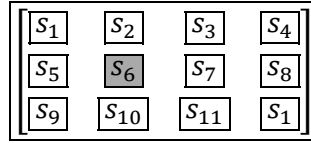


GRID WORLD PROBLEM/GAME



Bellman Eq

$$V_{i+1}(s) = \max_{a \in \mathbb{A}} \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_i(s')]$$

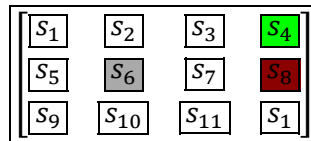
Value iteration algorithm

Algorithm
Start with $V_0(s) = 0 \forall s \in \{s_1, s_2, \dots, s_{12}\}$ For $i=1:H$ $\forall s \in \{s_1, s_2, \dots, s_{12}\}$: $V_{i+1}(s) \leftarrow \max_{a \in \mathbb{A}} \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_i(s')]$ $\pi_{i+1}(s) \leftarrow \arg \max_{a \in \mathbb{A}} \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_i(s')]$
This is called value-update or Bellman-update

$V_i^*(s)$ = expected sum of rewards accumulated starting from state s , acting optimally for i steps

$\pi_i^*(s)$ = optimal action when in state s and getting to act for i steps

THE GRID WORLD PROBLEM



i	$V_i(s_1)$	$V_i(s_2)$	$V_i(s_3)$	$V_i(s_4)$	$V_i(s_5)$	$V_i(s_6)$	$V_i(s_7)$	$V_i(s_8)$	$V_i(s_9)$	$V_i(s_{10})$	$V_i(s_{11})$	$V_i(s_{12})$
1	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	1	0	0	0	-1	0	0	0	0

For $i=2,3,\dots,H$ for the terminal states you write the rewards directly

(w/o computing the value using Bellman Eq)

$$R(s_4, a, s_4) = +1$$

$$R(s_8, a, s_8) = -1$$

When $i=3$:

i	$V_i(s_1)$	$V_i(s_2)$	$V_i(s_3)$	$V_i(s_4)$	$V_i(s_5)$	$V_i(s_6)$	$V_i(s_7)$	$V_i(s_8)$	$V_i(s_9)$	$V_i(s_{10})$	$V_i(s_{11})$	$V_i(s_{12})$
1	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	1	0	0	0	-1	0	0	0	0
3				1				-1				

$$V_{i+1}(s) \leftarrow \max_{a \in \mathbb{A}} \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_i(s')]$$

$\Rightarrow T(s_3, a, s_4)[R(s_3, a, s_4) + (0.9)V_2(s_4)] +$ $\Rightarrow: \uparrow T(s_3, a, s_3)[R(s_3, a, s_3) + (0.9)V_2(s_3)] +$ $\downarrow T(s_3, a, s_7)[R(s_3, a, s_7) + (0.9)V_2(s_7)]$	$\Rightarrow 0.8[0 + (0.9)(1)] +$ $\Rightarrow: \uparrow 0.1[0 + (0.9)(0)] +$ $\downarrow 0.1[0 + (0.9)(0)]$	0.72
$\uparrow T(s_3, a, s_3)[R(s_3, a, s_3) + (0.9)V_2(s_3)] +$ $\uparrow: \Rightarrow T(s_3, a, s_4)[R(s_3, a, s_4) + (0.9)V_2(s_4)] +$ $\leftarrow T(s_3, a, s_2)[R(s_3, a, s_2) + (0.9)V_2(s_2)]$	$\uparrow 0.8[0 + (0.9)(0)] +$ $\uparrow: \Rightarrow 0.1[0 + (0.9)(1)] +$ $\leftarrow 0.1[0 + (0.9)(0)]$	0.09
$\leftarrow T(s_3, a, s_2)[R(s_3, a, s_2) + (0.9)V_2(s_2)] +$ $\leftarrow: \uparrow T(s_3, a, s_3)[R(s_3, a, s_3) + (0.9)V_2(s_3)] +$ $\downarrow T(s_3, a, s_7)[R(s_3, a, s_7) + (0.9)V_2(s_7)]$	$\leftarrow 0.8[0 + (0.9)(0)] +$ $\leftarrow: \uparrow 0.1[0 + (0.9)(0)] +$ $\downarrow 0.1[0 + (0.9)(0)]$	0
$\downarrow T(s_3, a, s_7)[R(s_3, a, s_7) + (0.9)V_2(s_7)] +$ $\downarrow: \leftarrow T(s_3, a, s_2)[R(s_3, a, s_2) + (0.9)V_2(s_2)] +$ $\Rightarrow T(s_3, a, s_4)[R(s_3, a, s_4) + (0.9)V_2(s_4)]$	$\downarrow 0.8[0 + (0.9)(0)] +$ $\downarrow: \leftarrow 0.1[0 + (0.9)(0)] +$ $\Rightarrow 0.1[0 + (0.9)(1)]$	0.09

And take my word for it, the values of the other states are zero.

