

1.

$$\begin{aligned}
 (1): \quad V_{\pi}(s) &\leq q_{\pi}(s, \pi'(s)) \\
 &= E[R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s, A_t = \pi'(s)] \\
 &\leq E_{\pi'}[R_{t+1} + \gamma q_{\pi}(S_{t+1}, \pi'(S_{t+1})) | S_t = s] \\
 &= E_{\pi'}[R_{t+1} + \gamma E_{\pi'}[R_{t+2} + \gamma v_{\pi}(S_{t+2}) | S_{t+1}, A_t = \pi'(S_{t+1})] | S_t = s] \\
 &= E_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 v_{\pi}(S_{t+2}) | S_t = s] \\
 &\leq E_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 v_{\pi}(S_{t+3}) | S_t = s] \\
 &\vdots \\
 &\leq E_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s] \\
 &= E_{\pi'}[G_t | S_t = s] \\
 &= v_{\pi}(s)
 \end{aligned}$$

$$\text{由 } \pi'(a|s) \text{ 知, } \pi'(s) = \arg \max_{a \in A} q_{\pi}(s, a)$$

$$\text{故: } q_{\pi}(s, \pi'(s)) = \max_{a \in A} q_{\pi}(s, a)$$

$$\text{则: } v_{\pi}(s) = \sum_{a \in A} \pi(a|s) q_{\pi}(s, a) \leq \sum_{a \in A} \pi(a|s) q_{\pi}(s, \pi'(s)) = q_{\pi}(s, \pi'(s)) \sum_{a \in A} \pi(a|s)$$

$$\text{故, } v_{\pi}(s) \leq q_{\pi}(s, \pi'(s)).$$

$$\text{故由(1), } v_{\pi}(s) \leq v_{\pi'}(s).$$

(3):

$$q_{\pi}(s, \pi'(s)) = \sum_{a \in A} \pi'(a|s) q_{\pi}(s, a)$$

$$= \frac{\epsilon}{|A|} \sum_{a \in A} q_{\pi}(s, a) + (1-\epsilon) \max_{a \in A} q_{\pi}(s, a)$$

$$\geq \frac{\epsilon}{|A|} \sum_{a \in A} q_{\pi}(s, a) + (1-\epsilon) \sum_{a \in A} \frac{\pi(a|s) - \frac{\epsilon}{|A|}}{1-\epsilon} \cdot q_{\pi}(s, a)$$

$$= \frac{\epsilon}{|A|} \sum_{a \in A} q_{\pi}(s, a) - \frac{\epsilon}{|A|} \sum_{a \in A} q_{\pi}(s, a) + \sum_{a \in A} \pi(a|s) q_{\pi}(s, a)$$

$$= v_{\pi}(s).$$

故由 (1), $v_{\pi'}(s) \geq v_{\pi}(s)$.

4.

1) $4 \rightarrow 1$: $V(1) = V(4) + \alpha \cdot (-1 + r \cdot V(1) - V(4))$

$$\therefore V(4) = -\frac{1}{2}$$

$1 \rightarrow 4$: $V(1) = V(1) + \alpha \cdot (-1 + r \cdot V(4) - V(1))$

$$\therefore V(1) = -\frac{3}{4}$$

$4 \rightarrow 7$: $V(4) = V(4) + \alpha \cdot (-1 + r \cdot V(7) - V(4))$

$$\therefore V(4) = -\frac{3}{4}$$

$7 \rightarrow \text{terminal}$:

$$V(7) = V(7) + \alpha \cdot (-1 + r \cdot V(\text{terminal}) - V(7))$$

$$\therefore V(7) = -\frac{1}{2}$$

$$\therefore V = \begin{bmatrix} 0 & -\frac{3}{4} & 0 \\ 0 & -\frac{3}{4} & 0 \\ 0 & -\frac{1}{2} & 0 \end{bmatrix}$$

(2):

初始为4. 则由Q表知, 下一步行动为"下":

$R = -1$. $S' = 7$. 则: $A' = \text{"左"}$

$$Q(4, \downarrow) = Q(4, \downarrow) + \alpha [-1 + \gamma Q(7, \leftarrow) - Q(4, \downarrow)]$$

$$\therefore Q(4, \downarrow) = -3$$

令 $S \leftarrow 7$. $A \leftarrow \text{"左"}$

则: $R = -1$. $S' = 6$. 则: $A' = \text{"上"}$

$$Q(7, \leftarrow) = Q(7, \leftarrow) + \alpha [-1 + \gamma Q(6, \uparrow) - Q(7, \leftarrow)]$$

$$\therefore Q(7, \leftarrow) = -3$$

令 $S \leftarrow 6$. $A \leftarrow \text{"上"}$;

则: $R = -1$. $S' = 3$. 则: $A' = \text{"上"}$

$$Q(6, \uparrow) = Q(6, \uparrow) + \alpha [-1 + \gamma Q(3, \uparrow) - Q(6, \uparrow)]$$

$$\therefore Q(6, \uparrow) = -2$$

令 $S \leftarrow 3$. $A \leftarrow \text{"上"}$;

则: $R = -1$. $S' = \text{terminal}$.

$$Q(3, \uparrow) = Q(3, \uparrow) + \alpha [-1 + 0 - Q(3, \uparrow)]$$

$$\therefore Q(3, \uparrow) = -1$$

$S \leftarrow \text{terminal}$. 达到终止态,

此时Q表:

-4	-3	-1	-3	-4	-2	-4
-3	-3	-2	-4	-2	-3	-3
-4	-3	-4	-3	-2	-3	-4
-3	-2	-3	-3	-4	-3	-3

5,

(1) 状态: $S = \{ \overset{b}{\text{西双}}, \overset{c}{\text{普洱}}, \overset{d}{\text{墨江}}, \overset{e}{\text{元江}}, \overset{f}{\text{石屏}}, \overset{g}{\text{峨山}}, \overset{h}{\text{玉溪}}, \overset{i}{\text{晋宁}} \}$.

