

Lecture 8

Cezar Ionescu
06/07/2019

DEPARTMENT FOR
CONTINUING
EDUCATION



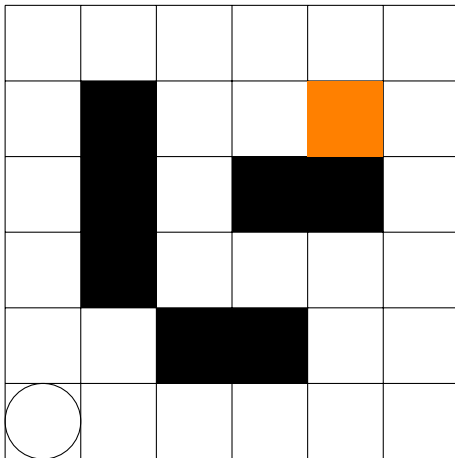
- Homework from 29/06/2019 due now!
- Please complete and hand in the declarations of authorship.

Questions?

Solution to homework from lecture 7

Dynamic programming

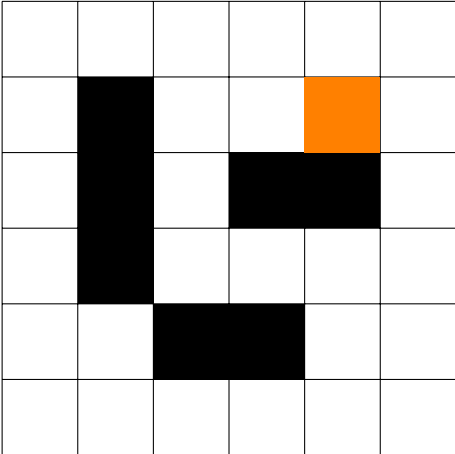
A maze



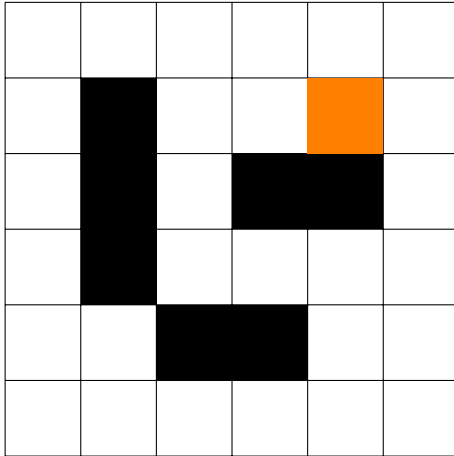
Value map

-4	-3	-2	-1	0	-1
-5		-1	0	1	0
-6		-2			-1
-7		-3	-4	-3	-2
-8	-9			-4	-3
-9	-8	-7	-6	-5	-4

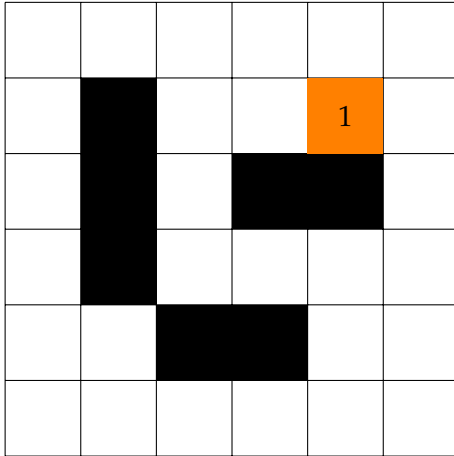
Dynamic programming



A 6x6 grid representing a search space. The grid contains several obstacles (black cells) and a goal cell (orange cell). The obstacles are located at the following coordinates (row, column): (1, 1), (1, 2), (1, 3), (1, 4), (1, 5), (1, 6), (2, 1), (2, 2), (2, 3), (2, 4), (2, 5), (2, 6), (3, 1), (3, 2), (3, 3), (3, 4), (3, 5), (3, 6), (4, 1), (4, 2), (4, 3), (4, 4), (4, 5), (4, 6), (5, 1), (5, 2), (5, 3), (5, 4), (5, 5), (5, 6), (6, 1), (6, 2), (6, 3), (6, 4), (6, 5), (6, 6). The goal cell is located at (2, 5).



