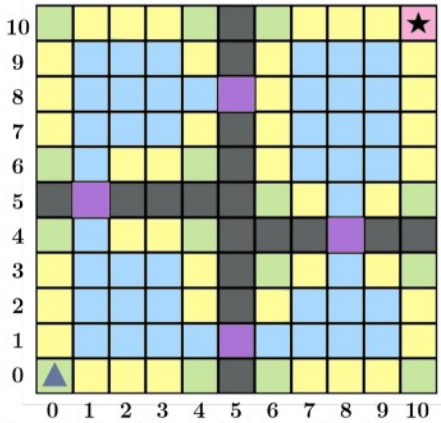


Q₁ = (a) Action = { left, right, up, down }.

State = (i, j); i, j ∈ [0, 10].

(i, j) ≠ (5, i), (j, 5), (k, 4) $\left\{ \begin{array}{l} i = 0, 2, 3, 4, 5, 6, 7, 9, 10 \\ j = 0, 2, 3, 4. \\ k = 6, 7, 9, 10 \quad (\text{except "Walls"}) \end{array} \right.$

b)



As show in figure :

it have 44 states that can move in 4 direction (blue),

40 states that can move in 3 direction (yellow)

(15+4) states that can move in 2 direction (green & purple)

and 1 goal state that return to the start state.

Therefore,

$$44 \times 4 \times 3 + 40 \times 4 \times 3 + 15 \times 2 \times 2 + 15 \times 2 \times 2 + 4 \times 2 \times 2 + 4 \times 2 \times 2 + 1 = 1199.$$

I think there are 1199 non-zero rows in this conditional probability table.

Q₂ = (a) episodic function: $G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T.$

with discounting: $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T.$

$$= \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1} = -\gamma^{T-t-1}$$

continuing function: $G'_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} = \gamma^k$

Because "episodic" has only 1 failure, it returns the last G_t , while "continuing" will have many failures and the return value will always be updated.

(b) $R_1 = R_2 = \dots = R_{t+k} = 0$; $R_{t+k+1} = 1$.

Then the G_t will always be 1, so it will never have improvement.

Q3: (a) $G_5 = R_6 + rG_6 = 0$

$G_4 = R_5 + rG_5 = 2$

$G_3 = R_4 + rG_4 = 3 + 0.5 \times 2 = 4$

$G_2 = R_3 + rG_3 = 6 + 0.5 \times 4 = 8$

$G_1 = R_2 + rG_2 = 2 + 0.5 \times 8 = 6$

$G_0 = R_1 + rG_1 = 1 + 0.5 \times 6 = 2$

(b) $G_1 = R_2 + rR_3 + r^2R_4 + \dots + r^nR_{n+2}$
 $= \sum_{i=0}^{\infty} r^i R_{n+2+i}$ ($R_1=2, R_2=R_3=\dots=R_n=7$)

$= \frac{1}{1-r} \times 7 = \frac{1}{1-0.9} \times 7 = 70$

$G_0 = R_1 + 0.9 \times 70 = 2 + 63 = 65$.

Q4: $\sum_{i=1}^{100} r^i = r + r^2 + \dots + r^{100} = \frac{(r^{100} - 1)r}{r-1}$

$$\begin{cases} G_0 = R_1 + \sum_{i=1}^{100} r^i R_{101-i} = 50 + \sum_{i=1}^{100} r^i \times (-1) = 50 - \frac{(r^{100} - 1)r}{r-1} \\ G'_0 = R'_1 + \sum_{i=1}^{100} r^i R'_{101-i} = (-50) + \sum_{i=1}^{100} r^i \times 1 = (-50) + \frac{(r^{100} - 1)r}{r-1} \end{cases}$$

When $G_0 > G'_0$: $50 - \frac{(r^{100} - 1)r}{r-1} > -50 + \frac{(r^{100} - 1)r}{r-1}$
 $50 > \frac{(r^{100} - 1)r}{r-1}$
 $-1.04 < r < 0.9844$

Thus, while $0 < r < 0.9844$, choose [UP]; else, choose [DOWN]

Q5: (a) $G_t = \sum_{i=0}^{\infty} r^i R_{t+k+i}$

$V_N(s) \doteq E_N[G_t | S_t = s] = E_N\left[\sum_{i=0}^{\infty} r^i R_{t+k+i} \mid S_t = s\right]$

