the target. However, variance of the estimation would be pretty large if the π' is not selected well. Generally, we

would try to estimate the target by

$$E_{\pi}(f(X)) \approx \sum_{i=1}^N \frac{w_i}{\sum_{j=1}^N w_j} f(x^i) \quad (3.7)$$

This could effectively reduce the estimated variance.

Back to MDP, since under policy, π

$$\begin{aligned} \Pr(a_t, s_{t+1}, \dots, s_T | s_t) &= \prod_{k=t}^{T-1} \pi(a_k | s_k) \Pr(s_{k+1} | s_k, a_k) \\ &= \prod_{k=t}^{T-1} \pi(a_k | s_k) P_{s_k, s_{k+1}}^{a_k} \end{aligned} \quad (3.8)$$

Similarly, under policy, π'

$$\Pr(a_t, s_{t+1}, \dots, s_T | s_t) = \prod_{k=t}^{T-1} \pi'(a_k | s_k) P_{s_k, s_{k+1}}^{a_k}. \quad (3.9)$$

Now, the weight then becomes

$$w = \frac{\prod_{k=t}^{T-1} \pi(a_k | s_k)}{\prod_{k=t}^{T-1} \pi'(a_k | s_k)}. \quad (3.10)$$

It has a remarkable property, i.e it doesn't need the model(conditional probability).

To estimate $V^{\pi}(s)$, let w_i be the weight of the i th episode which contains state s and the R^i represents the discount reward of the i th episode. Then the off-policy estimates $V^{\pi}(s) = \frac{\sum_i w_i R^i}{\sum_i w_i}$.

In practice, when we want to estimate a deterministic policy π , we would use the above approach to estimate $Q^{\pi}(s, a)$, and then use following algorithm to estimate a deterministic policy π :

Algorithm 7 Off-policy Monte Carlo Control

Input: behavior policy π' .

Output: deterministic target policy π .

- 1: $\forall s \in S, a \in A$, initialize $Q^\pi(s, a) = 0$, $N(s, a) = 0$, $D(s, a) = 0$ and π be any arbitrary deterministic policy.
 - 2: **while** True **do**
 - 3: Use policy π' to generate an episode : $s_0, a_0, r_1, s_1, a_1, \dots, r_T, s_T$.
 - 4: Let τ be the latest time t where $a_t \neq \pi(s_t)$
 - 5: **for** each pair s, a appearing in the episode after τ **do**
 - 6: Let t be the first occurrence of s, a after τ .
 - 7: $w = \prod_{k=t+1}^{T-1} \frac{1}{\pi'(a_k|s_k)}$.
 - 8: $N(s, a) = N(s, a) + wR_{t+1}$.
 - 9: $D(s, a) = D(s, a) + w$.
 - 10: $Q(s, a) = \frac{N(s, a)}{D(s, a)}$
 - 11: **end for**
 - 12: **for** each $s \in S$ **do**
 - 13: Let $\pi(s) = \underset{a}{\operatorname{argmax}} Q(s, a)$
 - 14: **end for**
 - 15: **end while**
-

Note that step 6 indicates we have to consider the process after s and a , so that we are actually estimating $Q^\pi(s, a)$ not $V^\pi(s)$.

3.3 Incremental Implementation

In last subsection, we state that we could use this estimation, $V^\pi = \frac{\sum w_i R^i}{\sum w_i}$, to estimate the value-state function. However, in practice, we could use incremental implementation to write down the estimation as

$$V_{n+1}^\pi = V_n^\pi + \frac{w_{n+1}}{W_{n+1}} [R^{n+1} - V_n^\pi] \quad (3.11)$$

where

$$W_{n+1} = W_n + w_{n+1} \quad (3.12)$$

Thus Algorithm 2 could be modified as

Algorithm 8 Off-policy Monte Carlo Control

Input: behavior policy π' .

