

Notes while reading Reinforcement learning an introduction (Sutton/Barto)

Willem

January 2013

Contents

1	Introduction	1
2	Mutli-armed bandits	3
3	Finite Markov Decision Process	5
4	Dynamic Programming	7
4.1	Exercise 4.8	7
5	Monte Carlo Methods	9
5.1	Exercises	9
5.1.1	Exercise 5.1 page 94	9
5.1.2	Exercise 5.2 page 94	9
5.1.3	Exercise 5.4 page 99	9
5.1.4	Exercise 5.5 page 105	9
5.1.5	Exercise 5.6 page 108	10
5.1.6	Exercise 5.7 page 108	10
5.1.7	Exercise 5.8 page 108	10
5.1.8	Exercise 5.11 page 111	11
6	TD Prediction	13
6.1	Exercises	13
6.1.1	Exercise 6.1	13
6.1.2	Exercise 6.2	13
6.1.3	Exercise 6.3	14
6.1.4	Exercise 6.4	14
6.1.5	Exercise 6.5	14
6.1.6	Exercise 6.6	14
6.1.7	Exercise 6.7	14

Chapter 1

Introduction

Chapter 2

Mutli-armed bandits

Chapter 3

Finite Markov Decision Process

Chapter 4

Dynamic Programming

4.1 Exercise 4.8

The reward is only obtained when the capital is above 99. When the capital is at 50, there is a 50% chance you can win the game. So this obviously is the optimal policy. When you reach 51: it would be rather odd to bet the entire capital, as you don't need to risk it all to reach 100. Bigger downside, but same upside. So the best course of action is to bet with 1, see if you can grow this above 50. If you lose it, you still have a 50% chance to win by betting it all.

Chapter 5

Monte Carlo Methods

5.1 Exercises

5.1.1 Exercise 5.1 page 94

The last 2 rows in the rear means you either have 21, or 20, which means the odd's are very good you will win. (hence high value function)

The last row on the left means the dealer has an ace, so it's at an advantage to get a higher score.

The front row's are higher on the upper diagram, as there is a usable ace. Which means that if you get a bad hit that put's you over 21. It can count as 1.

5.1.2 Exercise 5.2 page 94

As this is Markov process eg. The cards drawn are not exhaustible. The odds of winning on the second time your in the same state is just as good as the first time.

5.1.3 Exercise 5.4 page 99

The "Append G to Returns (S_t, A_t) would be replaced by increasing a count and added it as running average to some table.

5.1.4 Exercise 5.5 page 105

question: Consider an MDP with a single Non-terminal state and a single action that transitions back to the nonterminal state with probability p and transitions to the terminal state with probability

$p-1$. Let the reward be 1 on all transitions, and let $\gamma = 1$. Suppose you observe one episode that lasts 10 steps, with a return of 10. What are the first-visit and every visit estimators of the value of the non-terminal state.

10 Steps means 9 towards the non-terminal, and one towards the terminal. The rewards are all-way's the same so the final cost=10.

If $\gamma = 1$ then $G = G + \gamma R_{k+1}$ in every iteration.

In case of all visit the complete horizon counts 10 times in the non-terminal state, as the 10th time we leave the non-terminal state for good and enter the terminal state. $(1+2+3+4+5+6+7+8+9+10)/10 = 55/10 = 5.5$ So the value is 5.

In case of the first-visit, we only count the first visit which has a reward of 1.

5.1.5 Exercise 5.6 page 108

question: What is the equation analogous to (5.6) for action values $Q(s, a)$ instead of state values $V(s)$, again given returns generated using b?

$Q(s, a)$ is similar to $V(s)$, it takes the $V(s)$ given a certain step was taken first.

$$Q(s, a) = \frac{\sum_{t \in J(s, a)} \rho_{t+1:T(t)-1} G_t}{\sum_{t \in J(s, a)} \rho_{t+1:T(t)-1}} \quad (5.1)$$

5.1.6 Exercise 5.7 page 108

question: In learning curves such as those shown in Figure 5.3 error generally decreases with training as indeed happened for the ordinary importance-sampling method. But for the weighted importance-sampling method error first increased and then decreased. Why do you think this happened

If there are but a few samples, the bias will be the dominating error. And it will increase as more and more samples are added. Until there are so many samples, it starts to disappear.

