

Machine, Data, & Learning

24/01/2023

Monday

Aayem Ajay Sharma
Roll no: 2022121001

Bellman Equation:

$$U(s) = R(s, a) + \gamma \cdot \max \left[\sum_{s'} P(s'|s, a) U(s') \right]$$

Iteration - 0 :

$$U_0 = \begin{bmatrix} 0 & 0 & 1 & -1 \\ 1 & 0 & 0 & 0 \\ 2 & 0 & \text{WALL} & 0 \\ 3 & 0 & 0 & 0 \end{bmatrix}$$

Iteration - 1 :

$$\begin{aligned}
 U[0,0] &= -0.04 + 0.95 \max [0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, \text{up} \\
 &\quad 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 1, \text{down} \\
 &\quad 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, \text{left} \\
 &\quad 0.7 \times 1 + 0.15 \times 0 + 0.15 \times 0, \text{right}] \\
 &= -0.04 + 0.95 \times 0.07 \\
 &= 0.625
 \end{aligned}$$

$$U[0,1] = 1$$

$$U[0,2] = -1$$

$$U[1,0] = -0.04 + 0.95 \max \left[\begin{array}{l} 0.7x_0 + 0.15x_0 + 0.15x_0, \text{ up} \\ 0.7x_0 + 0.15x_0 + 0.15x_0, \text{ down} \\ 0.7x_0 + 0.15x_0 + 0.15x_0, \text{ Left} \\ 0.7x_0 + 0.15x_0 + 0.15x_0, \text{ right} \end{array} \right]$$

$$= -0.04 + 0.95 \times 0$$

$$= -0.04$$

$$U[1,1] = -0.04 + 0.95 \max \left[\begin{array}{l} 0.7x_1 + 0.15x_0 + 0.15x_0, \text{ up} \\ 0.7x_0 + 0.15x_0 + 0.15x_0, \text{ down} \\ 0.7x_0 + 0.15x_1 + 0.15x_0, \text{ left} \\ 0.7x_0 + 0.15x_0 + 0.15x_1, \text{ right} \end{array} \right]$$

$$= -0.04 + 0.95 \times 0.7$$

$$= 0.625$$

$$U[1,2] = -0.04 + 0.95 \max \left[\begin{array}{l} 0.7x_{-1} + 0.15x_0 + 0.15x_0, \text{ up} \\ 0.7x_0 + 0.15x_0 + 0.15x_0, \text{ down} \\ 0.7x_0 + 0.15x_{-1} + 0.15x_0, \text{ left} \\ 0.7x_0 + 0.15x_0 + 0.15x_{-1}, \text{ right} \end{array} \right]$$

$$= -0.04 + 0.95 \times 0$$

$$= -0.04$$

$$U[2,0] = -0.04 + 0.95 \times \max \left[\begin{array}{l} 0.7x_0 + 0.15x_0 + 0.15x_0, \text{ up} \\ 0.7x_0 + 0.15x_0 + 0.15x_0, \text{ down} \\ 0.7x_0 + 0.15x_0 + 0.15x_0, \text{ left} \\ 0.7x_0 + 0.15x_0 + 0.15x_0, \text{ right} \end{array} \right]$$

$$U[2,0] = -0.04 + 0.95 \times 0 \\ = -0.04$$

$U[2,1] = \text{WALL}$

$$U[2,1] = -0.04 + 0.95 \max [0.7x_0 + 0.15x_0 + 0.15x_0, \text{up} \\ 0.7x_0 + 0.15x_0 + 0.15x_0, \text{down} \\ 0.7x_0 + 0.15x_0 + 0.15x_0, \text{left} \\ 0.7x_0 + 0.15x_0 + 0.15x_0, \text{right}] \\ = -0.04 + 0.95 \times 0 \\ = -0.04$$

$$U[3,0] = -0.04 + 0.95 \max [0.7x_0 + 0.15x_0 + 0.15x_0, \text{up} \\ 0.7x_0 + 0.15x_0 + 0.15x_0, \text{left} \\ 0.7x_0 + 0.15x_0 + 0.15x_0, \text{down} \\ 0.7x_0 + 0.15x_0 + 0.15x_0, \text{right}] \\ = -0.04 + 0.95 \times 0 \\ = -0.04$$

$$U[3,1] = -0.04 + 0.95 \max [0.7x_0 + 0.15x_0 + 0.15x_0, \text{up} \\ 0.7x_0 + 0.15x_0 + 0.15x_0, \text{down} \\ 0.7x_0 + 0.15x_0 + 0.15x_0, \text{left} \\ 0.7x_0 + 0.15x_0 + 0.15x_0, \text{right}] \\ U[3,1] = -0.04 + 0.95 \times 0 \\ = -0.04$$

