# Question 3 : Exercise 3.15

In the given gridworld example, the signs aren't important but their magnitude difference is what matters.

We need to have higher values for goals as compared to edges and rest so that the agent learns to favor this goal state in order to maximize the overall reward.

Similarily, running to the edge should have the least value and stationary should be somewhere in between.

$$V_\pi(S) = E_\pi [q_t | S_t = S]$$

$$= E_\pi [R_{t+1} + VR_{t+2} + V^2 R_{t+3} \cdots | S_t = S]$$

Adding a constant 'c' to every reward gives

$$V_\pi^{new}(S) = E_\pi\left[R_{t+1}+C+V(R_{t+2}+C)+V^2(R_{t+3}+C)+\ldots\infty \Big| S_t=S\right]$$

$$= E_\pi\left[(R_{t+1}+VR_{t+2}+\ldots+C+VC+\ldots)\Big|S_t=S\right]$$

$$= E_\pi\left[(R_{t+1}+VR_{t+2}+V^2R_{t+3}+\ldots)\Big|S_t=S\right]$$

$$+ E_\pi\left[C+VC+V^2C+\ldots\Big|S_t=S\right]$$

(Using property of expectation)

$$= E_\pi\left[q_t\Big|S_t=S\right] + E_\pi\left[C(1+V+V^2+\ldots)\Big|S_t=S\right]$$

$$= V_\pi(S) + E_\pi\left[C(1+V+V^2+\ldots)\Big|S_t=S\right]$$

We can know that 'c' and gamas are constant, so we can take it out of the expected setl.

$$V_\pi^{new}(S) = V_\pi(S) + C\times\frac{1}{(1-V)} \quad \left[\begin{array}{c}\text{Sum to} \\ \text{infinity G.P}\end{array}\right]$$

$$\boxed{V_\pi^{new}(S) = V_\pi(S) + V_C} \quad \boxed{V_C = \frac{C}{1-V}}$$

Hence, proved

## Question 3: Exercise 3.16

Solving like the previous question, we get

$$V_A^{new}(S) = E_A\left[R_{J+1} + VR_{J+2} + \dots V^k R_T \mid S_J = S\right]$$

$$= E_A +$$

$$V_A^{new}(S) = E_A\left[R_{J+1} + VR_{J+2} \dots V^k R_T \mid S_J = S\right]$$

$$+ E_A\left[c + Vc + \dots V^k c \mid S_J = S\right]$$

$$V_A^{new}(S) = \mu_A(S) + E_A\left[c + Vc + \dots V^k c \mid S_J = S\right]$$

Here the sum is not till infinite but only till T.

$$E_A\left[c + Vc + \dots V^k \mid S_J = S\right] = c\left(1 + V + V^2 \dots V^k\right)$$

There are $(T - J - 1)$ terms.

$$\cdot c\left(\frac{1 - V^{T-J}}{1 - V}\right)$$

$$\therefore \mu_A^{new}(S) = \mu_A(S) + c\left(\frac{1 - V^{T-J}}{1 - V}\right)$$

Here, we can see that $\mu_c$ depends on a random variable T.

However, this will not affect the $U_s$ in one particular episode because $T$ is constant for that episode.

$$\therefore U_c = C\left(\frac{1-V^{T-t+1}}{1-V}\right)$$

## Question 5

We know that

$$U^*(s) = \max_a q^*(s,a) \qquad ①$$

At any state, the optimal value function will be the optimal action value function.

$$v^*(s) = \max_a \sum_{s'} p(s'|s,a) \cdot \left[ E[r|s,a,s'] + V v^*(s') \right]$$

$$x_1 + x_2 + x_3 = 0$$

$\Rightarrow$ But From ①, we have

$$v^*(s) = \max_a \sum_{s'} p(s'|s,a) \left[ E[r|s,a,s'] \right]$$

$$+ V \max_{a'} q^*_p(s',a') \right]$$

## Question 1: 3.4

In the table, value of $p(s'|s,a)$ is mentioned. We need to find the value of $p(s',r|s,a)$.

We know that
$$p(r,s'|a,s) = p(s'|s,a) \times p(r|s',s,a)$$

$$E[r|s,a,s'] = \sum_r r \times p(r|s',a,s')$$

This can be done as follow, when $s =$ low, and $a =$ search, and next state $s' = 1$, we can have two rewards, 1 and 0, so we will have to take expectation over both of these.

Scanned with CamScanner

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