Dynamic programming



Dynamic programming

				0	
			0	1	0

Dynamic programming

			-1	0	-1
		-1	0	1	0
					-1

Dynamic programming

		-2	-1	0	-1
		-1	0	1	0
		-2			-1
					-2

Dynamic programming

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

Dynamic programming

-4	-3	-2	-1	0	-1
		-1	0	1	0
		-2			-1
		-3	-4	-3	-2
				-4	-3
					-4

Dynamic programming

-4	-3	-2	-1	0	-1
-5		-1	0	1	0
		-2			-1
		-3	-4	-3	-2
				-4	-3
				-5	-4

Dynamic programming

-4	-3	-2	-1	0	-1
-5		-1	0	1	0
-6		-2			-1
		-3	-4	-3	-2
				-4	-3
			-6	-5	-4

Dynamic programming

-4	-3	-2	-1	0	-1
-5		-1	0	1	0
-6		-2			-1
-7		-3	-4	-3	-2
				-4	-3
		-7	-6	-5	-4

Dynamic programming

-4	-3	-2	-1	0	-1
-5		-1	0	1	0
-6		-2			-1
-7		-3	-4	-3	-2
-8				-4	-3
	-8	-7	-6	-5	-4

Dynamic programming

-4	-3	-2	-1	0	-1
-5		-1	0	1	0
-6		-2			-1
-7		-3	-4	-3	-2
-8	-9			-4	-3
-9	-8	-7	-6	-5	-4

Optimal policy

If we have the optimal value map, defining the optimal policy is easy:
Choose the action that results in the greatest sum of reward now and value in the next state.

Choose $\text{pol}(s)$ that maximises

$\text{Reward}(s, a) + \text{Value}(\text{sys}(s, a))$

where $\text{sys} : (\text{State}, \text{Action}) \leadsto \text{State}$ is the function that tells us what happens depending on the current state and chosen action.

Bellman equation

If $\text{pol}(s) = a^{\text{opt}}$ the optimal action in state s , then

$$\text{Val}(s) = \text{Rew}(s, a^{\text{opt}}) + \text{Val}(\text{sys}(s, a^{\text{opt}}))$$

Therefore

$$\text{Val}(s) = \max_a (\text{Rew}(s, a) + \text{Val}(\text{sys}(s, a)))$$

This is called *Bellman's equation*.

Direct adaptive optimal control

Model of the problem

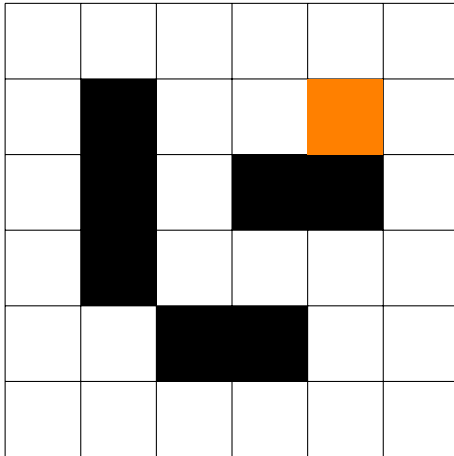
If we do not know the exact map of the maze, or the location of the goal, we cannot apply dynamic programming. We need an *adaptive method*. Alternatives:

- explore the maze and make a map of it, find the goal, and then apply dynamic programming (*indirect method*);
- learn the optimal value map directly! (*direct method*).

Reinforcement learning is *direct adaptive optimal control*.

Value iteration

Start with a randomly generated value map



Start with a randomly generated value map

0	0	0	0	0	0
0		0	0	0	0
0		0			0
0		0	0	0	0
0	0			0	0
0	0	0	0	0	0

Alternatives

0	0	0	0	0	0
0		0	0	0	0
0		0			0
0		0	0	0	0
0	0			0	0
0	0	0	0	0	0

Alternatives

A 6x6 grid with the following cell values and highlights:

0	0	0	0	0	0
0		0	0	0	0
0		0			0
0		0	0	0	0
0	0			0	0
0	0	0	0	0	0

Diagram details:

- The cell at row 2, column 3 (value 0) is circled in blue.
- Red arrows point from this cell to the cells at (2,1), (2,4), and (3,3).
- The cell at (2,5) is highlighted in orange.
- Cells at (2,2), (3,1), (3,2), (3,4), (3,5), (4,1), (4,2), (4,3), (4,4), (4,5), (5,3), (5,4), and (5,5) are blacked out.

Choose an action

A 6x6 grid world environment. The grid contains obstacles (black cells) and a goal state (orange cell). The current state is highlighted with a blue circle and a red arrow pointing down.

0	0	0	0	0	0
0		0	0	0	0
0		0			0
0		0	0	0	0
0	0			0	0
0	0	0	0	0	0

Obstacles (black cells) are located at: (1,1), (1,2), (1,3), (1,4), (2,1), (2,2), (2,3), (2,4), (3,2), (3,3), (3,4), (3,5), (4,2), (4,3), (4,4), (4,5).

