Surabhi S. Nath
_/_/_

Q1 | States: Low, High

Actions: search, Wait, Recharge

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

$$p(s', r | s, a) = p(s' | s, a) \times p(r | s's, a) \quad \text{since independent}$$

can be found from the table

we know $r(s, a, s') = \sum_r r(p(r | s, a, s'))$

Therefore, calculating these equations we get

we can get this from the table.

| a | s | s' | r | $p(s', r | s, a)$. |
|---|---|---|---|---|
| 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 | $[1-\alpha](r_{search})$ |
| low | search | high | 1 | $\alpha(r_{search})$ $[1-\beta]$ |
| low | search | low | 0 | $\beta(1 - r_{search})$ |
| low | search | low | 1 | $\beta r_{search}$ |
| high | wait | high | 0 | $r_{wait}$ |
| high | wait | low | 1 | $r_{wait}$ |
| low | wait | low | 1 | $r_{wait}$ |
| low | wait | low | 0 | $1 - r_{wait}$ |
| low | recharge | high | 0 | $1$ |

**Q3**

**a)** To show : Adding a constant $c$ to all rewards adds a constant $v_c$ to the values of all states and does not affect relative values of any states under any policy.

○ $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$.

Now adding constant $c$ to all rewards,

○ $G_t' = (R_{t+1} + c) + \gamma (R_{t+2} + c) + \gamma^2 (R_{t+3} + c) + \cdots$
$= \sum_{k=0}^{\infty} \gamma^k (R_{t+k+1} + c)$.

$$V_\pi'(s) = E_\pi [G_t' | S_t = c] = E_\pi \left[ \sum_{k=0}^{\infty} \gamma^k (R_{t+k+1} + c) \Big| S_t = s \right] \quad \forall s.$$

$$= E_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} + \sum_{k=0}^{\infty} \gamma^k c \Big| S_t = s \right]$$

∴ value of all states is increased by

$$= E_\pi \left[ \sum_{k=0}^{\infty} \gamma^k (R_{t+k+1}) \Big| S_t = s \right] +$$

$v_c = \boxed{\dfrac{c}{1-\gamma}}$ and relative values of $V(s)$ are

$$E_\pi \left[ \sum_{k=0}^{\infty} \gamma^k (c) \Big| S_t = s \right]$$

not affected. ∵ signs of the rewards are not important as

$$= E_\pi \left[ G_t \Big| S_t = s \right] + \boxed{\sum_{k=0}^{\infty} \gamma^k c}$$

adding a constant to make them all positive doesn't affect the relative $V$ values of the learning

const. $c \quad c \times \dfrac{1}{1-\gamma}$

**b) Episodic task**

In an episodic task, the number of steps are limited by $T$. / / /

Thus we get:

$$V'_\pi(s) = E_\pi[G_t | S_t = s] + \left(\sum_{k=0}^{T} \gamma^k c\right)$$

$$c\left(\frac{1 - \gamma^{T+1}}{1 - \gamma}\right) = V_c$$

$V_c$ here depends on $T$ ie the number of steps.

~~in the simulation~~

However, $T$ is fixed for an episode ~~task~~ and all $V_\pi(s)$ will hence be added by the same value. ∴ $V_c$ is constant within ~~across~~ episodes. Across episodes $T$ may change ∴ relative values across episodes can be different.

④ This can also be visualized as a special case of the previous formulation with $\infty$ steps. Here, reward after time instant $T$ will repeatedly be $0$. ∴ for a given $T$, it is same as the previous case.

---

**Q5** Equation for $V_*$ in terms of $q_*$

$$V_*(s) = \max_{a \in A(s)} q_*(s, a). \qquad ①$$

$$V_*(s) = \max_{a \in A(s)} \sum p(s', \gamma | s, a)\left(1 + \gamma V_*(s')\right) \qquad ②$$

Using ① and ② we get:

$$V_*(s) = \max_{a \in A(s)} \sum p(s', r | s, a)\left(1 + \gamma \max_{a' \in A(s')} q_*(s', a')\right)$$