Output: deterministic target policy π .

```
1:  $\forall s \in S, a \in A$ , initialize  $Q^{\pi'}(s, a) = 0$ ,  $N(s, a) = 0$ ,  $D(s, a) = 0$  and  $\pi$  be  
   any arbitrary deterministic policy.  
2: while True do  
3: Use policy  $\pi'$  to generate an episode :  $s_0, a_0, r_1, s_1, a_1, \dots, r_T, s_T$ .  
4: Let  $w = 1$  and  $R = 0$ .  
5:   for  $t = T - 1, T - 2, \dots, 0$  do  
6:      $R = \gamma R + r_{t+1}$ .  
7:      $D(s_t, a_t) = D(s_t, a_t) + w$ .  
8:      $Q(s_t, a_t) = Q(s_t, a_t) + \frac{w}{D(s_t, a_t)} [R - Q(s_t, a_t)]$   
9:      $\pi(s_t) = \underset{a}{\operatorname{argmax}} Q(s_t, a)$   
10:    if  $a_t \neq \pi(s_t)$  then  
11:      break  
12:    end if  
13:     $W = W \frac{1}{\pi(a_t|s_t)}$   
14:  end for  
15: end while
```

4 Temporal Difference

Recall that Monte Carlo approach estimate the target by *sampling*(Monte Carlo) and dynamic programming approach estimate the target by *bootstrapping*. Temporal difference(TD) approach combines the two methods to estimate the target.

In section 3, we see that the estimate could be done by

$$V^\pi(s) = V^\pi(s) + \alpha [R_{t+1} - V^\pi(s)] \quad (4.1)$$

for some α . Monte Carlo approach uses the *sampled return* to update the value-state function. Different from Monte Carlo approach, temporal difference method use the result estimated before to estimate, i.e.

$$V^\pi(s_t) = V^\pi(s_t) + \alpha [r_{t+1} + \gamma V^\pi(s_{t+1}) - V^\pi(s_t)]. \quad (4.2)$$

Theoretically, the difference between temporal difference and monte carlo approach could be identified by the following equations.

$$\begin{aligned} V^\pi(s_{t+1}) &= E_\pi [R_{t+1} | s_t = s] \\ &= E_\pi [r_{t+1} + \gamma V^\pi(s_{t+1}) | s_t = s] \end{aligned} \quad (4.3)$$

Monte Carlo approach estimates the function by the first equation of (4.3); Nevertheless, the temporal difference approach estimate it by the second one. The complete algorithm is as follows:

Algorithm 9 Temporal Difference

Input: target policy π' .

- 1: Initialize $V(s)$ arbitrary.
 - 2: Initialize s .
 - 3: **for** each step, s , of an episode **do**
 - 4: Let a be actopm given by state s and π .
 - 5: Take action a and observe s' and r .
 - 6: $V(s) = V(s) + \alpha[r + \gamma V(s') - V(s)]$
 - 7: $s' = s$
 - 8: **end for**
-

4.1 On-policy TD control–Sarsa algorithm

Sarsa algorithm uses the temporal difference to estimate the action-value function Q^π by the same way. And it uses the ϵ – *greedy* strategy.

Algorithm 10 Sarsa algorithm

Input: arbitrary target policy π and α .

- 1: Initialize $Q(s, a)$ arbitrarily.
 - 2: **while** True **do**
 - 3: Initialize s .
 - 4: Choose a from s and π derived from Q .
 - 5: **while** s is not terminal state **do**
 - 6: Take action a and observe s' and r' .
 - 7: Choose a' from s' and π derived from Q .
 - 8: $Q(s, a) = Q(s, a) + \alpha[r' + \gamma Q(s', a') - Q(s, a)]$
 - 9: $s = s'$ and $a = a'$.
 - 10: **end while**
 - 11: **end while**
-

4.2 Off-policy TD control–Q–learning

Sarsa algorithm follows the same policy π for every step. However, Q–learning directly estimate Q^* , instead of Q^π (Sarsa algorithm). Although both algorithm converge to optimal policy π^* , Q–learning dramatically simplifies the analysis of the algorithm.

Algorithm 11 Q-learning algorithm

Input: arbitrary target policy π and α .

```
1: Initialize  $Q(s, a)$  arbitrarily.  
2: while True do  
3:   Initialize  $s$ .  
4:   while  $s$  is not terminal state do  
5:     Choose  $a$  from  $s$  and  $\pi$  derived from  $Q$ .  
6:     Take action  $a$  and observe  $s'$  and  $r'$ .  
7:      $Q(s, a) = Q(s, a) + \alpha \left[ r' + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$ .  
8:      $s = s'$ .  
9:   end while  
10: end while
```