Goal state (orange cell) is located at: (1,4).

Current state (blue circle and red arrow) is at: (0,2).

Make a move

0	0	0	0	0	0
0		0	0	0	0
0		0			0
0		0	0	0	0
0	0			0	0
0	0	0	0	0	0

Changing the current value function

If we had the optimal value function, then

$$\text{Val}((3,6)) = \text{Rew}((3,6), \downarrow) + \text{Val}((3,5))$$

But we have

$$\text{Val}((3, 6)) = 0, \quad \text{Rew}((3,6), \downarrow) = -1, \quad \text{Val}((3,5)) = 0$$

Update rule

Idea from *supervised learning*:

```
Val((3,6)) <- Val((3,6)) +  
    δ * (Rew((3,6), ↓) + Val((3,5)) - Val((3,6)))
```

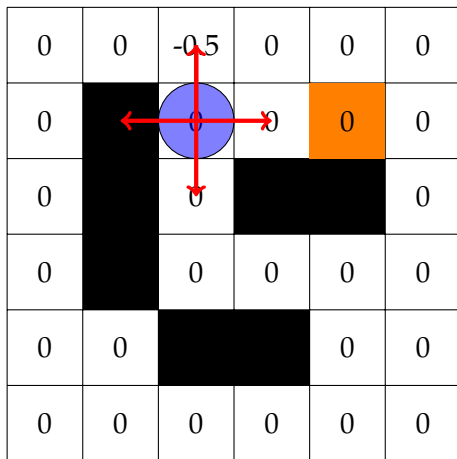
Therefore (say $\delta = 0.5$)

```
Val(3,6) <- 0 + 0.5 * (-1 + 0 - 0) = -0.5
```

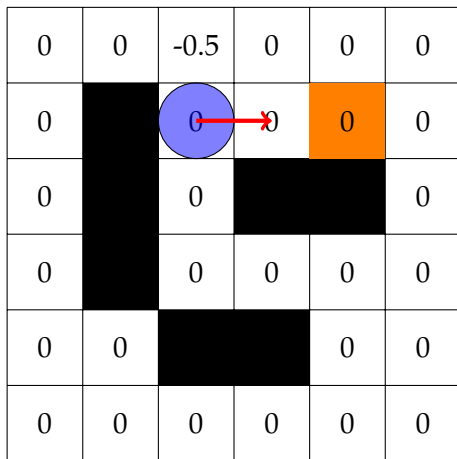

The new situation

0	0	-0.5	0	0	0
0		0	0	0	0
0		0			0
0		0	0	0	0
0	0			0	0
0	0	0	0	0	0

The new situation



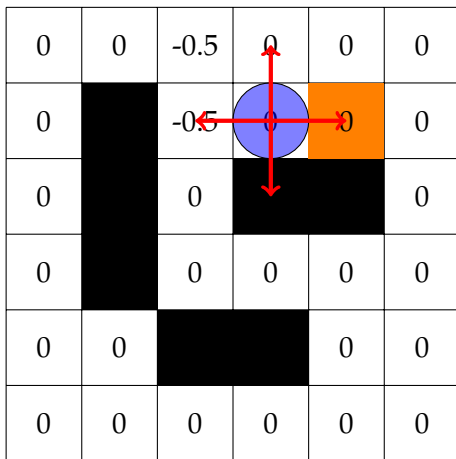
Pick a new direction



Move and update

0	0	-0.5	0	0	0
0		-0.5	0	0	0
0		0			0
0		0	0	0	0
0	0			0	0
0	0	0	0	0	0

And so on...



A good choice

0	0	-0.5	0	0	0
0		-0.5	0	0	0
0		0			0
0		0	0	0	0
0	0			0	0
0	0	0	0	0	0

A good choice

0	0	-0.5	0	0	0
0		-0.5	0.5	0	0
0		0			0
0		0	0	0	0
0	0			0	0
0	0	0	0	0	0

The game goes on

0	0	-0.5	0	0	0
0		-0.5	0.5	0	0
0		0			0
0		0	0	0	0
0	0			0	0
0	0	0	0	0	0

Discount factors

In this setting, we maximise the sum of rewards. Problem: if the number of steps is large, the values become difficult to estimate (huge absolute value).

