

Q1) x_{search} and x_{wait} are expectations of the reward while searching and waiting.
Therefore consider a probability distribution P_{search} that has expectation x_{search} and P_{wait} that has expectation x_{wait}

s	a	s'	x	$P(s', x s, a)$
1) high	search	high	x_1	$\alpha \cdot P_{search}(x_1)$
2) high	search	low	x_2	$(1-\alpha) P_{search}(x_2)$
3) low	search	high	-3	$(1-\beta) \cdot 1$
4) low	search	low	x_3	$\beta P_{search}(x_3)$
5) high	wait	high	x_4	$1 \cdot P_{wait}(x_4)$
6) low	wait	low	x_5	$1 \cdot P_{wait}(x_5)$

~~low~~ ~~average~~

For the rest of the rows such as, high wait low gives $P(s', x | s, a)$ as 0 because if you wait in high state, you cannot go to low state.
Similarly for,

classmate

low cost high $\rightarrow 0$ as transition probability is 0.

low recharge high $\rightarrow 0$ because reward is 0.

Explanation of table

Row 1 \rightarrow At state high, you take an action search but return to the same state.

The probability of this is α . The reward has expectation r_{search} hence we have π_1 that belongs to the distribution P_{search} with expectation r_{search} . $\therefore P(s', \pi | s, a)$ becomes transition probability times reward probability.

Row 2 \rightarrow π_2 belongs to the distribution of P_{search} and now from high cost move to low with probability $1-\alpha$. Hence $(1-\alpha) \cdot P(\pi_2)$

To check our answer for row 1 and row 2, we can sum $P(s', \pi | s, a)$ and it should be equal to 1.

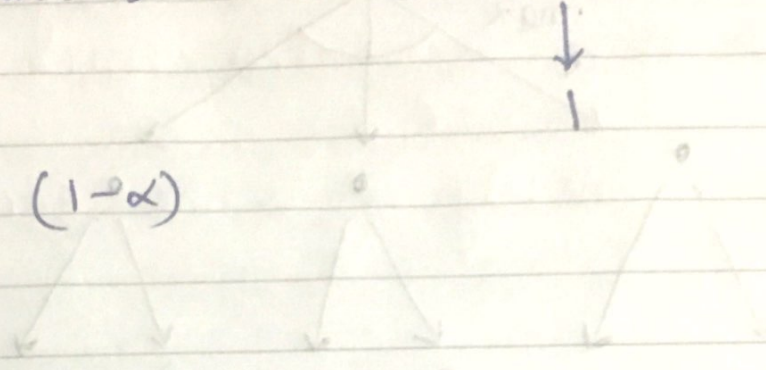
$$\sum \alpha P_{\text{search}}(\pi_1) + \sum (1-\alpha) P_{\text{search}}(\pi_2)$$

classmate

$$= \alpha \sum P_{\text{search}}(s_1) + (1-\alpha) \sum P_{\text{search}}(s_2)$$

$$= \alpha + (1-\alpha)$$

$$= 1$$



Row 3 \rightarrow At state low, action search and search high, gives a reward of -3 with probability 1, ~~here~~ and $1-\beta$ is the probability of going to state high, Hence $(1-\beta)1$

Row 4 $\rightarrow \beta \cdot P_{\text{search}}(s_3)$. stay at state low and get reward s_3 from P_{search} distribution.

summation of $(1-\beta) \cdot 1 + \beta \cdot \sum P_{\text{search}}(s_3) = 1$

Row 5 \rightarrow state is high, action is wait, state is high, reward is from distribution with expectation s_{wait} . Hence $1 \cdot P_{\text{wait}}(s_4)$

Row 6 \rightarrow Similarly as above but with starting state low, wait then come back to low with probability 1 and receive

classmate

reward 25 with probability $P_{wait}(25)$.
Hence $1 - P_{wait}(25)$.

Q3) Adding a constant c to all rewards, in case of continuing task

$$V_{\pi}(s) = E_{\pi} [G_t | s_t = s] \\ = E_{\pi} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | s_t = s]$$

Adding c to all rewards.

$$= E_{\pi} [(R_{t+1} + c) + \gamma(R_{t+2} + c) + \gamma^2(R_{t+3} + c) + \dots | s] \\ = E_{\pi} [(c + \gamma c + \gamma^2 c + \dots) + (R_{t+1} + \gamma R_{t+2} + \dots) | s_t = s] \\ = E_{\pi} \left[\frac{c}{1-\gamma} + G_t | s_t = s \right]$$

Expectation of constant is constant

$$V'_{\pi}(s) = \frac{c}{1-\gamma} + E_{\pi} [G_t | s_t = s]$$

$$V'_{\pi}(s) = \frac{c}{1-\gamma} + V_{\pi}(s)$$

Therefore all the values are just incremented by a factor of $\frac{c}{1-r}$ and hence ~~choosing the~~
~~actions~~ relatively it is the same thing.

