Roll No 2022113005

we are given that on hitting a boundary then the agent remains on the same cell

constants are dd given as follows

step cost = $-0.04$

$P$(going direction chosen) = $0.7$ $P$

$P$(going in $\perp$ direction) = $0.15$ $\frac{(1-P)}{2}$

$\gamma$ (Discount factor) = $0.95$

given a state the next utility is calculated as:

$$U_{t+1}(I) = \max_A [C(I, A)$$
$$+ \gamma \sum_J P(I|I, A) \cdot U_t(I)]$$

| 0 | -1 | +1 |
|---|---|---|
| 0 | 0 | 0 |
| 0 | $\boxed{0}^{(w)}$ | 0 |
| 0 | 0 | 0 |

(Bellman update eqn)

$A$ = action that we choose in state $A$

$C(I.A)$ = cost of taking action

☆ $P(J|I, A)$ = probab. of reaching a state $J$ given that agent chooses action $A$ in state $I$.

For all cells and for all directions we calculate and consider maximum in one

<u>1st Iteration</u> : episode

1. $U(0,0)$    we can go up, down, left, right

a) Up = $-0.04 + 0.95(0.7 \times 0 + 0.15 \times (-1)$
$+ 0.15(0))$

= $-0.1825$

b) Down → $-0.04 + 0.95(0.7 \times 0 + 0.15 \times (-1)$
$+ 0.15(0))$.

= $-0.1825$

c) Right $= -0.04 + 0.95(0.7(-1) + 0.15 \times 0 + 0.15...)$

$$= -0.705$$

d) Left $= -0.04 + 0.95(0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0)$

$$= -0.04$$

Among these maximum is $,) \; U_1(0,0) = -0.04$

$$\text{left}$$

2. $U_1(1,0) =$

a) up $\rightarrow -0.04 + 0.95(0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0)$

$$= -0.04$$

b) Down $\rightarrow -0.04 + 0.95(0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0)$

$$= -0.04$$

c) Right $\rightarrow -0.04 + 0.95(0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0)$

$$= -0.04 \quad \text{left is same}$$

$$U_1(1,0) = argman(-0.04, -0.04, 0.04, -0.04)$$

$$U_1(1,0) = -0.04$$

3. $U_1(1,1)$

a) up $= -0.04 + 0.95(0.7 \times (-1) + 0.15(0) \times ?)$

$$= -0.705$$

b) Down $= -0.04 + 0.95(0 + 0 + 0)$

$$= 0.04$$

c) Right $= -0.04 + 0.95(0.7 \times 0 + 0.15 \times 0 + 0.15 \times \left(\frac{-1}{}\right) + 0.15 \times 0)$

$$= -0.1825$$

d) Left $=$ similarly, $-0.1825$

$) \; U_1(1,1) = 0.04$

$U_1(1,2)$ @

$U \to (-0.04) + 0.95(0.7 \times 1 + 0.15 \times 0$
$\qquad\qquad\qquad\qquad\qquad + 0.15 \times 0)$
$\qquad\qquad : +0.675$

$D \to (-0.04) + 0.95(0.9(0) + 0 + 0) \qquad -0.04$

$R \to (0.04) + 0.95(0.7(0) + 0.15(1) + 0)$
$\qquad = +0.1025$

$L \to 0.1025$ similarly

$\Rightarrow U_1(1,2) = 0.625$

$U_1(2,0)$ $\qquad U \to -0.04 + 0.95(0.7 \times 0 + 0.15 \times 0$
$\qquad\qquad\qquad\qquad\qquad\qquad + 0.15 \times 0)$
$\qquad\qquad : -0.04.$

$D \to -0.04 + 0.95(0.4 \times 0 + 0.15 \times 0 + 0.15 \times 0) = 0.04$

$L \to -0.04$ similarly

$R \to -0.04$ similarly

$\Rightarrow U_1(2,0) = -0.04$

$U_1(2,2)$ $\qquad U \to -0.04 + 0.95(0 + 0 + 0)$
$\qquad\qquad\qquad\qquad : -0.04$

