Ex. 3.15

We know that, for an policy $\pi$ & state $s$,

$$v_\pi(s) = \mathbb{E}\left[ G_t \mid S_t = s \right]$$

$$\Rightarrow v_\pi(s) = \mathbb{E}\left[ R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \ldots \mid S_t = s \right]$$

using eq 3.8 in the chapter, ie,

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \ldots$$

Adding constant 'c' to all the rewards,

$$\Rightarrow v'_\pi(s) = \mathbb{E}\left[ (R_{t+1} + c) + \gamma(R_{t+2} + c) + \gamma^2(R_{t+3} + c) + \ldots \mid S_t = s \right]$$

$$\Rightarrow v'_\pi(s) = \mathbb{E}\left[ R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \ldots \mid S_t = s \right] + \mathbb{E}\left[ c\gamma^0 + c\gamma + c\gamma^2 + \ldots \mid S_t = s \right]$$

$$\Rightarrow v'_\pi(s) = v_\pi(s) + c\,\mathbb{E}\left[ \underbrace{\sum_{k=0}^{\infty} \gamma^k}_{constant} \mid S_t = s \right]$$

$$\Rightarrow v'_\pi(s) = v_\pi(s) + \frac{c}{1-\gamma} \qquad \left( \begin{array}{l} \text{since, } \gamma < 0 \text{ & it is an} \\ \text{infinite sum of G.P.} \end{array} \right)$$

Hence, by adding 'c' to all the rewards changes the value of of each state by $\frac{c}{1-\gamma}$. ~~That~~ Thus,

$$V_c = \frac{c}{1-\gamma}$$

## Ex. 3.16

Given : an episodic Task & adding $c$ to all the rewards

$$v_\pi(s) = \mathbb{E}\left[ G_t \mid S_t = s \right]$$

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

Now, adding $c$ to all the rewards,

$$v'_\pi(s) = \mathbb{E}\left[ \sum_{k=0}^{T-t-1} \gamma^k (R_{t+1+k} + c) \mid S_t = s \right]$$

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

$$= \mathbb{E}\left[ G_t \mid S_t = s \right] + c \left( \frac{\gamma^{T-t} - 1}{\gamma - 1} \right)$$

$$v'_\pi(s) = v_\pi(s) + \frac{c(1 - \gamma^{T-t})}{1 - \gamma}$$

For a given $t$, the increase is smaller, for a smaller value of $(T-t)$.

So states that are at a shorter distance to terminate will end up having a relatively smaller change.

# Ex. 3.4

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

The above table has been found out using the below
formulas & the given conditions : ~~such as~~

$$p(s' \mid s, a) = \sum_{r \in R} p(s', r \mid s, a)$$

$$r(s, a, s') = \sum_{r \in R} r \; \frac{p(s', r \mid s, a)}{p(s' \mid s, a)}$$

Some of the conditions are :

→ If the energy level is high, then a period of active search can always be completed without risk of depleting the battery.

→ A reward of −3 results whenever the robot has to be rescued.
→ No cans can be collected during a run home for recharging.

~~some~~ 1st row has been explained below :

$$P(high, 1 \mid high, search) + P(high, 0 \mid high, search) = P(high \mid high, search)$$

$$= \alpha$$

& $r(high, search, high) = r_{search}$

$$r_{search} = 1 \cdot \frac{P(high, 1 \mid high, search)}{P(high \mid high, search)} + 0 \cdot \frac{P(high, 0 \mid high, search)}{P(high \mid high, search)}$$

$$\Rightarrow P(high, 1 \mid high, search) = \alpha \, r_{search}$$

$$\Rightarrow P(high, 0 \mid high, search) = \alpha - \alpha \, r_{search}$$

8.5. $v_*(s) = \max_a q_*(s,a)$

$s \mapsto v_*(s)$

$a \mapsto a_*(s, a)$