

Q1 ~~Reasoning~~ average return ~~is~~

s	a	s'	r	$p(s, a s', r)$
high	search	high	r_{search}	α
high	search	low	r_{search}	$1 - \alpha$
low	search	high	r_{search}	$1 - \beta$
low	search	low	r_{search}	β
high	wait	high	r_{wait}	1
low	wait	low	r_{wait}	1
low	recharge	high	0	1

Solved using transition graph. As r_{search} , r_{wait} are ~~also~~ guaranteed rewards for fixed (s, s', a) triple, ~~$p(s, a | s', r)$~~
 $\therefore p(s' | s, a) = p(s', r(s, a, s') | s, a)$

Q3 ~~G_t~~ let R' be new rewards

$$G'_t = R'_t + \gamma R'_{t+1} + \dots = \sum_{k=0}^{\infty} \gamma^k (R'_{t+k+1})$$

$$= \sum_{k=0}^{\infty} \gamma^k (R_{t+k+1} + c)$$

$$= \sum_{k=0}^{\infty} \gamma^k (R_{t+k+1}) + \sum_{k=0}^{\infty} \gamma^k c$$

$$= G_t + \frac{c\gamma}{1-\gamma} \quad \therefore v_c = \frac{c}{1-\gamma}$$

$$\therefore v'_\pi(s) = E[G'_t | S_t = s] = E[G_t + v_c | S_t = s]$$

$$= E[G_t | S_t = s] + v_c$$

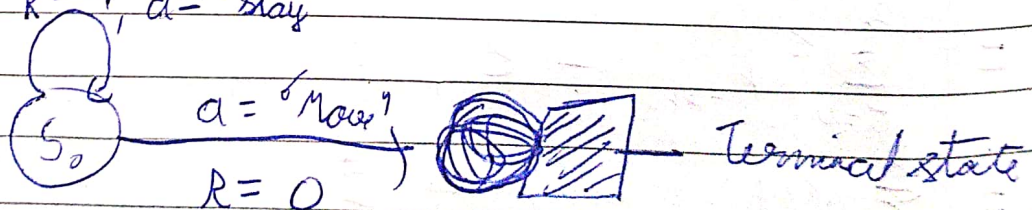
$$= v_{\pi}(s) + v_c$$

\therefore adding constant to all rewards has no effect on relative values of states under any policy.

$$\begin{aligned} \text{Q3b)} \quad G'_t &= \cancel{G_t} + \sum_{k=0}^T \gamma^k R_{t+k+1} = \sum_{k=0}^T \gamma^k (R_t + c) \\ &= \cancel{G_t} + \frac{c(1 - \gamma^{T+1})}{1 - \gamma} \end{aligned}$$

As there is still a factor of T which can depend on policy, there is an effect of the constant.

example: ~~consider a maze-solving~~
~~normally reward is 1 for~~
 $R = -1, a = \text{'stay'}$



In this case the best policy goes straight to terminal state.

However if we add $c = 2$, then the best policy 'stay' is forever looping.

~~Theorem~~

Q5) $Q_*(z) = \max_{a \in A(z)} Q_*(z, a)$

Optimal policy makes best decision greedily

Q

Q