$$U[3,2] = -0.04 + 0.95 \max [$$

+ 0.7x0 + 0.15x0 + 0.15x0, up
 0.7x0 + 0.15x0 + 0.15x0, down
 0.7x0 + 0.15x0 + 0.15x0, left
 0.7x0 + 0.15x0 + 0.15x0, right]

$$= -0.04 + 0.95x0$$

$$U_1 = \begin{bmatrix} 0.625 & 1 & -1 \\ -0.04 & 0.625 & -0.04 \\ -0.04 & WALL & -0.04 \\ -0.04 & -0.04 & -0.04 \end{bmatrix}$$

Iteration 2:

$$U[0,0] = -0.04 + 0.95 \max [$$

(up) 0.7x0.625 + 0.15x0.625 + 0.15x1
 (down) 0.7x0.04 + 0.15x0.625 + 0.15x1
 (left) 0.7x0.625 + 0.15x0.625 + 0.15x-0.04
 (right) 0.7x1 + 0.15x0.625 + 0.15x-0.04]

$$U[0,0] = -0.04 + 0.95(0.7x1 + 0.15x0.625 + 0.15x-0.04)$$

$$= -0.04 + 0.95 \times 0.768$$

$$= -0.04 + 0.748$$

$$= 0.708$$

$$U[0,1] = +1 \quad (\text{Terminating state})$$

$$U[0,2] = -1 \quad (\text{Terminating state})$$

$$U[1,1] = -0.04 + 0.95 \times \text{max of}$$

$$1) 0.7 \times 1 + 0.15 \times -0.04 + 0.15 \times -0.04, \text{ up}$$

$$2) 0.7 \times 0.625 + 0.15 \times -0.04 + 0.15 \times -0.04, \text{ down}$$

$$3) 0.7 \times -0.04 + 0.15 \times 1 + 0.15 \times 0.625, \text{ left}$$

$$4) 0.7 \times -0.04 + 0.15 \times 1 + 0.15 \times 0.625, \text{ right}$$

$$= -0.04 + 0.95(0.7 \times 1 + 0.15 \times -0.04 + 0.15 \times -0.04)$$

$$= -0.04 + 0.95 \times 0.688$$

$$= -0.04 + \cancel{0.64} + 0.654$$

$$= \cancel{0.64} + 0.614.$$

$$\boxed{U[1,1] = 0.614}$$

$$U[1,2] = -0.04 + 0.95 \times \text{max of}$$

$$1) 0.7 \times -1 + 0.15 \times 0.625 + 0.15 \times -0.04, \text{ up}$$

$$2) 0.7 \times -0.04 + 0.15 \times 0.625 + 0.15 \times -0.04, \text{ down}$$

$$3) 0.7 \times 0.625 + 0.15 \times -1 + 0.15 \times -0.04, \text{ left}$$

$$4) 0.7 \times -0.04 + 0.15 \times -1 + 0.15 \times -0.04, \text{ right}$$

$$U[1,2] = -0.04 + 0.95(0.7 \times 0.625 + 0.15 \times -1 + 0.15 \times -0.04)$$

$$= -0.04 + 0.95 \times 0.282$$

$$= -0.04 + 0.267$$

$$= 0.227$$

$$U[1,2] = 0.227$$

$$U[1,0] = -0.04 + 0.95 \times \max \text{ of}$$

$$(up) 1) 0.7 \times 0.625 + 0.15 \times -0.04 + 0.15 \times 0.625$$

$$(down) 2) 0.7 \times -0.04 + 0.15 \times -0.04 + 0.15 \times 0.625$$

$$(left) 3) 0.7 \times -0.04 + 0.15 \times -0.04 + 0.15 \times 0.625$$

$$(right) 4) 0.7 \times 0.625 + 0.15 \times -0.04 + 0.15 \times 0.625$$

$$U[1,0] = -0.04 + 0.95 \times (0.7 \times 0.625 + 0.15 \times -0.04 + 0.15 \times 0.625)$$

