

Homework 2

10 5 10

Ans 3

Ex 3.15 →

9 10 15 20

- (a) No, sign of the rewards doesn't matter, only the intervals between them matter. This is because when we add a constant R^c to all the possible rewards, ~~at~~ another constant V^c is added to all $V_\pi(s)$ $s \in S$. Thus relative ^{Expected} reward among states remain same. We can prove this in the following way using just Bellman equation → (b)

$$V_\pi(s) = \sum_a \pi(a|s) \sum_{s', r} P(s', r|s, a) (r + \gamma V_\pi(s'))$$

now we add constant c to r

$$\Rightarrow V_\pi(s) = \sum_a \pi(a|s) \sum_{s', r} P(s', r|s, a) (r + c + \gamma V_\pi(s'))$$

$$\Rightarrow V_\pi(s) = \sum_a \pi(a|s) \left(\sum_{s', r} P(s', r|s, a) (r + \gamma V_\pi(s')) + \sum_{s', r} P(s', r|s, a) x c \right)$$

$$= \sum_a \pi(a|s) \sum_{s', r} P(s', r|s, a) (r + \gamma V_\pi(s')) + c$$

$$\sum_a \pi(a|s) x \left(\sum_a \sum_{s'} \sum_r \pi(a|s) P(s', r|s, a) \right) = 1$$

- (b) $G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$
adding constant c to reward.

$$G_t = \sum_{k=0}^{\infty} \gamma^k (R_{t+k+1} + c)$$

$$= \sum_{k=0}^{\infty} R_{t+k+1} + c \sum_{k=0}^{\infty} \gamma^k$$

$$\leftarrow \rightarrow \Sigma$$

$$G_t = \sum_{k=0}^{\infty} \gamma^k (R_{t+k+1}) + \frac{C}{1-\gamma}$$

$$V(s) = E[G_t | s_t = s]$$

$$= E \left[\sum_{k=0}^{\infty} \gamma^k (R_{t+k+1}) + \frac{C}{1-\gamma} \mid s_t = s \right]$$

$$V(s) = \frac{C}{1-\gamma} + E \left[\sum_{k=0}^{\infty} \gamma^k (R_{t+k+1}) \mid s_t = s \right]$$

$$\therefore V_C = \frac{C}{1-\gamma}$$

We can clearly see that adding constant C to all the rewards adds $\frac{C}{1-\gamma}$ to

$V(s)$. H.O.P

Ex 13.16

G_t would ~~def~~ have an effect in the case of episodic task as in that case if the length of the episode is T

then, $\sum_{k=0}^{T-1} \gamma^k$ would be added, that is

$$\text{equal to } \frac{C(1-\gamma^T)}{1-\gamma}$$

thus ~~the~~ reward would change more as we progress towards the end of an episode.

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

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

$$Q_*(s, a) = \max_{\pi} Q_{\pi}(s, a).$$

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

Q1

s	a	s'	γ	$p(s' \gamma s,a)$
high	search	low	0	$p_1 = 1 - \alpha - \gamma_{\text{search}}$
"	"	"	1	$p_2 = \gamma_{\text{search}}$
"	"	"	-3	0
"	"	high	0	$p_3 = \alpha - \gamma_{\text{search}}$
"	"	"	1	$p_4 = \gamma_{\text{search}}$
"	"	"	-3	0
low	"	low	0	$p_5 = \beta - \gamma_{\text{search}}$
"	"	"	1	$p_6 = \gamma_{\text{search}}$
"	"	"	-3	0
"	"	high	0	0
"	"	"	1	0
"	"	"	-3	0 $1 - \beta$
high	wait	low	0	0
"	"	"	1	0
"	"	"	-3	0
"	"	high	0	$p_7 = 1 - \gamma_{\text{wait}}$
"	"	"	1	$p_{10} = \gamma_{\text{wait}}$
"	"	"	-3	0
low	"	low	0	$p_{11} = 1 - \gamma_{\text{wait}}$
"	"	"	1	$p_{12} = \gamma_{\text{wait}}$
"	"	"	-3	0
"	"	high	0	0
"	"	"	1	0
"	"	"	-3	0
low	exchange	high	0	1
"	"	"	1	0
"	"	"	-3	0
"	"	low	0	0
"	"	"	1	0
"	"	"	-3	0

Now, we calculate p_i with the help of given ~~equation~~ information

$$p_1 + p_2 = 1 - \alpha, \quad p_1 \times 0 + p_2 \times 1 = \gamma_{\text{search}}$$

$$\Rightarrow p_2 = \gamma_{\text{search}} \quad p_1 = 1 - \alpha - \gamma_{\text{search}}$$

$$p_3 + p_4 = \alpha, \quad p_3 \times 0 + p_4 \times 1 = \gamma_{\text{search}}$$

$$\Rightarrow p_4 = \gamma_{\text{search}}, \quad p_3 = \alpha - \gamma_{\text{search}}$$

$$p_5 + p_6 = \beta, \quad p_5 \times 0 + p_6 \times 1 = \gamma_{\text{search}}$$

$$\Rightarrow p_6 = \gamma_{\text{search}} \quad p_5 = \cancel{\gamma_{\text{search}}} \beta \quad \beta - \gamma_{\text{search}}$$

$$\cancel{p_7 + p_8 = 1 - \beta}$$

$$p_9 + p_{10} = 1, \quad p_9 \times 0 + p_{10} \times 1 = \gamma_{\text{wait}}$$

$$\Rightarrow p_{10} = \gamma_{\text{wait}} \quad p_9 = 1 - \gamma_{\text{wait}}$$

$$p_{11} + p_{12} = 1, \quad p_{11} \times 0 + p_{12} \times 1 = \gamma_{\text{wait}}$$

$$\Rightarrow p_{12} = \gamma_{\text{wait}} \quad p_{11} = 1 - \gamma_{\text{wait}}$$