

## Value Iteration:

Finding optimal policies of MDPs when we have transition probabilities and Reward functions

Bellman Equation gives us a recursive definition of the optimal value

$$V^*(s) = \max_{a \in A} \sum_{s' \in S} P(s'|s, a) [R(s, a, s') + \gamma V^*(s')]$$

Dynamic programming

In Value iteration convert this equation into an update rule (recursive def.)

Iterative process → Iterate Bellman update until convergence  
Refine  $V^*(s)$

- Start with  $V_0(s) = 0$  for all states  $s$
- Iterate Bellman update until convergence:

$$V_i(s) \leftarrow \max_{a \in A} \sum_{s' \in S} P(s'|s, a) [R(s, a, s') + \gamma V_i(s')]$$

iteration number

as  $i \rightarrow \infty$

we get  $V^*(s)$

Example MDP

living reward: every time step in env is not penalized  
Or given a reward

3				$+1$	reward given in terminal state
2				$-1$	
1	2	3	4		living reward = 0

$\gamma = 0.9$

transition where 80% correct action is executed ← noise = 0.2

probability

0<sup>th</sup> iteration :-

- ① Initializing the estimate of optimal value function  
start with  $V_0(S) = 0$

3	0	0	0	0
2	0	0	0	0
1	0	0	0	0

$V_0 \uparrow$

Plug  $V_0$  into Bellman update rule, we get

$$V_{t+1}(s) \leftarrow \max_{a \in A} \sum_{s' \in S} P(s'|s, a) [R(s, a, s') + \gamma V_t(s')]$$

3	0	0	0	+1
2	0	0	0	-1
1	0	0	0	0

$V_1 \uparrow$

$V_2 \downarrow$

3		?	+1
2			-1
1			

$V_2(3,4) \leftarrow 1$  } terminal  
 $V_2(2,4) \leftarrow -1$  } states  
 no new rewards  
 since  $(3,4)$  &  $(2,4)$  are  
 terminal states  
 so set  $V_2(3,4)$  &  $V_2(2,4)$

$$v_2(\langle 3,3 \rangle) \leftarrow \sum_{s' \in S} P(s'| \langle 3,3 \rangle, \text{right}) [r(\langle 3,3 \rangle, \text{right}, s') + 0.9 v_1(s')]$$

$$\leftarrow 0.8 [0 + 0.9 \times 1] + 0.1 (0 + 0.9 \times 0) + 0.1 (0 + 0.9 \times 0)$$

$$= 0.72.$$

At  $\langle 3,3 \rangle$  move towards right -

end up at  $\langle 3,4 \rangle$  or  $\langle 2,3 \rangle$

At  $\langle 3,3 \rangle$  move towards left

end up at  $\langle 3,2 \rangle$  or  $\langle 2,3 \rangle$

At  $\langle 3,3 \rangle$  move up

end up at  $\langle 3,3 \rangle$

Consider moving right from  $\langle 3,3 \rangle$

optimal action = 'right'

3	0	0	0	+1
2	0	0	0	-1
1	0	0	0	0

0	0	0.72	+1
0	X	0	-1
0	0	0	0

$v_3$	0	0.52	0.78	+1
$\hookrightarrow$	0	X	0.43	-1
	0	0	0	0

- \* Information propagates outward from terminal states
- \* Eventually all states have correct value estimates

Bellman Eq, is a property of optimal value function  
 Use Bellman function as an update rule and we end up with optimal value function

Value Iteration converges to optimal value function

Pseudocode of Value Iteration

- ① start with  $V_0(s) = 0$
- ② Iterate until convergence:

$$V_{i+1}(s) \leftarrow \max_{a \in A} \sum_{s' \in S} P(s'|s, a) [R(s, a, s') + \gamma V_i(s')]$$

$$V_3(s) \leftarrow \max_{a \in A} \sum_{s' \in S} P(s'|s, a) [R(s, a, s') + \gamma V_2(s')]$$

## Value Iteration:

- ↳ In value iteration, because of the max operation, the equations are not linear anymore.
- ↳ Use an iterative procedure to solve them.
- ↳ Initialize the utility of every state = 0 ( $\delta = 0.5$ )
- ↳ loop through states using Bellman equation

$$v(s) = r(s) + \gamma \max_a \left[ \sum_{s'} P(s'|s,a) v(s') \right]$$

when  $s=0$

$$\begin{aligned} v(0) &= r(0) + \gamma \max_a \left[ \sum_{s'} P(s'|s=0, a) v(s') \right] \\ &= \gamma(0) + \gamma \max \left[ \begin{array}{l} \sum_{s'} P(s'|s=0, a=up) v(s'), \\ \sum_{s'} P(s'|s=0, a=left) v(s'), \\ \sum_{s'} P(s'|s=0, a=down) v(s'), \\ \sum_{s'} P(s|s=0, a=right) v(s') \end{array} \right] \\ &= \gamma(0) + \gamma \max \left[ \begin{array}{l} 0.8 v(0) + 0.1 v(0) + 0.1 v(1) \\ 0.8 v(0) + 0.1 v(4) + 0.1 v(0) \\ 0.8 v(4) + 0.1 v(1) + 0.1 v(0) \\ 0.8 v(1) + 0.1 v(0) + 0.1 v(4) \end{array} \right] \\ &\approx -0.04 + 0.5 \times \max \left[ \begin{array}{l} 0 \\ 0 \\ 0 \\ 0 \end{array} \right] = -0.04 \end{aligned}$$

in-place procedure, when we see  $v(0)$  then  
 $v(0) = -0.04$  instead of 0 [prev value of  $v(0)$ ]

Now when  $s=1$ ,

$$v(1) = \gamma(1) + \gamma \max \left[ \begin{array}{l} 0.8 v(1) + 0.1 v(0) + 0.1 v(2) \\ 0.8 v(0) + 0.1 v(1) + 0.1 v(1) \\ 0.8 v(1) + 0.1 v(2) + 0.1 v(0) \\ 0.8 v(2) + 0.1 v(1) + 0.1 v(1) \end{array} \right]$$

$$= -0.04 + 0.5 \times \max \left[ \begin{array}{l} 0.1 \times (-0.04) \\ 0.8 \times (-0.04) \\ 0.1 \times (-0.04) \\ 0 \end{array} \right] = -0.04$$

Repeat for states 2 - 11, we get this utility

$s=0$	$s=1$	$s=2$	$s=3$
$v = -0.04$	$v = -0.04$	$v = -0.04$	$v = 1.0$
$s=4$	$s=5$	$s=6$	$s=7$
$v = -0.04$		$v = -0.04$	$v = -1.0$
$s=8$	$s=9$	$s=10$	$s=11$
$v = -0.04$	$v = -0.04$	$v = -0.042$	$v = -0.0421$



utility Value after the first iteration

utility value after 11 iterations

$s=0$	$s=1$	$s=2$	$s=3$
$v = 0.0897$	$v = 0.3147$	$v = 0.8093$	$v = 1.999$
$s=4$	<del><math>s=5</math></del>	$s=6$	$s=7$
$v = -0.0046$		$v = 0.1935$	$v = -1.9990$
$s=8$	$s=9$	$s=10$	$s=11$
$v = -0.0456$	$v = -0.0301$	$v = 0.0324$	$v = -0.0698$

Get optimal policy :-

$$\pi(s) \leftarrow \operatorname{argmax}_a \left[ \sum_{s'} p(s'|s,a) \gamma(s') \right]$$

get optimal policy by applying equation  
for each state