$$= -0.04 + 0.95 \times 0.525$$

$$= -0.04 + 0.499$$

$$= 0.459.$$

$$U[1,0] = 0.459$$

$$U[2,0] = -0.04 + 0.95 \times \max \text{ of}$$

$$(up) 1) 0.7 \times -0.04 + 0.15 \times -0.04 + 0.15 \times -0.04$$

$$(down) 2) 0.7 \times -0.04 + 0.15 \times -0.04 + 0.15 \times -0.04$$

$$(left) 3) 0.7 \times -0.04 + 0.15 \times -0.04 + 0.15 \times -0.04$$

$$(right) 4) 0.7 \times -0.04 + 0.15 \times -0.04 + 0.15 \times -0.04$$

$$\begin{aligned}
 U[2,0] &= -0.04 + 0.95x(0.7x - 0.04 + 0.15x - 0.04 + 0.15x - 0.04) \\
 &= -0.04 + 0.95x - 0.04 \\
 &= 1.95x - 0.04 \\
 &= -0.078
 \end{aligned}$$

$$U[2,0] = -0.078$$

$U[2,1]$ = WALL

$$U[2,2] = -0.04 + 0.95x \text{ max of}$$

$$(up) 1) 0.7x - 0.04 + 0.15x - 0.04 + 0.15x - 0.04$$

$$(down) 2) 0.7x - 0.04 + 0.15x - 0.04 + 0.15x - 0.04$$

$$(left) 3) 0.7x - 0.04 + 0.15x - 0.04 + 0.15x - 0.04$$

$$(right) 4) 0.7x - 0.04 + 0.15x - 0.04 + 0.15x - 0.04$$

$$U[2,2] = -0.04 + 0.95x(0.7x - 0.04 + 0.15x - 0.04 + 0.15x - 0.04)$$

$$= -0.04 + 0.95x - 0.04$$

$$= 1.95x - 0.04$$

$$= -0.078$$

$$U[2,2] = -0.078$$

$$U[3,0] = -0.04 + 0.95x \text{ max of}$$

$$(up) 1) 0.7x - 0.04 + 0.15x - 0.04 + 0.15x - 0.04$$

$$(down) 2) 0.7x - 0.04 + 0.15x - 0.04 + 0.15x - 0.04$$

$$(left) 3) 0.7x - 0.04 + 0.15x - 0.04 + 0.15x - 0.04$$

$$(right) 4) 0.7x - 0.04 + 0.15x - 0.04 + 0.15x - 0.04$$

$$\begin{aligned}
 V[3,0] &= -0.04 + 0.95(0.7x - 0.04 + 0.15x - 0.04 + 0.15x - 0.04) \\
 &= -0.04 + 0.95x - 0.04 \\
 &= -0.04 + 1.95x \\
 &= -0.078
 \end{aligned}$$

$$V[3,1] = -0.04 + \text{max of}$$

- (up) 1) $0.7x - 0.04 + 0.15x - 0.04 + 0.15x - 0.04$
- (down) 2) $0.7x - 0.04 + 0.15x - 0.04 + 0.15x - 0.04$
- (left) 3) $0.7x - 0.04 + 0.15x - 0.04 + 0.15x - 0.04$
- (right) 4) $0.7x - 0.04 + 0.15x - 0.04 + 0.18x - 0.04$

$$\begin{aligned}
 V[3,1] &= -0.04 + 0.95x(0.7x - 0.04 + 0.15x - 0.04 + 0.15x - 0.04) \\
 &= -0.04 + 0.95x - 0.04 \\
 &= 1.95x - 0.04 \\
 &= -0.078
 \end{aligned}$$

$V[3,1] = -0.078.$

$$V[3,2] = -0.04 + 0.95x \text{ max of}$$

- (up) 1) $0.7x - 0.04 + 0.15x - 0.04 + 0.15x - 0.04$
- (down) 2) $0.7x - 0.04 + 0.18x - 0.04 + 0.15x - 0.04$
- (left) 3) $0.7x - 0.04 + 0.15x - 0.04 + 0.15x - 0.04$
- (right) 4) $0.7x - 0.04 + 0.15x - 0.04 + 0.15x - 0.04$

$$\therefore V[3,2] = -0.04 + 0.95x(0.7x - 0.04 + 0.15x - 0.04 + 0.15x - 0.04)$$

$$\begin{aligned}
 \therefore U[3,2] &= -0.04 + 0.95x - 0.04 \\
 &= 1.95x - 0.04 \\
 &= -0.078 \\
 U[3,2] &= -0.078.
 \end{aligned}$$

$$\therefore U_2 = \begin{bmatrix} 0.708 & 1 & -1 \\ 0.459 & 0.614 & 0.227 \\ -0.078 & WALL & -0.078 \\ -0.078 & -0.078 & -0.078 \end{bmatrix}$$

The values of the grid match with output of the code.*

* $U[2,1] = 0.614$ in the dry run, but in the code it's $U[1,1] = 0.613$.

This is probably due to how Python handles & approximates floats].