To alleviate this problem, we introduce a *discount factor*: $0 < \beta < 1$

Instead of maximising

$$\text{Rew}(s_0, a_0) + \text{Rew}(s_1, a_1) + \text{Rew}(s_2, a_2) + \dots$$

we maximise

$$\text{Rew}(s_0, a_0) + \beta * \text{Rew}(s_1, a_1) + \beta^2 * \text{Rew}(s_2, a_2) + \dots$$

where $s_{t+1} = \text{sys}(s_t, \text{pol}(s_t))$

Bellman's equation

We choose $pol(s)$ that maximises

$$Rew(s, a) + \beta * Val(sys(s, a))$$

Therefore

$$Val(s) = \max_a (Rew(s, a) + \beta * Val(sys(s, a)))$$

Value iteration and function approximation

$$Val(s) = \max_a (Rew(s, a) + \beta * Val(sys(s, a)))$$

Bellman's equation allows us to iteratively approximate the optimal value function.

$$Val(s) \leftarrow Val(s) + \delta * (\max_a (Rew(s, a) + \beta * Val(sys(s, a)))) - Val(s))$$

The optimal value function Val is unique (but there might be many policy functions that realise it).

Value iteration

$$Val(s) \leftarrow Val(s) + \delta * (\max_a (Rew(s, a) + \beta * Val(s_{next}(s, a)))) - Val(s)$$

Problem: we require knowledge of the reward function.

We could use supervised learning to approximate the reward function.

But now we have two function approximations and an optimisation at every step!

Q -learning

Watkins' idea

The reward function only enters the picture when combined with **Val**, e.g.:

$$\text{Val}(s) = \max_a (\text{Rew}(s, a) + \beta * \text{Val}(\text{sys}(s, a)))$$

Chris Watkins (1989): maybe it's simpler to learn the combination of **Rew** and **Val**!

$$Q(s, a) = \text{Rew}(s, a) + \beta * \text{Val}(\text{sys}(s, a))$$

Q-learning and Val

$$Q(s, a) = \text{Rew}(s, a) + \beta * \text{Val}(\text{sys}(s, a))$$

We have $\text{Val}(s) = \max_a Q(s, a)$ therefore, knowing Q is sufficient for determining Val .

Q-learning and the optimal policy

$$Q(s, a) = \text{Rew}(s, a) + \beta * \text{Val}(\text{sys}(s, a))$$

The optimal policy is the one that maximises

$$\text{Rew}(s, a) + \beta * \text{Val}(\text{sys}(s, a))$$

Therefore, the optimal action in a state s is the one that maximises the Q function:

$$a^{\text{opt}} = \arg \max_a Q(s, a)$$

Recursive equation for Q-learning

$$Q(s, a) = \text{Rew}(s, a) + \beta * \text{Val}(\text{sys}(s, a))$$

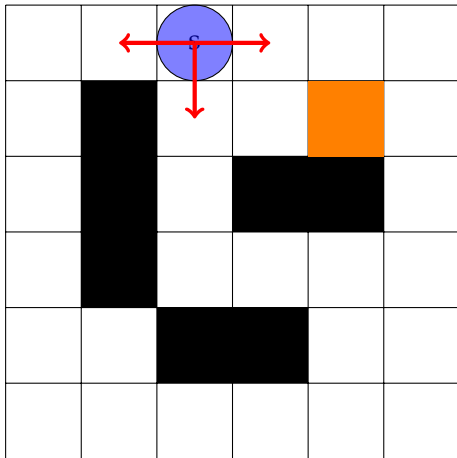
$$\text{Val}(s) = \max_a Q(s, a)$$

Therefore

$$Q(s, a) = \max_a (\text{Rew}(s, a) + \beta * \max_x Q(\text{sys}(s, a), x))$$

We have a recursive equation for the Q function.