$D \to -0.04$ similarly

$R \to -0.04$ similarly

$L \to -0.04$ similarly

$\Rightarrow U_1(2,2) = -0.04.$

7. $U_1(3,0) =$ $\qquad\qquad\qquad\qquad +0.95(0 + 0 + 0)$
$\qquad\qquad\qquad U \to -0.04 + 0 + 0 = -0.04$

$\qquad\qquad\qquad D \to -0.04$

$\qquad\qquad\qquad L \to -0.04$ $\Big\}$ similar

$\qquad\qquad\qquad R \to -0.04$

$\Rightarrow U_1(3,0) = -0.04$

8. $U_1(3,1) =$ $\qquad U \to -0.04 + 0.95(0 + 0 + 0)$
$\qquad\qquad\qquad\qquad = -0.04$

$\qquad\qquad\qquad D \to -0.04$

$\qquad\qquad\qquad L \to -0.04$

$\qquad\qquad\qquad R \to -0.04$

$\Rightarrow U_1(3,1) = -0.04$

9. $U_1(3, 2)$

$$U \to -0.04 + 0.95(0.7 \times 0 + 0.15 \times ... + 0.15 \times 0)$$

$$= -0.04$$

$$D, L, R = -0.04 \text{ similarly}$$

$\to U_1(3, 2) = -0.04$.

Now, the grid has been updated as follows.

| | | |
|---|---|---|
| $-0.04$ | $-1$ | $+1$ |
| $-0.04$ | $-0.04$ | $0.625$ |
| $-0.04$ | $0$ | $-0.04$ |
| $-0.04$ | $-0.04$ | $-0.04$ |

we see that wall
and the final
states remain
unchanged.

Iteration #2

We will use the Bellman update eqⁿ again

$$U_2(a, b) = \text{argmax} \left[ \text{step cost} + \gamma \sum_J P(J \mid (a, b), A) \cdot U_{(J)} \right]$$

1. $U_2(0, 0) \to U_p = -0.04 + 0.95((-0.04 \times 0) + (-1)(0.15))$

$$= -0.2148.$$

$$\text{Down} = (0.04) + 0.95((-0.04 \times 0.7) + -1 \times 0.15 + (0.04)(0.15))$$

$$= -0.2148.$$

$$\text{Right} = (-0.04) + 0.95((-0.7 \times (-1)) + (0.04 \times 0.15) + (0.04 \times 0.15))$$

$$= -0.71604$$

$$\text{Left} = (-0.04) + 0.95((0.7 \times (-0.04)) + -0.04 \times 0.15 + -0.04 \times 0.15)$$

$$= -0.078.$$

$U_2(0, 0) = \text{max of all these} = -0.078$

1)$_2$ (1,0).

$U \to -0.04 + 0.95 ( 0.7 x - 0.04 + (-0.04 \times 0.15)$
$\qquad\qquad + 1 - 0.04 \times 0.15))$

$\qquad = -0.078$

$D \to -0.04 + 0.95 [(0.7x - 0.04) + (-0.04 \times 0.15)$
$\qquad\qquad\qquad + (-0.04 \times 0.15)]$

$\qquad = -0.078$

$P \to -0.04 \times 0.95 [0.7 \times 0.4) + (0.15 \times (-0.04)) +$
$\qquad\qquad (0.15 \times (0.04))]$

$\qquad = -0.678$

$L \to$ ~~countnity~~ $= -0.04 \times 0.95 [(0.7 \times (-0.04) +$
$\qquad\qquad\qquad\qquad\qquad 0.15 \times (0.04) +$
$\qquad\qquad\qquad\qquad\qquad 0.15 \times -0.04)]$

$\qquad = -0.078$

Considering man of all these, $U_2 (1,0) = -0.078$.

3. $U_2 (1,1) =$

$U \to -0.04 + 0.95 [(0.7 \times (1)) + (0.15 \times -0.04 \neq$
$\qquad\qquad\qquad (0.15 \times 0.625)]$.

