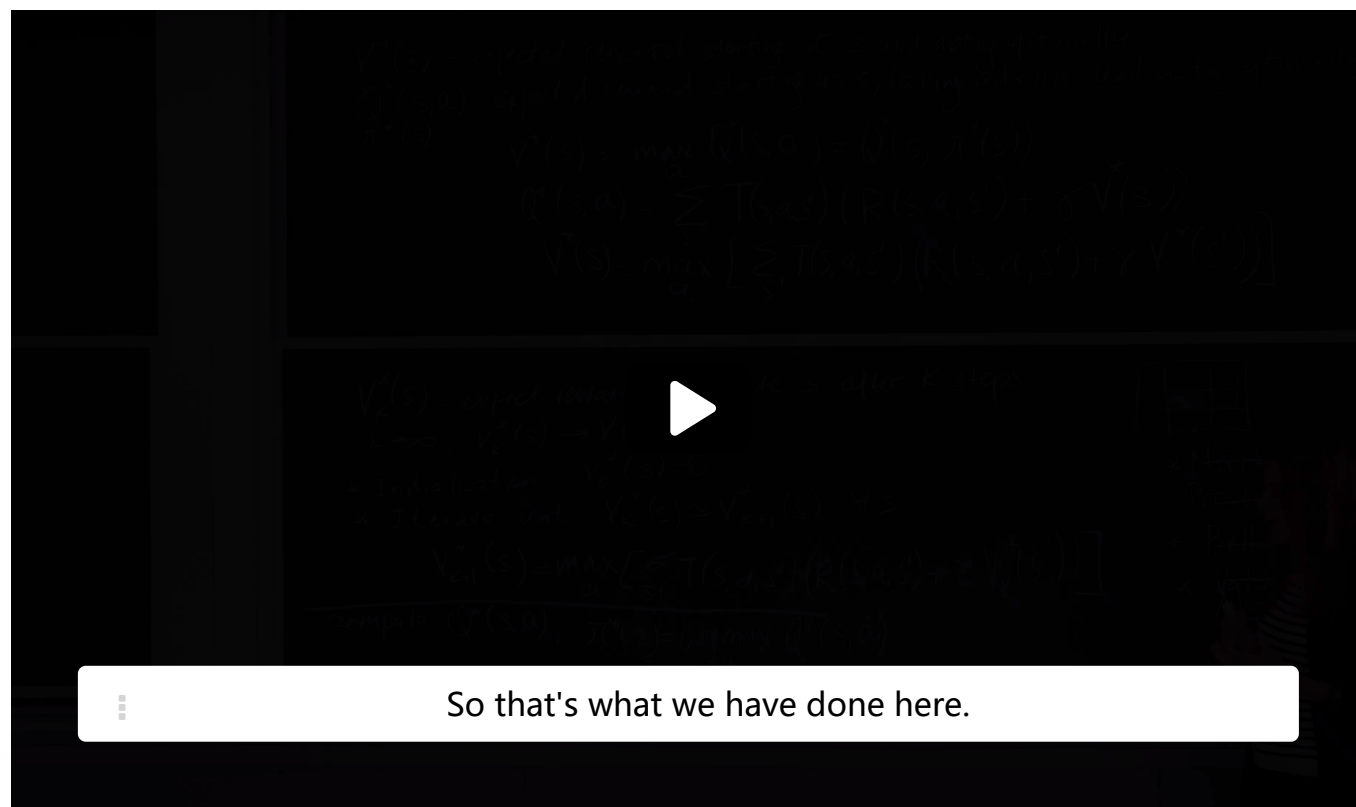


## 7. Value Iteration

### Value Iteration



a Q star S A. What it says for a given state is go select, check all the Q star S a's and select

the action which maximizes it.

So what eventually I go to after this algorithm?

After this algorithm, when it converged, I computed the Q values, and then I computed the policy.

And now I know how to act in my MDP.

So that's what we have done here.



End of transcript. Skip to the start.