行动: $A = \{ \text{向上}, \text{向下} \}$

(2) 转移矩阵

	b	c	d	e	f	g	h	i
b	1							
c	0.5		0.5					
d		0.5		0.5				
e			0.5		0.5			
f				0.5		0.5		
g					0.5		0.5	
h						0.5		0.5
i								1

$$(3) \quad v_{\pi}(c) = \frac{1}{2} \cdot (0 + v_{\pi}(b)) + \frac{1}{2} \cdot (0 + v_{\pi}(d))$$

$$v_{\pi}(d) = \frac{1}{2} (0 + v_{\pi}(c)) + \frac{1}{2} (0 + v_{\pi}(e))$$

$$v_{\pi}(e) = \frac{1}{2} (0 + v_{\pi}(f)) + \frac{1}{2} (0 + v_{\pi}(d))$$

$$v_{\pi}(f) = \frac{1}{2} (0 + v_{\pi}(g)) + \frac{1}{2} (0 + v_{\pi}(e))$$

$$v_{\pi}(g) = \frac{1}{2} (0 + v_{\pi}(h)) + \frac{1}{2} (0 + v_{\pi}(f))$$

$$v_{\pi}(h) = \frac{1}{2} (0 + v_{\pi}(i)) + \frac{1}{2} (0 + v_{\pi}(g))$$

b, i 为终止态. 则 $v_{\pi}(b) = v_{\pi}(i) = 0$

由上述方程组解得:

$$V_{\pi}(c) = 100, V_{\pi}(d) = 200.$$

$$V_{\pi}(e) = 300, V_{\pi}(f) = 400$$

$$V_{\pi}(g) = 500, V_{\pi}(h) = 600$$

(4): 由于 b 和 i 为终止点, 故在这 2 个状态不会有后续行动.

$$\text{例: } q_{\pi}(c, \uparrow) = 0 + 1 \cdot 1 \cdot V_{\pi}(d) = 200.$$

$$q_{\pi}(c, \downarrow) = 0 + 1 \cdot 1 \cdot V_{\pi}(b) = 0$$

$$q_{\pi}(d, \uparrow) = 0 + V_{\pi}(e) = 300.$$

$$q_{\pi}(d, \downarrow) = 0 + V_{\pi}(c) = 100.$$

$$q_{\pi}(e, \uparrow) = 0 + V_{\pi}(f) = 400$$

$$q_{\pi}(e, \downarrow) = 0 + V_{\pi}(d) = 200$$

$$q_{\pi}(f, \uparrow) = 0 + V_{\pi}(g) = 500$$

$$q_{\pi}(f, \downarrow) = 0 + V_{\pi}(e) = 300$$

$$q_{\pi}(g, \uparrow) = 0 + V_{\pi}(h) = 600$$

$$q_{\pi}(g, \downarrow) = 0 + V_{\pi}(f) = 400$$

$$q_{\pi}(h, \uparrow) = 700 + V_{\pi}(i) = 700$$

$$q_{\pi}(h, \downarrow) = 0 + V_{\pi}(g) = 500$$

5. 可以采用 SARSA 算法来解决这个问题, 令一个学习过程如下:

首先初始状态所在地 (非终止点). 初始化 $Q(s, a)$ 表格

之后, 对于初始态, 根据 Q 表以及 ϵ -贪心策略选择

行动 a

之后进入循环迭代,

根据 A 获得单步行动回报 R 以及下一状态 S' .

根据 S' 贪心, 选择下一动作 A' .

之后更新 $Q(S, A)$:

$$Q(S, A) += \alpha [R + \gamma Q(S', A') - Q(S, A)],$$

更新 $S \leftarrow S', A \leftarrow A'$.

直到达到终点.

不断进行上述学习过程, 则 $Q(S, A)$ 会逐渐收敛.

利用最后得到的 $Q(S, A)$ 表格, 获知从群

初始位置, 即可由贪心策略给出一条达到目标地的最优方案.

(6) 可以采用动态规划的思想求解.

同时可以利用异步迭代的方法.

对于大系统当前的状态 S , 及迭代次数 k .

$$\text{首先进行价值迭代: } V_{k+1}(S) \leftarrow \max_A \left[R_S + \gamma \sum_{S'} P_{S'}^A V_k(S') \right]$$

之后立刻对 $V_{k+1}(S)$ 更新, 并用确定贪婪策略改进 $\pi(S)$.

这样大系统每走一步便更新之前的状态价值函数 $V(S)$

以及策略 $\pi(S)$.

可以较方便地选择何时停止更新, 并以此时的 $V(S)$ 近似状态价值.

当前策略 $\pi(S)$ 近似为最优策略.

故可解决状态过多的问题.