

Biographical Learning: Exercise RL3

RL3.1: Bellman equation

Policy

- actions: right, up, left, down
- For each action a , policy: $\pi(a|s) = p(A_t=a|S_t=s) = 0.25$
- every action has same probability

Probability $p(s', r | s, a)$

$s = (1, 1)$:

$$p(s', r | s = (1, 1), a = \text{'right'}) = \begin{cases} 1 & \text{for } s' = (1, 2) \text{ and } r = 0 \\ 0 & \text{for all other values of } s' \text{ and } r \end{cases}$$

$$p(s', r | s = (1, 1), a = \text{'up'}) = \begin{cases} 1 & \text{for } s' = (0, 1) \text{ and } r = 0 \\ 0 & \text{for ...} \end{cases}$$

$$p(s', r | s = (1, 1), a = \text{'left'}) = \begin{cases} 1 & \text{for } s' = (1, 0) \text{ and } r = 0 \\ 0 & \text{for ...} \end{cases}$$

$$p(s', r | s = (1, 1), a = \text{'down'}) = \begin{cases} 1 & \text{for } s' = (2, 1) \text{ and } r = 0 \\ 0 & \text{for ...} \end{cases}$$

$s = (0, 2)$:

$$p(s', r | s = (0, 2), a = \text{'right'}) = \begin{cases} 1 & \text{for } s' = (0, 3) \text{ and } r = 0 \\ 0 & \text{for ...} \end{cases}$$

$$p(s', r | s = (0, 2), a = \text{'up'}) = \begin{cases} 1 & \text{for } s' = (0, 1) \text{ and } r = -1 \\ 0 & \text{for ...} \end{cases}$$

$$p(s', r | s = (0, 2), a = \text{'left'}) = \begin{cases} 1 & \text{for } s' = (0, 1) \text{ and } r = 0 \\ 0 & \text{for ...} \end{cases}$$

$$p(s', r | s = (0, 2), a = \text{'down'}) = \begin{cases} 1 & \text{for } s' = (1, 2) \text{ and } r = 0 \\ 0 & \text{for ...} \end{cases}$$

$s = A$:

$$p(s', r | s = A, a = \begin{matrix} \text{'right' or} \\ \text{'up' or} \\ \text{'left' or} \\ \text{'down'} \end{matrix}) = \begin{cases} 1 & \text{for } s' = (1, 1) \text{ and } r = 1 \\ 0 & \text{for ...} \end{cases}$$

Bellman equation

$$v_{\pi}(s) = \sum_a \pi(a|s) \sum_{s' \in S} p(s', r|s, a) [r + \gamma v_{\pi}(s')],$$

for all $s \in S$

$s = (1, 1)$:

$$\begin{aligned} v_{\pi}(s=(1, 1)) &= 0.25 \cdot (0 + \gamma v_{\pi}(s'=1, 2)) \\ &\quad + 0.25 \cdot (0 + \gamma v_{\pi}(s'=0, 1)) \\ &\quad + 0.25 \cdot (0 + \gamma v_{\pi}(s'=1, 0)) \\ &\quad + 0.25 \cdot (0 + \gamma v_{\pi}(s'=2, 1)) \end{aligned}$$

$$\begin{aligned} \text{with } \gamma = 0.9 | &= 0.25 \cdot 0.9 (2.3 + 0.2 + 1.5 + 0.8) \\ &= 2.9925 // \end{aligned}$$

$s = (0, 2)$:

$$\begin{aligned} v_{\pi}(s=(0, 2)) &= 0.25 \cdot (0 + \gamma v_{\pi}(s'=0, 3)) \\ &\quad + 0.25 \cdot (-1 + \gamma v_{\pi}(s'=0, 2)) \\ &\quad + 0.25 \cdot (0 + \gamma v_{\pi}(s'=0, 1)) \\ &\quad + 0.25 \cdot (0 + \gamma v_{\pi}(s'=1, 2)) \\ &= 0.25 \cdot 0.9 (5.3 + 4.9 + 1.1 + 2.3) - 0.25 \\ &= 4.53 // \end{aligned}$$

$s = A$:

$$\begin{aligned} v_{\pi}(s=A) &= 4 \cdot 0.25 \cdot (10 + \gamma v_{\pi}(s'=(5, 1))) \\ &= 10 + 0.9 \cdot (-1.3) \\ &= 7.13 // \end{aligned}$$

RLJ. 2: Compute the values of states using
for e.g.

Policy 1

6	7	8	9
7	8	9	10

Policy 2

6	7	8	9
5	9	6	10

$$v(s) = 0.5 \cdot (-1) + 0.5 \cdot 4 + 0.5 \cdot (-1) + 0.5 \cdot 10 \\ = -0.5 + 2 - 0.5 + 5 = 6$$

→ Policy 1 is better: each state in Policy 1 has
than
a higher value or similar value as
the corresponding state in Policy 2

Policy 3

a	d	8	9
=	=		
5	6		
b	c		
=	=	9	10
5	6		

$$c = 0.5 \cdot (-1 + 3) + 0.5 \cdot (-1 + 8)$$

$$= 0.5 \cdot (3 + 7)$$

$$s = c - 1$$

$$d = 0.5 \cdot (-1 + 8) + 0.5 \cdot (-1 + 0.5 \cdot (s + 2))$$

$$= 0.5 \cdot (6 + 0.5 \cdot (s + 2))$$

$$= 0.5 \cdot (6 + 0.5s + 2.5)$$

$$= 0.5 \cdot (9.5 + 0.5s)$$

$$a = d - 1$$

c in b:

$$b = 0.5 \cdot (6 + 7) - 1 = 0.5 \cdot 13 - 1 = 0.5 \cdot 12 = 6$$

$$0.5s = 2.5 \quad s = 5$$

$$b = 5$$

b in c:

$$b = c - 1 \quad c = 5$$

$$c = 5 + 1 \quad c = 6$$

$$c = 6 \quad c = 6$$

b in d :

$$d = 0.5(0.5 + 0.5b)$$

$$= 4.75 + 0.25b$$

$$= 4.75 + 0.25 \cdot 5$$

$$= 4.75 + 1.25$$

$\Rightarrow b$

d in a :

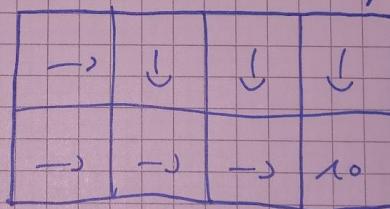
$$a = d - 1$$

$$= 6 - 1$$

$$= 5$$

Optimal policy

- Policy is best out of the three
- It is also optimal: optimizes number of steps to reach the goal but other equivalent policies are there, e.g.



Bellmann optimality equation

v^*	6	2	8	3
	7	8	9	10

$$\max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v^*(s')]$$

- maximal for policies with movements in row 0 either right or down and movement in row 1 right
- If u with policy 1 but not policy 3
 $(s=(1,1) \text{ has 0.5 probability of going left})$