$V_N(s)(c) \doteq E_N\left[\sum_{i=0}^{\infty} r^i (R_{t+k+i} + c) \mid S_t = s\right]$

$= E_N\left[\sum_{i=0}^{\infty} r^i R_{t+k+i} + \sum_{i=0}^{\infty} r^i c \mid S_t = s\right]$

$= E_N\left[\sum_{i=0}^{\infty} r^i R_{t+k+i} \mid S_t = s\right] + E_N\left[\sum_{i=0}^{\infty} r^i c \mid c = c\right]$

$$\begin{aligned}
&= E_{\pi} \left[\sum_{i=0}^{\infty} \gamma^i R_{t+k+1} + \sum_{i=0}^{\infty} \gamma^i c \mid S_t = s \right] \\
&= E_{\pi} \left[\sum_{i=0}^{\infty} \gamma^i R_{t+k+1} \mid S_t = s \right] + E_{\pi} \left[\sum_{i=0}^{\infty} \gamma^i c \mid S_t = s \right] \\
&= V_{\pi}(s) + \sum_{i=0}^{\infty} \gamma^i c \\
&= V_{\pi}(s) + \frac{c}{1-\gamma}
\end{aligned}$$

(b) If we make the reward after each action a positive number, then the result will remain the same, it will not move. Since the first time it already got a positive result. The requirement to add a constant without bringing an effect is that the reward remains negative after adding the constant, and if it becomes positive after adding the constant, it will have the effect I mentioned before.

Q6: (a)
$$\begin{aligned}
V_{\pi}(s) &= \frac{1}{4} \times 1 \times (0 + 0.9 \times 2.3) + \frac{1}{4} \times 1 \times (0 + 0.9 \times 0.4) \\
&\quad + \frac{1}{4} \times 1 \times (0 + 0.9 \times (-0.4)) + \frac{1}{4} \times 1 \times (0 + 0.9 \times 0.7) \\
&= \frac{1}{4} \times 0.9 \times 2.3 + \frac{1}{4} \times 0.9 \times 0.7 \\
&= 0.5175 + 0.1575 = 0.675
\end{aligned}$$

(b)
$$V_{\pi}(s) = \frac{1}{2} \times 1 \times (0 + 0.9 \times 19.8) + \frac{1}{2} \times 1 \times (0 + 0.9 \times 19.8) = 17.82.$$

Q7: (a) I guess the value function $= \frac{1}{2} \times 0 + \frac{1}{2} \times 1 = \frac{1}{2}.$

Verify:
$$\begin{aligned}
V_{\pi}(s) &= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma V_{\pi}(s')] \\
&= \frac{1}{2} \times 1 \times [1 + 1 \times 0] + \frac{1}{2} \times 1 \times [0 + 1 \times 0] \\
&= \frac{1}{2} \times 1 + \frac{1}{2} \times 0 = \frac{1}{2}.
\end{aligned}$$

(b)
$$V_{\pi}(s)(A) = \frac{1}{2} \times 1 \times [0 + 1 \times 0] + \frac{1}{2} \times 1 \times [0 + 1 \times V_{\pi}(s)(B)] = \frac{1}{2} V_{\pi}(s)(B)$$

$$V_{\pi}(s)(B) = \frac{1}{2} \times 1 \times [0 + 1 \times V_{\pi}(s)(A)] + \frac{1}{2} \times 1 \times [0 + 1 \times V_{\pi}(s)(C)] = \frac{1}{2} V_{\pi}(s)(A) + \frac{1}{2} V_{\pi}(s)(C)$$

$$V_{\pi}(s)(C) = \frac{1}{2} \times 1 \times [0 + 1 \times V_{\pi}(s)(B)] + \frac{1}{2} \times 1 \times [0 + 1 \times V_{\pi}(s)(D)] = \frac{1}{2} V_{\pi}(s)(B) + \frac{1}{2} V_{\pi}(s)(D)$$

$$V_{\pi}(s)(D) = \frac{1}{2} \times 1 \times [0 + 1 \times V_{\pi}(s)(C)] + \frac{1}{2} \times 1 \times [0 + 1 \times V_{\pi}(s)(E)] = \frac{1}{2} V_{\pi}(s)(C) + \frac{1}{2} V_{\pi}(s)(E)$$