#### Video

[Download video file](#)

#### Transcripts

[Download SubRip \(.srt\) file](#)

[Download Text \(.txt\) file](#)

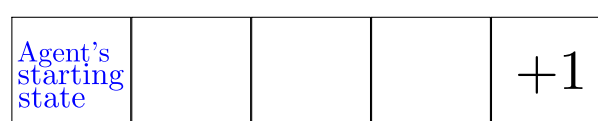


Recall from lecture the **value iteration update rule** :

$$V_{k+1}^*(s) = \max_a \left[ \sum_{s'} T(s, a, s') (R(s, a, s') + \gamma V_k^*(s')) \right],$$

where  $V_k^*(s)$  is the expected reward from state  $s$  after acting optimally for  $k$  steps.

Recall the example discussed in the lecture.



An agent is trying to navigate a one-dimensional grid consisting of 5 cells. At each step, the agent has only one action to choose from, i.e. it moves to the cell on the immediate right. When it reaches the rightmost cell, it receives a reward of +1 and comes to a halt.

Let  $V^*(i)$  denote the value function of state  $i$ , the  $i^{th}$  cell starting from left.

Let  $V_k^*(i)$  denote the value function estimate at state  $i$  at the  $k^{th}$  step of the value iteration algorithm. Let  $V_0^*(i)$  denote the initialization of this estimate.

Use the discount factor  $\gamma = 0.5$ .

We will write the functions  $V_k^*$  as arrays below, i.e. as  $\begin{bmatrix} V_k^*(1) & V_k^*(2) & V_k^*(3) & V_k^*(4) & V_k^*(5) \end{bmatrix}$ .

Initialize by setting  $V_0^*(i) = 0$  for all  $i$ :

$$V_0^* = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \end{bmatrix}.$$

Then, using the value iteration update rule, we get

$$V_1^* = \begin{bmatrix} 0 & 0 & 0 & 0 & 1 \end{bmatrix},$$

$$V_2^* = \begin{bmatrix} 0 & 0 & 0 & 0.5 & 1 \end{bmatrix}$$

**Note (Aug 22):** Note that as soon as the agent takes the first action to reach cell 5, it halts and does not take any more action, so we set  $V_{k+1}^*(5) = V_k^*(5)$  for all  $k \geq 1$ .

## Value Function Update

1/1 point (graded)

Run the 3<sup>rd</sup> iteration of the value iteration algorithm to get  $V_3^*$  and answer the following questions:

Enter the value of  $V_3^*$  as an array  $\begin{bmatrix} V_3^*(0) & V_3^*(1) & V_3^*(2) & V_3^*(3) & V_3^*(4) \end{bmatrix}$ .

(For example, type  $[0,2,0,3,4]$  for the array  $\begin{bmatrix} 0 & 2 & 0 & 3 & 4 \end{bmatrix}$ .)

[0,0,1/4,1/2,1]

✔ Answer: [0, 0, 0.25, 0.5, 1]

### Solution:

Note that a non-zero reward is obtained only in state  $s_4$  when transitioning to  $s_5$ .

The 3<sup>rd</sup> step of the value iteration could be worked out as follows:

$$V_3^*(1) = 0 + \gamma * V_2^*(2)$$

$$V_3^*(1) = 0 + 0.5 * 0 = 0$$

$$V_3^*(2) = 0 + \gamma * V_2^*(3)$$

$$V_3^*(2) = 0 + 0.5 * 0 = 0$$

$$V_3^*(3) = 0 + \gamma * V_2^*(4)$$

$$V_3^*(3) = 0 + 0.5 * 0.5 = 0.25$$

$$V_3^*(4) = 0 + \gamma * V_2^*(5)$$

$$V_3^*(4) = 0 + 0.5 * 1 = 0.5$$

and  $V_3^*(5) = V_2^*(5) = 1$

The same computation for the rest of the states.

Submit

You have used 1 of 3 attempts

**i** Answers are displayed within the problem

## Number of steps till convergence

1/1 point (graded)

Enter below the number of steps it takes starting from  $V_0^*$  for the value function updates to converge to the optimal value function  $V^*$ :

✔ Answer: 5

**Solution:**

Note that after the  $5^{th}$  step, the reward from the rightmost cell in the grid gets propagated to the leftmost state after which the value function estimate  $V_k^*$  stops updating. Hence, for this example it takes 5 steps for the value function estimate to converge to the optimal value function.

You have used 1 of 2 attempts

---

**i** Answers are displayed within the problem

---

## Discussion

**Topic:** Unit 5 Reinforcement Learning (2 weeks) :Lecture 17. Reinforcement Learning 1 / 7. Value Iteration