

Reinforcement Learning

An Introductory Note

Jingye Wang

✉ wangjy5@shanghaitech.edu.cn

Spring 2020

Contents

1	Introduction	3
2	Review of Basic Probability	5
2.1	Interpretation of Probability	5
2.2	Transformations	5
2.3	Limit Theorem	5
2.4	Sampling & Monte Carlo Methods	6
2.5	Basic Inequalities	8
2.6	Concentration Inequalities	10
2.7	Conditional Expectation	12
3	Bandit Algorithms	14
3.1	Bandit Models	14
3.2	Stochastic Bandits	14
3.3	Greedy Algorithms	15
3.4	UCB Algorithms	16
3.5	Thompson Sampling Algorithms	17
3.6	Gradient Bandit Algorithms	18
4	Markov Chains	20
4.1	Markov Model	20
4.2	Basic Computations	20
4.3	Classifications	21

CONTENTS	2
4.4 Stationary Distribution	22
4.5 Reversibility	22
4.6 Markov Chain Monte Carlo	23
5 Markov Decision Process	25
5.1 Markov Reward Process	25
5.2 Markov Decision Process	26
5.3 Dynamic Programming	28
6 Model-Free Prediction	33
6.1 Monte-Carlo Policy Evaluation	33
6.2 Temporal-Difference Learning	35
7 Model-Free Control	37
7.1 On Policy Monte-Carlo Control	37
7.2 On Policy Temporal-Difference Control: Sarsa	39
7.3 Off-Policy Temporal-Difference Control: Q-Learning	40
8 Value Function Approximation	41
8.1 Semi-gradient Method	41
8.2 Deep Q-Learning	43
9 Policy Optimization	46
9.1 Policy Optimization Theorem	46
9.2 REINFORCE: Monte-Carlo Policy Gradient	49
9.3 Actor-Critic Policy Gradient	51
9.4 Extension of Policy Gradient	52

7 Model-Free Control

In the last section, we introduce Monte-Carlo (MC) and Temporal Difference (TD) methods to evaluate the value function of a policy. With the idea of generalized policy iteration (GPI) we mentioned in section 5.3.2, we now consider how the value function can be used in control, that is, to find the optimal policy.

For an MDP with large state space, the way we sample episodes matters a lot. We define *behavior policy* that determines which action to take and *target policy* that determines the best policy we have so far, then based on the two concepts, we have two learning form:

- **On-policy** learning: the target policy and the behavior policy are the same, which means the action we will take follows the optimal policy we find;
- **Off-policy** learning: the target policy and the behavior policy are different, which means the action we will take is independent to the optimal policy we find.

On-policy may not ensure the enough exploration of the state space as when we update the policy greedily we may discard those states with great potential. Compared with on-policy learning, off-policy learning is more powerful and general, though it is often of greater variance and is slower to converge.

7.1 On Policy Monte-Carlo Control

Like MC evaluation we mentioned in section 6.1, MC control also requires fully episodes to improve the policy. It is natural to alternate between evaluation and improvement on an episode-by-episode basis. The following algorithms are based on the implementation of MC, and the differences lie in how they generate samples.

Exploring Starts: *To obtain diverse episodes, a naive ideal is initializing each episode randomly:*

Random starts usually are not common in reality. In the iteration part, we can leverage ε -Greedy Exploration to avoid the unlikely assumption of exploring starts.

ε -Greedy Exploration: *For a state with m actions, there is non-zero probability to try them:*

$$\pi(a|s) = \begin{cases} \frac{\varepsilon}{m} + 1 - \varepsilon, & a = \arg \max_{a' \in \mathcal{A}} q(s, a') \\ \frac{\varepsilon}{m}, & \text{otherwise} \end{cases}$$

which means we will choose the greedy action with probability $1 - \varepsilon$ and choose an action at random with probability ε .

Policy improvement theorem: *For any ε -greedy policy π , the ε -greedy policy π' with respect to q^π is an improvement, i.e., $v_{\pi'}(s) \geq v_\pi(s)$ for all $s \in \mathcal{S}$.*

Algorithm 13 Monte-Carlo Control

```

1: initialize  $\pi(s) \in \mathcal{A}(s)$ (arbitrarily), for all  $s \in \mathcal{S}$ ; initialize  $return(s) \leftarrow$  an empty list for all  $s \in \mathcal{S}$ ;
2: for true do:
    # variants for this alg. can be start with  $s_0 \in \mathcal{S}, a_0 \in \mathcal{A}(s_0)$  randomly,  $\varepsilon$ -greedy, etc.
3:   Generate a complete episode  $\tau = (s_0, a_0, r_1, \dots, s_{T-1}, a_{T-1}, r_T)$ ;
4:    $G \leftarrow 0$ ;
5:   for  $t = T - 1, T - 2, \dots, 0$  do:
6:      $G \leftarrow \gamma G + r_{t+1}$ ;
7:     if  $s_t, a_t$  appears in  $(s_0, a_0, s_1, a_1, \dots, s_{t-1}, a_{t-1})$  then:
8:       Append  $G$  to  $return(s_t, a_t)$ ;
9:        $q(s_t, a_t) \leftarrow \text{average}(return(s_t, a_t))$ ;
10:       $\pi(s_t) \leftarrow \arg \max_{a'} q(s_t, a')$ ;
11:     end if
12:   end for
13: end for

```

Proof:

$$\begin{aligned}
q^\pi(s, \pi'(s)) &= \sum_{a \in \mathcal{A}} \pi'(a|s) q^\pi(s, a) \\
&= \frac{\varepsilon}{m} \sum_{a \in \mathcal{A}} q^\pi(s, a) + (1 - \varepsilon) \max_a q^\pi(s, a) \\
&\geq \frac{\varepsilon}{m} \sum_{a \in \mathcal{A}} q^\pi(s, a) + (1 - \varepsilon) \sum_{a \in \mathcal{A}} \frac{\pi(a|s) - \frac{\varepsilon}{m}}{1 - \varepsilon} q^\pi(s, a) \\
&= \sum_{a \in \mathcal{A}} \pi(a|s) q^\pi(s, a) \\
&= v^\pi(s).
\end{aligned} \tag{1}$$

□

(The above proof is given in the Chapter 5 of the book [1]. After expanding the sum, however, it can be shown the inequality should be equality. Hmm...)