5.1.7 Exercise 5.8 page 108

question: The results with Example 5.5 and shown in Figure 5.4 used a first-visit MC method. Suppose that instead an every-visit MC method was used on the same problem. Would the variance

of the estimator still be infinite? Why or why not? A first Visit MC has less terms than a every Visit MC. All terms have a positive value, so it would also go to infinite.

5.1.8 Exercise 5.11 page 111

If the target policy is a greedy deterministic policy, and the loop is broken off if $\pi(S_t) \neq A_t$. Then $\pi(A_t|S_t) = 1$ by definition.

Chapter 6

TD Prediction

6.1 Exercises

6.1.1 Exercise 6.1

$$V_{t+1}(s_t) = \alpha[R_{t+1} + \gamma V_t(s_{t+1}) - V_t(s_t)] + V_t(s_t) \quad (6.1)$$

The difference between the value function at time t and $t+1$ is defined by equation 6.1.

The equality $G_t = R_{t+1} + \gamma G_{t+1}$ still holds. However the monte carlo error is slightly different in every iteration. $G_t - V_t(s_t)$ becomes $G_{t+1} - V_{t+1}(s_{t+1})$ in the next iteration. As the value function now changes at iteration t , with a difference of $d_t = \alpha[R_{t+1} + \gamma V_t(s_{t+1}) - V_t(s_t)]$.

$$G_{t+1} - V_t(s_{t+1}) = G_{t+1} - V_{t+1}(s_{t+1}) - d_{t+1} \quad (6.2)$$

$$error = - \sum_{k=t+1}^{T-1} \gamma^{k-t} d_{k-1} \quad (6.3)$$

In conclusion the different factor is equation 6.3.

6.1.2 Exercise 6.2

If (as explained in the example of the hint) a part of the statespace is already well estimated. Then the TD prediction will be very good as you enter those states and if your path ends on one of those states. So you only have lesser predictions while in an unexplored part.

The Monte Carlo approach would still need to evaluate through the already well estimated part. Which is rather slow.

6.1.3 Exercise 6.3

The change on a value function is defined by: $\alpha[R_{t+1} + \gamma V_t(s_{t+1}) - V_t(s_t)] = 0.1[0 + 0 - 0.5] = -0.05$ if $V_t(s_{t+1}) = 0$ so it ends on the left terminal state. And $\alpha = 0.1$ and $V_t(A) = 0.5$.

6.1.4 Exercise 6.4

The TD algo is over-fitting when $\alpha > 0.05$ we could try to make it a bit smaller. But at $\alpha = 0.05$ it seems to flatten out nicely, so I would not expect better results.

A similar story with the MC method, this time at $\alpha 0.02$ we get a nice flat tail. It's not as clear as with the TD method, but that's due the larger variance on the MC method.

So no, I would not expect any changes in results if more samples were ran with different values for α .

6.1.5 Exercise 6.5

Overfitting, the step is too large so TD cannot find the optimal values. But keeps over/under estimating every time it runs through an episode.

6.1.6 Exercise 6.6

You setup the bellman optimality equation, and the pick a method to solve it. As this is a rather simple example, you could just manually solve the equation.

$$\begin{aligned}
 V(A) &= 0.5V(B) \\
 V(B) &= 0.5V(A) + 0.5V(C) \\
 V(C) &= 0.5V(B) + 0.5V(D) \\
 V(D) &= 0.5V(C) + 0.5V(E) \\
 V(E) &= 0.5V(D) + 0.5
 \end{aligned} \tag{6.4}$$

This seems like the simplest way to do it, as it's small.

6.1.7 Exercise 6.7

The normal on-policy TD(0) update looks like $V(s_t) = V(s_t) + \alpha[R_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$. I would expect that $\alpha = \frac{\rho}{\sum_t \rho_t}$ as it becomes a weighted average due too the importance sampling.