$\qquad = -0.6216375$

$D \to -0.64 + 0.95 [(0.625 \times 0.7) + (-1 \times 0.15)$
$\qquad\qquad\qquad + (-0.04 \times 0.15)]$

$\qquad = 0.227425$

Right, similarly $= 0.277425$.

$L \to -0.04 + 0.95 [(-0.04 \times 0.7) + (-1 \times 0.15)$
$\qquad\qquad\qquad + (-0.04 \times 0.15)]$

$\qquad = -0.2148$.

$U_2 (1,1) = $ man $\{ -0.6236, 0.227425, -0.2148\}$

$\qquad = 0.227425$

4. $U_2 (1,2)$. $\qquad U \to -0.04 + 0.95 [(0.7 \times 1) + (0.15 \times 0.625)$
$\qquad\qquad\qquad\qquad\qquad + (0.15 \times -0.04)]$

$\qquad\qquad\qquad\qquad = 0.7683625$

$D \rightarrow -0.04 + 0.95 [(0.7 \times 0.04) + (6.625 \times 0.15) + 0.15$

$= 0.0167625$

$R \rightarrow -0.04 + 0.95 [(0.7 \times 6.625) + (0.15 \times 1) + (0.15 \times -0$

$= 0.512425$

$L \rightarrow -0.04 + 0.95 [(0.7 \times -0.04) + (6.15 \times 1) + (6.15 \times -0$

$= 0.0902$

$\Rightarrow U_2(1,2) = 0.7083625$

5. $U_2(2,0)$  $\quad U \rightarrow -0.04 + 0.95 (-0.04 \times 0.7 +$

$0 \times 0.15 (\times 0.(50))$

$= -0.04 + 0.95 (-0.04)$

$= -6.078$

$D \rightarrow -0.048 + 6.95 (-0.04) = -0.078$

Similarly $L, R = -0.078$

$\Rightarrow U_2(2,0) = -0.078$

6. $U_2(2,2)$

$U \rightarrow -0.04 + 0.95 [(6.625 \times 0.7) + (0.04 \times 0.15) + (0.04 \times 0.15)$

$= 0.364225$

$D \rightarrow -0.04 + 0.95 [(-0.04 \times 0.7) + (6.93 \times -0.04) + (-0.04 \times$

$= -0.078$

$L \rightarrow = 0.04 + 0.95 [(0.7 \times -0.04) + (0.15 \times 0.625)$

$+ (6.15 \times 0.625)]$

$= 0.0167625$

Similarly, $R \rightarrow 0.0167625$

So $U_1(2,2) = $ Max over all directions

$= 0.364225$

$U_2 (3,0)$

$$U \rightarrow -0.04 + 0.95[(0.7 \times -0.04) + (0.15 \times -0.04) + (0.15 \times -0.04)]$$

$$= -0.078$$

Similarly

same spot ↑

$R \rightarrow -0.078$

$$L, D \rightarrow -0.04 + 0.95[(-0.7 \times -0.04) + \begin{pmatrix} (0.04 \times 0.15) \\ + (0.04 \times 0.15) \end{pmatrix}]$$

$\Rightarrow U_2 (3,0) = -0.078.$

8. $U_2 (3,1), U_2 (3,2)$

$$U, L, D, R = -0.078.$$

$$-0.04 + 0.95[(-0.7 \times -0.04) + (-0.04 \times 0.15) + (0.04 \times 0.15)]$$

( Because on pumping into walls, it returns to same place which has same value as adjacent cells

$\Rightarrow U_2 (3,1), U_2 (3,2) = -0.078.$

After iteration 2

| -0.078 | -1 | +1 |
|--------|--------|--------|
| -0.078 | 0.2274 | 0.768 |
| -0.078 | 0 | 0.364225 |
| -0.078 | -0.078 | -0.078 |

As we can see, wall and reward/punishment final states remain unchanged. On comparison with the result generated in the second iteration with the computer is the same as what we calculated with the computer program.