## Exercise 3.15

6.

$$G = R_1 + \gamma R_2 + \gamma^2 R_3 + \cdots + \gamma^T R_T.$$

$$G' = (R_1 + c)^\gamma + (R_2 + c) + \cdots \rightarrow \gamma^T (R_T + c)$$

$$= (R_1 + \gamma R_2 + \cdots \gamma^T R_T) + c + \gamma c + \gamma^2 c$$
$$+ \cdots + \gamma^T c.$$

$$= G + \frac{c \cdot (1 - \gamma^T)}{1 - \gamma} \quad \text{as } T \to \infty$$

$$= G + \frac{c}{1 - \gamma} \quad \text{hence it is the doesn't affect}$$

the relative ordering.

$$V_c > \frac{c}{1 - \gamma}.$$

Exercise 3.16

Continuing the previous question, we got

$$G' = G + \frac{c(1-r^T)}{1-r}$$

Now how long the episode will last will

also play a role in comparing.

The longer the episode last, the more

reward we will get, as $(1-r^T)$ ~~decreases~~

increases with increase in $T$ as $0 < r < 1$.

$$G^{\Phi'}_{T=1} = G + \frac{c(1-r)}{1-r}$$

$$= G + c$$

$$G'_{T=2} = G + \frac{c \cdot (1-r^2)}{1-r} = G + c(1+r)$$

$$G_{T=2} > G_{T=1}.$$

**Q5 Ans:**

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

P1. Since $r(s, a, s')$ is an expectation we need to figure out the ~~order sentered~~ ~~for each~~ probability of finding a can or ~~not~~ finding a can.

$p \leftarrow$ find a can $q \leftarrow$

$r_{search} = $ ~~find~~ $p \times 1 + (1-p) \times 0$

$$= p.$$

∴ probability of finding a can is $\beta \cdot r_{search}$
~~not~~ " " not " " " " $1 - r_{search}$

Given $p(s', r \mid s, a) = \alpha \times r_{search}$.

~~$p(s', r \mid s = find-can, a) = \alpha\,r_{search}$~~
~~$p(s', r \mid s = not\,find\,ca$~~

$p(s \mid s, a) = \alpha$ So,

~~$p(s', r \to 0 \mid s, a)$~~ ~~$+$~~ $p(s', r \to 1 \mid s, a) \to p(s' \mid s, a)$

$\geq$ ~~$p(s', r \to 0 \mid s, a)$~~ $= \alpha - \alpha\,r_{search}$.

Doing this for others.

| s | a | s' | r | $P(s', r \mid s, a)$ |
|---|---|----|---|---------------------|
| h | s | h | 1 | $\alpha\, r_{search}$ |
| h | s | h | 0 | $\alpha - \alpha\, r_{search}$ |
| h | s | l | 1 | $(1-\alpha)\, r_{search}$ |
| l | s | h | -3 | $1 - \beta$ |
| l | s | l | 1 | $\beta \times r_{search}$ |
| l | s | l | 0 | $\beta - \beta \times r_{search}$ |
| h | s | l | 1 | $(1-\alpha) - (1-\alpha)\, r_{search}$ |
| h | w | h | 1 | $r_{wait}$ |
| l | w | l | 1 | $r_{wait}$ |
| l | w | l | 0 | $1 - r_{wait}$ |
| l | r | h | 0 | $1$ |
| h | w | h | 0 | $1 - r_{wait}$ |