

Dis 5)

$$\sum_{a \in A(s_t)} \pi(a | s_t) q_{\pi}(s_t, a) = V_{\pi}(s_t)$$

Now When $\pi = \pi^*$ optimal

then $V_{\pi} = V_*$, $q_{\pi} = q_*$

and π^* assigns non-zero probabilities only to those ~~actions~~ actions ($a' \in A(s_t)$) which ~~assign~~ assign the maximum value to $q_*(s_t, a)$

\therefore we can write this as:-

$$V_*(s_t) = \sum_{a \in A(s_t)} \pi^*(a | s_t) q_*(s_t, a)$$

$$= \sum_{a \in A(s_t)} \pi^*(a | s_t) \max(q_*(s_t, a))$$

const. so take out of summation

$$= \max_{a \in A(s_t)} (q_*(s_t, a)) \sum_{a \in A(s_t)} \pi^*(a | s_t)$$

$$\sum_{a \in A(s_t)} \pi^*(a | s_t) = 1$$

here 'a' is that value of 'a' that maximizes $q_*(s_t, a)$ and that max value of $q_*(s_t, a)$ is assigned to $V_*(s_t)$

$$\Rightarrow V_*(s_t) = \max_{a \in A(s_t)} (q_*(s_t, a))$$

such that $a \in A(s_t)$

Ans 1)

current state (s)	actions $a \in A(s)$	next state (s')	rewards (r)	$p(s', r s, a)$
high	search	high	r_{search}	κ
high	search	low	r_{search}	$1 - \kappa$
low	search	high	-3	$1 - \beta$
low	search	low	r_{search}	β
high	wait	high	r_{wait}	1
low	wait	low	r_{wait}	1
low	recharge	high	0	1

The table was obtained by using the figure of "transition graph" given in the book for that problem.

Also $\sum_{s', r} p(s', r | s, a) = 1 \forall s, a; s \in S \text{ and } a \in A(s)$. So it is right ~~problem~~.

Ans 3) (3.15) and 3.16)

Signs of the rewards are irrelevant because ~~that~~ they can be changed by adding appropriate const. 'c' to all rewards and rewards are something that we decide. Now the intervals of the values ~~are~~ are important as they show the relative preferences of states. And those intervals can get affected in different ways by adding a constant 'c', based on the task is episodic or continuous.

for Episodic :-

$$V'_\pi(s) = E_\pi [G'_t | S_t = s], \quad G_t = R_{t+1} + R_{t+2} + \dots + R_T$$

$$G'_t = \sum_{k=1}^M (\gamma^k (R_{t+k+1} + c)) \rightarrow \text{terminates after some } M \text{ steps}$$

$$G'_t = \underbrace{\sum_{k=1}^M \gamma^k R_{t+k+1}}_{= G_t} + c \left(\frac{1 - \gamma^M}{1 - \gamma} \right) \rightarrow \text{GP formula}$$

$$\begin{aligned} \Rightarrow V'_\pi(s) &= E_\pi \left[G_t + c \left(\frac{1 - \gamma^M}{1 - \gamma} \right) \mid S_t = s \right] \\ &= E_\pi [G_t \mid S_t = s] + E \left[c \left(\frac{1 - \gamma^M}{1 - \gamma} \right) \mid S_t = s \right] \\ V'_\pi(s) &= \boxed{V_\pi(s) + c \left(\frac{1 - \gamma^M}{1 - \gamma} \right)} \quad \text{const.} \end{aligned}$$

So, if M is smaller, then $c \left(\frac{1 - \gamma^M}{1 - \gamma} \right)$ will be smaller, so states closer to termination will see smaller increase in their preference. $M < \infty$ with probability = 1 and $0 < \gamma < 1$

for continuous task :- $G'_t = G_t + \sum_{k=0}^{\infty} \gamma^k C = G_t + \frac{C}{1-\gamma}$

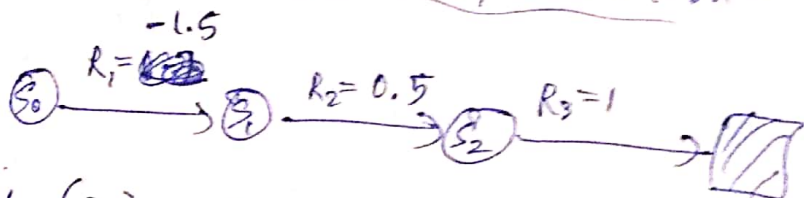
$$V'_a(s) = E_{\pi} \left[G_t + \frac{C}{1-\gamma} \mid S_t = s \right]$$

$$V'_a(s) = V_a(s) + \left(\frac{C}{1-\gamma} \right)$$

So, In this case, each values in offset by the same constant $\frac{C}{1-\gamma}$, so, the relative preferences of states doesn't change.

Example of changing the effect in Episodic task:-

Suppose given



So here $V_a(S_0) = 0$, $V_a(S_1) = -1.5$, $V_a(S_2) = 1$

now let $C = 100$ and $\gamma = 0.1$ so, $V'_a = V_a + \frac{C(1-\gamma^4)}{1-\gamma}$

$$\text{so, for } V'_a(S_0) = 100 \left(\frac{1 - (0.1)^3}{0.9} \right) + 0 = \underline{111}$$

$$V'_a(S_1) = 100 \left(\frac{1 - (0.1)^2}{0.9} \right) + 1.5 = \underline{111.5}$$

$$V'_a(S_2) = 100 \left(\frac{1 - (0.1)}{0.9} \right) + 1 = \underline{101}$$

So, $V_a(S_2) > V_a(S_0)$ but $V'_a(S_2) < V'_a(S_0)$