

Actions : { Left, Right, Up, Down, Nothing }

Immediate rewards, in this example depend only on states, not actions.

Example : $R([0,0], \text{Right}) = R([0,0], \text{left}) =$

$R([0,0], \text{Up}) = R([0,0], \text{Down}) = 0$

$R([0,1], \text{Right}) = R([0,1], \text{left}) =$

$R([0,1], \text{Up}) = R([0,1], \text{Down}) = -5$

(0,0)	(0,1)	(0,2)
0	-5	+10
(1,0)	(1,1)	(1,2)
0	0	0
(2,0)	(2,1)	(2,2)
0	0	0

The transition probabilities for every state and action are given in the description.

State Space on rewards.

Initial Policy π^0

$\gamma = 0.9$

State	0,0	0,1	0,2	1,0	1,1	1,2	2,0	2,1	2,2
$\pi^0(\text{state})$	Right	Right	Right	Right	Right	Right	Right	Right	Right

Evaluating Policy π^0

Initial values

S	$V^0(S)$	$V^1(S)$
(0,0)	0	$R([0,0], \pi^0([0,0])) + \gamma [T([0,0], \pi^0([0,0]), [0,1]) \cdot V^0([0,1]) + T([0,0], \pi^0([0,0]), [1,0]) \cdot V^0([1,0])]$ $= 0 + 0.9 \times [0.9 \times 0 + 0.1 \times 0] = 0$
(0,1)	0	$R([0,1], \pi^0([0,1])) + \gamma \times [T([0,1], \pi^0([0,1]), [0,2]) \cdot V^0([0,2]) + T([0,1], \pi^0([0,1]), [0,0]) \cdot V^0([0,0]) + T([0,1], \pi^0([0,1]), [1,1]) \cdot V^0([1,1])]$ $= -5 + 0.9 \times [0.9 \times 0 + 0.05 \times 0 + 0.05 \times 0]$ $= -5$
(0,2)	0	$R([0,2], \pi^0([0,2])) + \gamma \times [T([0,2], \pi^0([0,2]), [0,2]) \cdot V^0([0,2])]$ $= +10 + 0.9 \times [1 \times 0] = +10$
(1,0)	0	