

RL Assignment 2

Vishaal Udandarao
20/6/19

- 1) We know all the states as given in set $S \Rightarrow$
 $S = \{ \text{high, low} \}$

We also know the action sets for each of the individual states \Rightarrow

$$A(\text{high}) = \{ \text{search, wait} \}$$

$$A(\text{low}) = \{ \text{search, wait, recharge} \}$$

Now, out of all the 4-tuples, we only need to consider those for which $p(s', r | s, a) > 0$.

Therefore, we notice that for the following 4-tuples, $p(s', r | s, a) = 0$ as that combination of a 4-tuple cannot occur \Rightarrow

i) $p(s' = \text{low}, r | s = \text{high}, a = \text{wait}) = 0$

ii) $p(s' = \text{high}, r | s = \text{low}, a = \text{wait}) = 0$

iii) $p(s' = \text{low}, r | s = \text{low}, a = \text{recharge}) = 0$.

Therefore, for all other combinations, we know that $p(s', r | s, a) > 0$ and hence we can write the table as follows \Rightarrow

s	a	s'	r	$p(s', r s, a)$
high	search	high	r_{search}	α
high	search	low	r_{search}	$1 - \alpha$
low	search	high	-3	$1 - \beta$
low	search	low	r_{search}	β
high	wait	high	r_{wait}	1
low	wait	low	r_{wait}	1
low	recharge	high	0	1

3) Exercise 3.15

The signs of the reward are ~~not~~ not important at an absolute scale but are important relatively. That is, we can change the value of the reward by adding a large positive constant value to all rewards. This would not affect the value function. However, if the signs of the rewards are inverted, then the value function will change because then negative rewards which ~~are~~ were positive rewards initially would mean that in the particular state, that action will not be preferred, when in fact it should be preferred. Therefore, the signs ~~are~~ of the rewards will only matter if they are not inverted for the value-function computation.

We know, state-value function \Rightarrow

$$V_{\pi}(s) = E_{\pi} [G_t | S_t = s]$$

$$= E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right]$$

Now, if add a constant c to every reward \Rightarrow

$$V_{\pi}^1(s) = E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k (R_{t+k+1} + c) | S_t = s \right]$$

$$= E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right] +$$

$$E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k c | S_t = s \right]$$

$$= E_{\pi} [G_t | S_t = s] + E_{\pi} \left[c \sum_{k=0}^{\infty} \gamma^k | S_t = s \right]$$

Now, the second term is clearly a constant and hence can be written as \Rightarrow

$$V_c = E_{\pi} \left[c \sum_{k=0}^{\infty} \gamma^k | S_t = s \right]$$

$$= c \sum_{k=0}^{\infty} \gamma^k$$

Now, we know $0 \leq \gamma \leq 1 \Rightarrow \sum_{k=0}^{\infty} \gamma^k = \frac{1}{1-\gamma}$

$$\Rightarrow V_c = c \left(\frac{1}{1-\gamma} \right)$$

Therefore, we can write $V'_\pi(s) \Rightarrow$

$$V'_\pi(s) = E_\pi[G_t | S_t = s] + V_c$$

$$\boxed{V'_\pi(s) = V_\pi(s) + V_c}$$

Therefore, $V'_\pi(s)$ is a sum of $V_\pi(s)$ and a constant V_c . Thus, the relative values of all states under any policy remain unaffected.

Exercise 3.16

For the case of an episodic task, we know state-value function as \Rightarrow

~~$$V_\pi(s) = E_\pi[G_t | S_t = s]$$~~

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

$$= E_\pi \left[\sum_{k=0}^{T-1} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

Here, T denotes no. of steps to terminal state.

Following a similar procedure as in Exercise 3.15 \Rightarrow

on adding constant c to all rewards \Rightarrow

$$V'_\pi(s) = E_\pi[G_t | S_t = s] + E_\pi \left[c \sum_{k=0}^{T-1} \gamma^k \mid S_t = s \right]$$

If we write it in the form as in exercise 3.15 \Rightarrow

$$V'_\pi(s) = V_\pi(s) + V_c \text{ where } \Rightarrow$$

$$V_c = E_\pi \left[c \sum_{k=0}^{T-1} \gamma^k \mid S_t = s \right]$$

$$= c \sum_{k=0}^{T-1} \gamma^k$$

$$= c \left[\frac{1 - \gamma^{T+1}}{1 - \gamma} \right] \quad (\text{sum of a G.P.})$$

$$\therefore V_c = C \left[\frac{1 - \gamma^{T+1}}{1 - \gamma} \right]$$

We clearly see that V_c is stochastic as it is a random variable dependent on T .

However, for one particular episode, T remains the same and hence, relatively for all the states, the term V_c is a constant since T is a constant for all states, in a given episode.

Therefore, within a given episode, the state-value functions will remain unaffected relatively.

5) We know,

$$V_*(s) = \max_{\pi} V_{\pi}(s)$$

~~$$V_*(s) = \max_{a \in A(s)} V_{\pi_*(s,a)}(s)$$~~

$$V_*(s) = \max_{a \in A(s)} q_{\pi_*(s,a)}(s,a)$$

That is $\Rightarrow V_*(s)$ is the maximum value of optimal action value functions over all actions.

$$\Rightarrow V_*(s) = \max_{a \in A(s)} q_{\pi_*(s,a)}(s,a)$$

$$\Rightarrow V_*(s) = \max_{a \in A(s)} E[G_t | S_t = s, A_t = a]$$

$$\Rightarrow V_*(s) = \max_{a \in A(s)} \sum_{s', r} p(s', r | s, a) [\gamma + \gamma V_*(s')]$$

Now, we also know \Rightarrow

$$V_*(s') = \max_{a' \in A(s')} q_{\pi_*}(s', a')$$

$$\Rightarrow V_*(s) = \max_{a \in A(s)} \sum_{s', r} p(s', r | s, a) [\gamma + \gamma \max_{a' \in A(s')} q_{\pi_*}(s', a')]$$