Q-learning in action

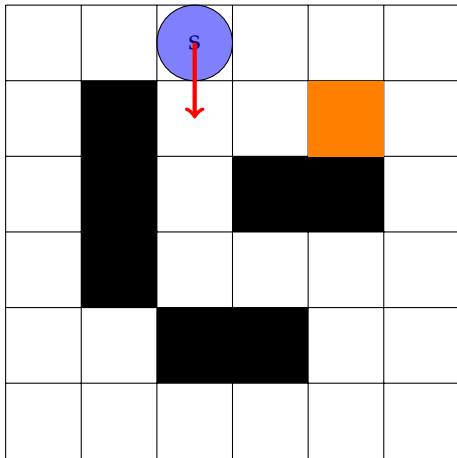


$$Q(s, \leftarrow) = -10$$

$$Q(s, \downarrow) = -8$$

$$Q(s, \rightarrow) = -9$$

Q-learning in action

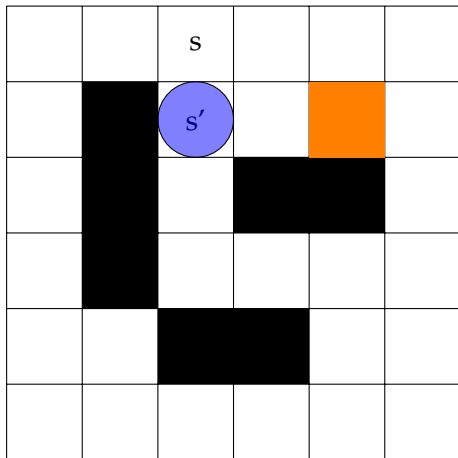


$$Q(s, \leftarrow) = -10$$

$$Q(s, \downarrow) = -8$$

$$Q(s, \rightarrow) = -9$$

Q-learning in action



$$Q(s, \downarrow) = -8$$

$$Rew(s, \downarrow) = -1$$

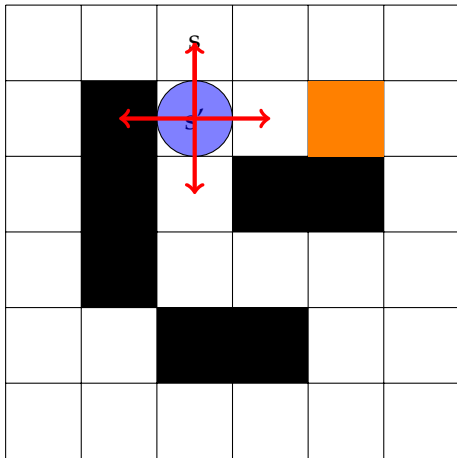
$$Q(s', \leftarrow) = -\infty$$

$$Q(s', \downarrow) = -8$$

$$Q(s', \rightarrow) = -5$$

$$Q(s', \uparrow) = -10$$

Q-learning in action



$$Q(s, \downarrow) = -8$$

$$Rew(s, \downarrow) = -1$$

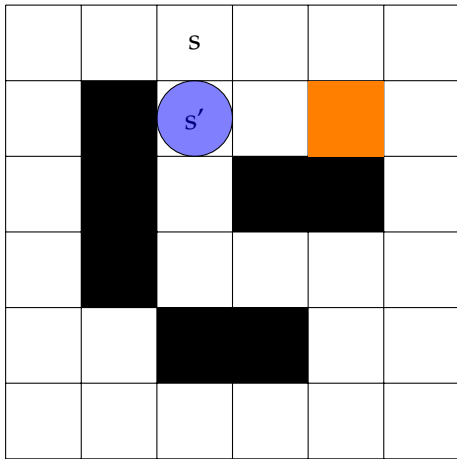
$$Q(s', \leftarrow) = -\infty$$

$$Q(s', \downarrow) = -8$$

$$Q(s', \rightarrow) = -5$$

$$Q(s', \uparrow) = -10$$

Q-learning in action



$$Q(s, \downarrow) = -8$$

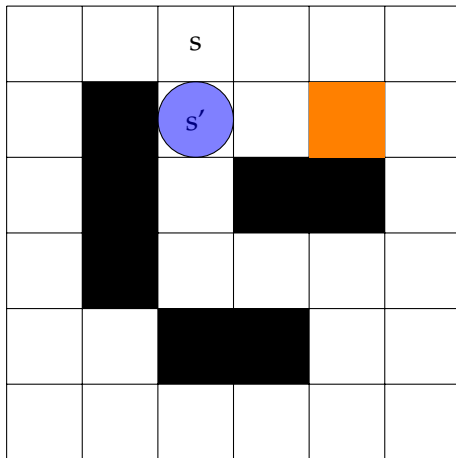
$$Rew(s, \downarrow) = -7$$

$$Val(s') =$$

$$\max(-\infty, -8, -5, -10)$$

$$Rew(s, \downarrow) + \beta Val(s') =$$

Q-learning in action



$$Q(s, \downarrow) = -8$$

$$Rew(s, \downarrow) + \beta Val(s') = -6$$

$$Q(s, \downarrow) \leftarrow -8 + \delta(-6 - (-8)) = -6$$

We have assumed a deterministic system, but in most cases this is unrealistic. E.g., the maze environment could take us (with a small probability) in a different direction than we chose. In that case, we need to maximise the *expected value* of the sum of rewards. *Q*-learning can be applied to these cases almost without any change!

Conclusions

Reinforcement learning applications: games, RoboCup Soccer, self-driving cars, etc. See also discussion in Week 1 (robot-catching arm, self-driving car, barriers on emissions...) Reinforcement learning seems to make the most out of very little: is it the road to Artificial General Intelligence? Chris Watkins

(<http://www.cs.rhul.ac.uk/~chrisw/>):

I have long felt the standard model of RL is deceptively attractive, but limited.