# RL Homework 2

**Submitted By:**
**Snehal Gupta**
**2016201**

**Q1**

**Q1.**

given:

$S = \{ low, high \}$

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

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

search at high — leaves high with probability $\alpha$
— reduces it to low with probability $1-\alpha$

search at low — leaves low with probability $\beta$
— depletes the battery with $1-\beta$ [Robot is rescued and battery is recharged to high]

Reward — 1 for each can collected
— $-3$ when got rescued

$R_{search}$ — Expected no. of cans the robot will collect after searching.

$R_{wait}$ — Expected no. of cans the robot will collect while waiting.

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

This has been obtained using

$$p(s',r\mid s,a) = P\left[ S_{t+1}=s',\ R_{t+1}=r \mid S_t=s,\ A_t=a \right]$$

**Q2**

According to Bellman equation for $V_\pi$,

$$V_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) \left[ r + \gamma V_\pi(s') \right]$$

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

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

$$V_\pi(s) - \sum_a \pi(a|s) \sum_{s',r} \gamma\, p(s',r|s,a)\, V_\pi(s')$$

$$= \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)\, r$$

For $s = s'$

$$\text{Coeff of } V_\pi(s) = \left(1 - \sum_a \pi(a|s) \sum_{s',r} \gamma\, p(s',r|s,a)\right)$$

For $s \neq s'$

$$\text{Coeff of } V_\pi(s) = \sum_{s', s \neq s'} \left( \sum_a \pi(a|s) \sum_{s',r} \gamma\, p(s',r|s,a) \right)$$

$$\text{Coeff of } V_\pi(s) \times V_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)\, r$$

# Q3

## Q3

### Ex 3.15

Signs of the rewards are not important and only the intervals between them are important. This is because, rewards can be made of the same sign by adding/subtracting a large positive constant c from all the rewards. This leads to increase/decrease in the value function by a constant, which does not affect the algorithm.

To Prove: Adding a constant $c$ to all the rewards adds a constant $V_c$ to the value of the states, and, thus does not affect the relative values of any states under any policies.

#### Proof:

We know that,

$$G_t = \sum_{k=0}^{\infty} \varphi^k R_{t+k+1} \quad \text{and} \quad V_\pi(s) = E_\pi[G_t \mid S_t = s, A_t = a]$$

After adding constant $c$ to all the rewards,

$$G_t' = \sum_{k=0}^{\infty} \varphi^k (R_{t+k+1} + c)$$

---

$$V_\pi'(s) = E[G_t' \mid S_t = s]$$

$$= E\left[\sum_{k=0}^{\infty} \varphi^k (R_{t+k+1} + c) \mid S_t = s\right]$$

$$= E\left[\sum_{k=0}^{\infty} \varphi^k R_{t+k+1} + \sum_{k=0}^{\infty} \varphi^k c \mid S_t = s\right]$$

$$= E\left[\sum_{k=0}^{\infty} \varphi^k R_{t+k+1} \mid S_t = s\right] +$$

$$\underbrace{E\left[\sum_{k=0}^{\infty} \varphi^k c \mid S_t = s\right]}_{\text{constant term}}$$

$$= V_\pi(s) + \sum_{k=0}^{\infty} \varphi^k c$$

$$V_\pi'(s) = V_\pi(s) + \underbrace{\frac{c}{1-\varphi}}_{\text{constant term } V_c} \quad [\text{since } 0 \leq \varphi < 1]$$

We can observe that the value function increases only by a constant $V_c$ and hence, does not affect the relative values of any states under any policies.

---

### Ex 3.16

In case of episodic task,
Let terminal time be T.

Solving the equation,

$$V_\pi'(s) = E[G_t' \mid S_t = s]$$

$$= E\left[\sum_{k=0}^{T} \varphi_k (R_{t+k+1} + c) \mid S_t = s\right]$$

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

$$E\left[\sum_{k=0}^{T} \varphi_k c \mid S_t = s\right]$$

$$= V_\pi(s) + E\left[\sum_{k=0}^{T} \varphi_k c \mid S_t = s\right]$$

Here, $V_c = E\left[\sum_{k=0}^{T} \varphi_k c \mid S_t = s\right]$ is a function of T and T is a random variable that normally varies from episode to episode. Different episodes will have different value-functions.

Now $G_t$ is $G_t + c\left(\dfrac{1 - \varphi^T}{1-\varphi}\right)$

$\Rightarrow$ It will increase $V_\pi$ when T increases.

---

### Example

Consider an episodic task with one state $S$ and two actions $A_1$ and $A_2$.

$A_1 \rightarrow$ Agent goes to terminal state with reward 1
$A_2 \rightarrow$ Agent goes back to $S$ with reward 0.

On adding 1 to each reward,
when $A_2$ is performed forever, return is $\dfrac{1}{1-\varphi}$
which can be bigger than 2 if $\varphi < \dfrac{1}{2}$.

**Q5**

$$V^*(s) = \max_{a \in A(s)} q_{\pi^*}(s,a)$$

$$= \max_a E_{\pi^*}\left[G_t \mid S_t = s, A_t = a\right]$$

$$= \max_a E_{\pi^*}\left[R_{t+1} + \varphi G_{t+1} \mid S_t = s, A_t = a\right]$$

$$= \max_a E\left[R_{t+1} + \varphi V^*(S_{t+1}) \mid S_t = s, A_t = a\right] \quad \text{—①}$$

Also,

$$q^*(s,a) = E\left[R_{t+1} + \varphi V^*(S_{t+1}) \mid S_t = s, A_t = a\right] \quad \text{—②}$$

Substituting ② into ①

$$V^*(s) = \max_a q^*(s,a)$$