In case of episodic task,

$$V_{\pi}(s) = E_{\pi} \left[\underbrace{R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{k-1} R_{t+k}}_{G_t} \mid s_t = s \right]$$

Adding c

$$V'_{\pi}(s) = E_{\pi} \left[(R_{t+1} + c) + \gamma(R_{t+2} + c) + \dots + \gamma^{k-1}(R_{t+k} + c) \mid s_t = s \right]$$

$$V'_{\pi}(s) = E_{\pi} \left[(c + r + \dots + r^{k-1} c) + G_t \mid s_t = s \right]$$

$$V'_{\pi}(s) = E_{\pi} \left[\frac{c(1-r^k)}{1-r} + G_t \mid s_t = s \right]$$

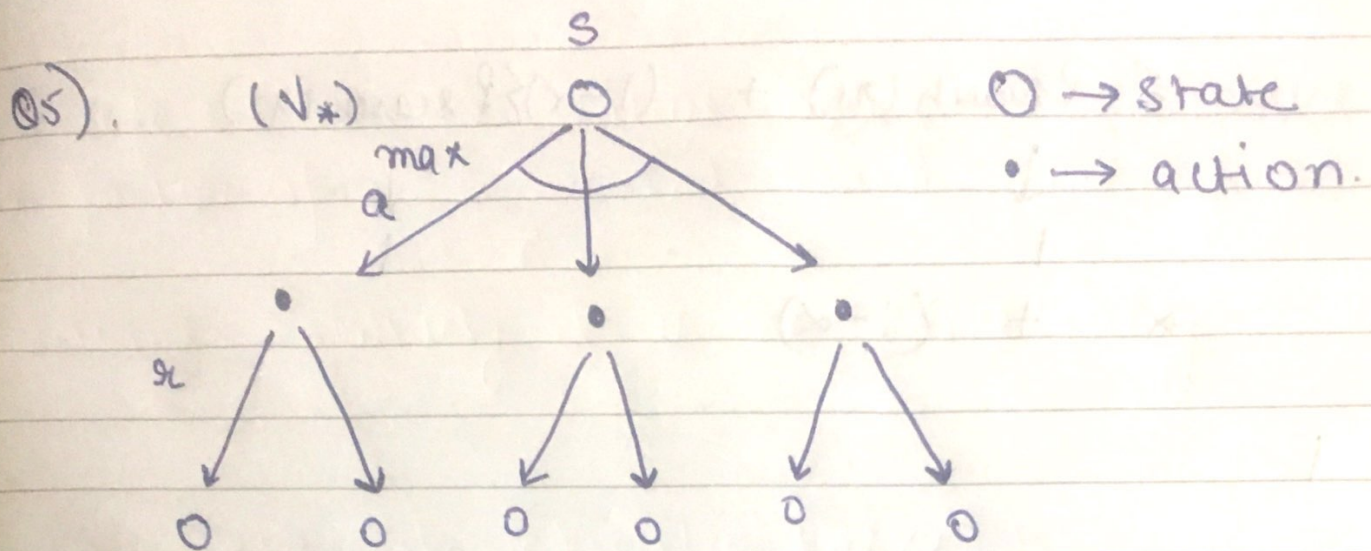
we get this additional term of $\frac{c(1-r^k)}{1-r}$

that is a random variable because k is random
 k is the time step at which terminating
 state is reached. Hence this depends on
 classmate

policy we choose and the state we begin in. If k is more then this means we started from a state that is far from the terminating state and hence $1 - r^k$ is ~~less~~ more ($0 < r < 1$) and therefore $v_{\pi}(s)$ gets a ~~smaller~~ larger constant.

If k is less, $1 - r^k$ is less and hence $v_{\pi}(s)$ gets a smaller constant.

So states that are closer will relatively have smaller value.



When you are at state s , choose an action a . You reach the • position. From there on choose the optimal ~~set~~ action value function i.e. $q_*(s, a)$. For each action

you took from state s , you will have an optimal further direction given by $q_*(s, a)$ $\forall a \in A(s)$.

Out of these, choose the best action and this gives the value v^*

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