$$V_{\pi(s)}(D) = \frac{1}{2} \times 1 \times [0 + 1 \times V_{\pi(s)}(C)] + \frac{1}{2} \times 1 \times [0 + 1 \times V_{\pi(s)}(E)] = \frac{1}{2} V_{\pi(s)}(C) + \frac{1}{2} V_{\pi(s)}(E)$$

$$V_{\pi(s)}(E) = \frac{1}{2} \times 1 \times [0 + 1 \times V_{\pi(s)}(D)] + \frac{1}{2} \times 1 \times [1 + 1 \times 0] = \frac{1}{2} V_{\pi(s)}(D) + \frac{1}{2}$$

$$\begin{cases} V_A = \frac{1}{6}; & V_D = 4V_A = \frac{2}{3}; \\ V_B = 2V_A = \frac{1}{3}; & V_E = 5V_A = \frac{5}{6}. \\ V_C = 3V_A = \frac{1}{2}; \end{cases}$$

$$(c) \quad V_{\pi(s)}(i) = \frac{i}{n-1} \quad (i = 1, 2, 3, \dots, (n-2))$$

Q8: (a) $S_{high} = \{\text{search, wait}\}$.

$$S_{low} = \{\text{search, wait, recharge}\}.$$

$$v(s) = E[R_{t+1} + \gamma V(s_{t+1}) | S_t = s].$$

$$\begin{aligned} v_{high} &= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma V_{\pi}(s')] \\ &= \pi(\text{search} | \text{high}) [\gamma \text{search} + \alpha \times \gamma \times v_{high} + (1-\alpha) \times \gamma \times v_{low}] \\ &\quad + \pi(\text{wait} | \text{high}) [1 \times (\gamma \text{wait} + \gamma \times v_{high})] \end{aligned}$$

$$\begin{aligned} v_{low} &= \pi(\text{search} | \text{low}) [\gamma \text{search} + \beta \times \gamma \times v_{low} + (1-\beta) \times \gamma \times v_{high}] \\ &\quad + \pi(\text{wait} | \text{low}) [1 \times (\gamma \text{wait} + \gamma \times v_{low})] \\ &\quad + \pi(\text{recharge} | \text{low}) [0 + 1 \times \gamma \times v_{high}] \end{aligned}$$

$$\begin{aligned} (b) \quad v_{high} &= 1 \times [0 + 0.8 \times 0.9 \times v_{high} + 0.2 \times 0.9 \times v_{low}] \\ &= 10 + 0.72 v_{high} + 0.18 v_{low} \end{aligned}$$

$$\begin{aligned} v_{low} &= 0.5 \times [3 + 0.9 \times v_{low}] + 0.5 \times 0.9 \times v_{high} \\ &= 1.5 + 0.45 v_{low} + 0.45 v_{high}. \end{aligned}$$

$$\begin{cases} 0.28 v_{high} = 10 + 0.18 v_{low} \\ 0.55 v_{low} = 1.5 + 0.45 v_{high} \end{cases}$$

Thus $\begin{cases} v_{high} = 77.04 \\ v_{low} = 1.5 \end{cases}$

$$\text{Thus } \begin{cases} v(\text{high}) = 79.04 \\ v(\text{low}) = 67.40. \end{cases}$$

$$\text{Q9: (a) } V_N(s) \doteq E_N[G_t | S_t = s]$$

$$= E_N\left[\sum_{k=0}^{\infty} r^k R_{t+k+1} \mid S_t = s\right]$$

$$q_N(s, a) \doteq E_N[G_t | S_t = s, A_t = a]$$

$$= E_N\left[\sum_{k=0}^{\infty} r^k R_{t+k+1} \mid S_t = s, A_t = a\right]$$

$$\left. \begin{array}{l} V_N(s) \doteq E_N[G_t | S_t = s] \\ q_N(s, a) \doteq E_N[G_t | S_t = s, A_t = a] \end{array} \right\} V_N(s) = \sum_a \pi(a|s) q(s, a)$$

$$(b) q_N(s, a) \doteq E_N\left[\sum_{k=0}^{\infty} r^k R_{t+k+1} \mid S_t = s, A_t = a\right]$$

$$= \sum_{s', r} p(s', r | s, a) [r + \gamma V_N(s')]$$

$$(c) V_N(s) = \sum_{a'} \pi(a'|s) q(s', a')$$

$$q_N(s, a) = \sum_{s', r} p(s', r | s, a) \left[r + \gamma \left(\sum_{a'} \pi(a'|s') q(s', a') \right) \right]$$