

2018101036
MDL-Assign 2 (Part 1)
BHAVESH SHUKLA

Using value iteration algorithm, we have:

$$U_t(i) = \max [R(i, a) + \gamma \sum_j P(j|i, a) \times U_t(j)]$$

(where $U_t(i)$ is utility of state i on time t .)

→ This eqⁿ is called Bellman Update Eqⁿ

States :

S_0	S_1	S_2	$S_3 + 10$
-------	-------	-------	------------

 where S_3 is absorbant state

Actions → $[l, r]$ l = left & r = right
 l → move left with $p=0.8$, right with $p=0.2$
 r → move right with $p=0.8$, left with $p=0.2$

Bellman factor = $0.01 = \delta$

Gamma = $\gamma = 0.25$

t_0 → Iteration 0

$t_0 =$

0	0	0	+10
---	---	---	-----

t_1 → Iteration 1

$$S_0 = \max(-1 + \delta(0 \times 0.8 + 0 \times 0.2), -1 + \delta(0 \times 0.8 + 0 \times 0.2))$$

$$S_0 = \max(-1, -1) = -1$$

$$S_1 = \max(-1 + \delta(0 \times 0.8 + 0 \times 0.2), -1 + \delta(0 \times 0.8 + 0 \times 0.2))$$

$$S_1 = \max(-1, -1) = -1$$

$$S_2 = \max(-1 + \delta(0 \times 0.8 + 10 \times 0.2), -1 + \delta(10 \times 0.8 + 0 \times 0.2))$$

$$= \max(-0.5, 1) = 1$$

$$\Rightarrow t_1 = \begin{bmatrix} -1 & -1 & 1 & 10 \end{bmatrix}, \delta = 1 \quad \boxed{2.0000000000000000}$$

$t_2 \rightarrow$ Iteration 2

$$S_0 = \max(-1 + \gamma(-1 \times 0.8 + -1 \times 0.2), -1 + \gamma(-1 \times 0.2 + -1 \times 0.8))$$

$$= \max(-1 + \gamma(-1), -1 + \gamma(-1))$$

$$= -1 + \gamma = -1.25$$

$$S_1 = \max(-1 + \gamma(-1 \times 0.8 + 1 \times 0.2), -1 + \gamma(1 \times 0.8 + (-1) \times 0.2))$$

$$= \max(-1 + \gamma(0.6), -1 + \gamma(0.6))$$

$$= -1 + \gamma(0.6) = -0.85$$

$$S_2 = \max(-1 + \gamma(-1 \times 0.8 + 10 \times 0.2), -1 + \gamma(10 \times 0.8 + (-1) \times 0.2))$$

$$= \max(-1 + \gamma(7.2), -1 + \gamma(7.8))$$

$$= -1 + \gamma(7.8) = +0.95$$

$$\Rightarrow t_2 = \begin{bmatrix} -1.25 & -0.85 & +0.95 & 10 \end{bmatrix}, \delta = 0.25$$

$t_3 \rightarrow$ Iteration 3

$$S_0 = \max(-1 + \gamma(-1.25 \times 0.8 + -0.85 \times 0.2), -1 + \gamma(-0.85 \times 0.8 + -1.25 \times 0.2))$$

$$= \max(-0.71, -1.2325) \max(-1 + \gamma(-1.41), -1 + \gamma(-0.945))$$

$$= -1 + \gamma(-0.945)$$

$$= -1.2325$$

$$S_1 = \max(-1 + \gamma(-1.25 \times 0.8 + -0.95 \times 0.2), -1 + \gamma(0.95 \times 0.8 + -1.25 \times 0.2))$$

$$= \max(-1.2975, -0.8725) = -0.8725$$

$$S_2 = \max(-1 + 0.25(-0.85 \times 0.8 + 10 \times 0.2), -1 + \gamma(10 \times 0.8 + -0.85 \times 0.2))$$

$$S_2 = \max(-1 + 0.33, -1 + 0.25(8 - 0.17)) = \max(0.67, 2.00)$$

$$S_2 = \max \{-0.87, 0.9575\} = 0.9575$$

$$\Rightarrow t_3 = [-1.2325 \mid -0.8725 \mid 0.9575 \mid 10] \quad \delta = 0.0225$$

$t_4 \rightarrow$ Iteration 4

$$\begin{aligned} S_0 &= \max \left\{ -1 + 0.25(0.8 \times (-1.233) + 0.2 \times (-0.873)), \right. \\ &\quad \left. -1 + 0.25(0.2 \times (-1.233) + 0.8 \times (-0.873)) \right\} \\ &= \max \{-1 + 0.25(-1.61, -0.945)\} = -1 + 0.25(-0.945) \\ &= -1.236 \end{aligned}$$

$$\begin{aligned} S_1 &= \max \left\{ -1 + 0.25(0.8 \times (-1.233) + 0.2 \times (0.958)), \right. \\ &\quad \left. -1 + 0.25(0.2 \times (-1.233) + 0.8 \times (0.958)) \right\} \\ &= \max \{-1 + 0.25 \times (-0.794), -1 + 0.25 \times 0.52\} \\ &= -0.873 \end{aligned}$$

$$\begin{aligned} S_2 &= \max \left\{ -1 + 0.25(0.8 \times (-0.873) + 0.2 \times 10), \right. \\ &\quad \left. -1 + 0.25(0.2 \times (-0.873) + 0.8 \times 10) \right\} \\ &= 0.958 \end{aligned}$$

$$\Rightarrow t_4 = [-1.236 \mid -0.873 \mid 0.958 \mid 10]$$

$$\delta = 0.003625$$

Max function results in move right

Ans) \Rightarrow Policy = {Move Right, Move Right, Move Right, No Action}

1) Reward is always 10

2) Movement on extremes is same as staying in same cell.

2019101036