When i=4:

i	$V_i(s_1)$	$V_i(s_2)$	$V_i(s_3)$	$V_i(s_4)$	$V_i(s_5)$	$V_i(s_6)$	$V_i(s_7)$	$V_i(s_8)$	$V_i(s_9)$	$V_i(s_{10})$	$V_i(s_{11})$	$V_i(s_{12})$
1	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	1	0	0	0	-1	0	0	0	0
3	0	0	0.72	1	0	0	0	-1	0	0	0	0
4				1		0		-1				

$$V_{i+1}(s) \leftarrow \max_{a \in \mathbb{A}} \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_i(s')]$$

For s=2

$\Rightarrow T(s_2, a, s_3)[R(s_2, a, s_3) + (0.9)V_3(s_3)] +$ $\Rightarrow: \uparrow T(s_2, a, s_2)[R(s_2, a, s_2) + (0.9)V_3(s_2)] +$ $\downarrow T(s_2, a, s_2)[R(s_2, a, s_2) + (0.9)V_3(s_2)]$	$\Rightarrow 0.8[0 + (0.9)(0.72)] +$ $\Rightarrow: \uparrow 0.1[0 + (0.9)(0)] +$ $\downarrow 0.1[0 + (0.9)(0)]$	0.5184 max
$\uparrow T(s_2, a, s_2)[R(s_2, a, s_2) + (0.9)V_3(s_2)] +$ $\uparrow: \Rightarrow T(s_2, a, s_3)[R(s_2, a, s_3) + (0.9)V_3(s_3)] +$ $\Leftarrow T(s_2, a, s_1)[R(s_2, a, s_1) + (0.9)V_3(s_1)]$	$\uparrow 0.8[0 + (0.9)(0)] +$ $\uparrow: \Rightarrow 0.1[0 + (0.9)(0.72)] +$ $\Leftarrow 0.1[0 + (0.9)(0)]$	0.0648
$\Leftarrow T(s_2, a, s_1)[R(s_2, a, s_1) + (0.9)V_3(s_1)] +$ $\Leftarrow: \uparrow T(s_2, a, s_2)[R(s_2, a, s_2) + (0.9)V_3(s_2)] +$ $\downarrow T(s_2, a, s_2)[R(s_2, a, s_2) + (0.9)V_3(s_2)]$	$\Leftarrow 0.8[0 + (0.9)(0)] +$ $\Leftarrow: \uparrow 0.1[0 + (0.9)(0)] +$ $\downarrow 0.1[0 + (0.9)(0)]$	0
$\downarrow T(s_2, a, s_2)[R(s_2, a, s_2) + (0.9)V_3(s_2)] +$ $\downarrow: \Leftarrow T(s_2, a, s_1)[R(s_2, a, s_1) + (0.9)V_3(s_1)] +$ $\Rightarrow T(s_2, a, s_3)[R(s_2, a, s_3) + (0.9)V_3(s_3)]$	$\downarrow 0.8[0 + (0.9)(0)] +$ $\downarrow: \Leftarrow 0.1[0 + (0.9)(0)] +$ $\Rightarrow 0.1[0 + (0.9)(0.72)]$	0.0648