For ε -greedy, when we explore, we choose actions at random without regard to their estimated values, which may waste a chance on the action that we can sure it has a bad performance. To focus more on those states with higher value, we can refer to *Boltzmann exploration*.

Boltzmann Exploration: *Boltzmann exploration also known as Gibbs sampling and soft-max, it chooses an action based on its estimated value:*

$$\pi(a|s) = \frac{e^{\beta \hat{q}(s, a)}}{\sum_{a'} e^{\beta \hat{q}(s, a')}},$$

where $\hat{q}(s, a)$ is the estimation of the value of being in state s and taking action a , β is a tunable parameter.

7.2 On Policy Temporal-Difference Control: Sarsa

As we mentioned before, Temporal-Difference(TD) learning has several advantages over Monte-Carlo (MC) such as lower variance, online, incomplete sequences. The difference between TD control and TD learning is similar to that between MC control and MC learning. What we need to do is applying TD to $Q(S, A)$ and updating the policy every time-step.

Sarsa: *The name Sarsa is inspired by the exploitation on the quintuple $s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1}$. Consider transitions from state-action pair to state-action pair, and learn the values of state-action pairs:*

$$q(s_t, a_t) \leftarrow q(s_t, a_t) + \alpha[r_{t+1} + \gamma q(s_{t+1}, a_{t+1}) - q(s_t, a_t)].$$

Algorithm 14 TD Sarsa

```

1: initialize  $\pi(s) \in \mathcal{A}(s)$ (arbitrarily), for all  $s \in \mathcal{S}$ ; initialize  $q(s, a) \in \mathcal{R}$ (arbitrarily) for all possible pairs
   in  $\mathcal{S} \times \mathcal{A}$ ; initialize  $q(\text{terminal}, \cdot) = 0$ ;
2: for true do:
   # variants for this alg. can be start with  $s_0 \in \mathcal{S}, a_0 \in \mathcal{A}(s_0)$  randomly,  $\varepsilon$ -greedy, etc.
3:   Generate a start state  $s$ ;
4:   Choose action  $a \leftarrow \pi(s)$ ;
5:   while  $s$  is not terminal do:
6:      $r, s' \leftarrow \text{environment}(s, a)$ ;
7:      $a' \leftarrow \pi(s')$ 
8:      $q(s, a) \leftarrow q(s, a) + \alpha[r + \gamma q(s', a') - q(s, a)]$ ;
9:     Update  $\pi$  according to  $q(s, a)$ ;
10:     $s \leftarrow s'$ ;
11:     $a \leftarrow a'$ ;
12:   end while
13: end for
```

There are also n -step Sarsa and Sarsa(λ) which can be derived from TD methods.

n -step Sarsa: *Define the n -step Q -return as*

$$q_t^{(n)} = r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{n-1} r_{t+n} + \gamma^n q(s_{t+n}, a_{t+n}),$$

then n -step Sarsa updates $q(s, a)$ towards the n -step Q -return

$$q(s_t, a_t) \leftarrow q(s_t, a_t) + \alpha(q_t^{(n)} - q(s_t, a_t)).$$

Like what we mentioned in n -step TD learning, when $n \rightarrow \infty$, Sarsa \rightarrow MC.

7.3 Off-Policy Temporal-Difference Control: Q-Learning

The development of Q-learning is a big breakout in the early days of Reinforcement Learning.

For the target policy, Q-learning updates it in a *greedy* way:

$$\pi(s_{t+1}) = \arg \max_{a'} q(s_{t+1}, a').$$

For the behavior policy, it will be updated in an ε -*greedy* way on $q(s, a)$.

The state-action value will be updated as

$$q(s, a) \leftarrow q(s, a) + \alpha[r + \gamma \max_{a'} q(s', a') - q(s, a)].$$

Algorithm 15 Q-Learning

```

1: initialize  $\pi(s) \in \mathcal{A}(s)$ (arbitrarily), for all  $s \in \mathcal{S}$ ; initialize  $q(s, a) \in \mathcal{R}$ (arbitrarily) for all possible pairs
   in  $\mathcal{S} \times \mathcal{A}$ ; initialize  $q(\text{terminal}, \cdot) = 0$ ;
2: for true do:
   # variants for this alg. can be start with  $s_0 \in \mathcal{S}, a_0 \in \mathcal{A}(s_0)$  randomly,  $\varepsilon$ -greedy, etc.
3:   Generate a start state  $s$ ;
4:   while  $s$  is not terminal do:
5:     Choose action  $a \leftarrow \pi(s)$ ;
6:      $r, s' \leftarrow \text{environment}(s, a)$ ;
7:      $q(s, a) \leftarrow q(s, a) + \alpha[r + \gamma \max_{a'} q(s', a') - q(s, a)]$ ;
8:     Update  $\pi$  according to  $q(s, a)$ ;
9:      $s \leftarrow s'$ ;
10:  end while
11: end for

```

One should notice that different with TD Sarsa in line 8, the action Q-learning uses to update $q(s, a)$ in line 7 is not the same as the one it truly takes.