For s=3

$\Rightarrow T(s_3, a, s_4)[R(s_3, a, s_4) + (0.9)V_3(s_4)] +$ $\Rightarrow: \uparrow T(s_3, a, s_3)[R(s_3, a, s_3) + (0.9)V_3(s_3)] +$ $\downarrow T(s_3, a, s_7)[R(s_3, a, s_7) + (0.9)V_3(s_7)]$	$\Rightarrow 0.8[0 + (0.9)(1)] +$ $\Rightarrow: \uparrow 0.1[0 + (0.9)(0.72)] +$ $\downarrow 0.1[0 + (0.9)(0)]$	0.7848 max
$\uparrow T(s_3, a, s_3)[R(s_3, a, s_3) + (0.9)V_3(s_3)] +$ $\uparrow: \Rightarrow T(s_3, a, s_4)[R(s_3, a, s_4) + (0.9)V_3(s_4)] +$ $\Leftarrow T(s_3, a, s_2)[R(s_3, a, s_2) + (0.9)V_3(s_2)]$	$\uparrow 0.8[0 + (0.9)(0.72)] +$ $\uparrow: \Rightarrow 0.1[0 + (0.9)(1)] +$ $\Leftarrow 0.1[0 + (0.9)(0)]$	0.6084
$\Leftarrow T(s_3, a, s_2)[R(s_3, a, s_2) + (0.9)V_3(s_2)] +$ $\Leftarrow: \uparrow T(s_3, a, s_3)[R(s_3, a, s_3) + (0.9)V_3(s_3)] +$ $\downarrow T(s_3, a, s_7)[R(s_3, a, s_7) + (0.9)V_3(s_7)]$	$\Leftarrow 0.8[0 + (0.9)(0)] +$ $\Leftarrow: \uparrow 0.1[0 + (0.9)(0.72)] +$ $\downarrow 0.1[0 + (0.9)(0)]$	0.0648
$\downarrow T(s_3, a, s_7)[R(s_3, a, s_7) + (0.9)V_3(s_7)] +$ $\downarrow: \Leftarrow T(s_3, a, s_2)[R(s_3, a, s_2) + (0.9)V_3(s_2)] +$ $\Rightarrow T(s_3, a, s_4)[R(s_3, a, s_4) + (0.9)V_3(s_4)]$	$\downarrow 0.8[0 + (0.9)(0)] +$ $\downarrow: \Leftarrow 0.1[0 + (0.9)(0)] +$ $\Rightarrow 0.1[0 + (0.9)(1)]$	0.09

For $s=7$

$\Rightarrow T(s_7, a, s_8)[R(s_7, a, s_8) + (0.9)V_3(s_8)] +$ $\Rightarrow: \uparrow T(s_7, a, s_3)[R(s_7, a, s_3) + (0.9)V_3(s_3)] +$ $\downarrow T(s_7, a, s_{11})[R(s_7, a, s_{11}) + (0.9)V_3(s_{11})]$	$\Rightarrow 0.8[0 + (0.9)(-1)] +$ $\Rightarrow: \uparrow 0.1[0 + (0.9)(0.72)] +$ $\downarrow 0.1[0 + (0.9)(0)]$	-0.2016
$\uparrow T(s_7, a, s_3)[R(s_7, a, s_3) + (0.9)V_3(s_3)] +$ $\uparrow: \Rightarrow T(s_7, a, s_8)[R(s_7, a, s_8) + (0.9)V_3(s_8)] +$ $\Leftarrow T(s_7, a, s_7)[R(s_7, a, s_7) + (0.9)V_3(s_7)]$	$\uparrow 0.8[0 + (0.9)(0.72)] +$ $\uparrow: \Rightarrow 0.1[0 + (0.9)(-1)] +$ $\Leftarrow 0.1[0 + (0.9)(0)]$	0.4284 max
$\Leftarrow T(s_7, a, s_7)[R(s_7, a, s_7) + (0.9)V_3(s_7)] +$ $\Leftarrow: \uparrow T(s_7, a, s_3)[R(s_7, a, s_3) + (0.9)V_3(s_3)] +$ $\downarrow T(s_7, a, s_{11})[R(s_7, a, s_{11}) + (0.9)V_3(s_{11})]$	$\Leftarrow 0.8[0 + (0.9)(0)] +$ $\Leftarrow: \uparrow 0.1[0 + (0.9)(0.72)] +$ $\downarrow 0.1[0 + (0.9)(0)]$	0.0648
$\downarrow T(s_7, a, s_{11})[R(s_7, a, s_{11}) + (0.9)V_3(s_{11})] +$ $\downarrow: \Leftarrow T(s_7, a, s_7)[R(s_7, a, s_7) + (0.9)V_3(s_7)] +$ $\Rightarrow T(s_7, a, s_8)[R(s_7, a, s_8) + (0.9)V_3(s_8)]$	$\downarrow 0.8[0 + (0.9)(0)] +$ $\downarrow: \Leftarrow 0.1[0 + (0.9)(0)] +$ $\Rightarrow 0.1[0 + (0.9)(-1)]$	-0.09

i	$V_i(s_1)$	$V_i(s_2)$	$V_i(s_3)$	$V_i(s_4)$	$V_i(s_5)$	$V_i(s_6)$	$V_i(s_7)$	$V_i(s_8)$	$V_i(s_9)$	$V_i(s_{10})$	$V_i(s_{11})$	$V_i(s_{12})$
1	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	1	0	0	0	-1	0	0	0	0
3	0	0	0.72	1	0	0	0	-1	0	0	0	0
4	0	0.5184	0.7848	1	0	0	0.4284	-1	0	0	0	0

$s_prime_list_generator(state, action)$

For each state-action pair there is a corresponding s-prime-list