# Lecture 8

Cezar Ionescu
06/07/2019

DEPARTMENT FOR
CONTINUING
EDUCATION

UNIVERSITY OF
OXFORD

# Administrative

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

# Questions?

# Solution to homework from lecture 7

```
y₁(v₁₁, v₂₁) = f(x₁*v₁₁ + x₂*v₂₁)     = -2
y₂(v₁₁, v₂₁) = f(x₁*v₁₂ + x₂*v₂₂)     =  2
o(w₁, w₂, y₁, y₂) = f(y₁*w₁ + y₂*w₂) =  2
E (w₁, w₂, v₁₁, ...) = (t - o)²

∂ E / ∂ v₁₁ = E'(o) * ∂ o / ∂ v₁₁
            = 2*(t - o)*(-1) * ∂ o / ∂ v₁₁
            = 2*(-2)*(-1)*f'(z)*∂ g / ∂ v₁₁
            = 4*1*∂ (y₁*w₁ + y₂*w₂) / ∂ v₁₁
            = 4 * ∂ (y₁*w₁ + y₂*w₂) / ∂ y₁ * ∂ y₁ / ∂ v₁₁
            = 4 * w₁ * f'(z) * ∂ (x₁*v₁₁ + x₂*v₂₁) / ∂ v₁₁
            = 4 * x₁ = -4
```

# Dynamic programming

# A maze

| -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

# Dynamic programming

# Dynamic programming

# 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

# 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 | -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 `pol(s)` that maximises

`Reward(s, a) + Value(sys(s, a))`

where `sys : (State, Action) ~> State` is the function that tells us what happens depending on the current state and chosen action.

# Review

Ingredients:

```
sys : (State, Action) -> State
rew : (State, Action) -> ℝ
pol : State -> Action
s₀ ∈ State -- given
```

Problem: Find `pol` that maximises $\Sigma_0^n\ \text{rew}(s_i,\ \text{pol}(s_i))$, where

```
s₍ᵢ₊₁₎ = sys(sᵢ, aᵢ)
```

# Value function

```
Val : ℕ -> State -> ℝ
```

$$Val_i(s) = \max_{pol} \Sigma_i^n \, rew(s_i, pol(s_i))$$

Note:

$Val_0(s_0)$ is the maximal value for the entire problem.

# Bellman equation

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

`Val`$_i$`(s) = Rew(s, a`$^{opt}$`) + Val`$_{i+1}$`(sys(s, a`$^{opt}$`))`

Therefore

`Val`$_i$`(s) = max`$_a$` (Rew(s, a) + Val`$_{i+1}$`(sys(s, a)))`

This is called *Bellman's equation*.
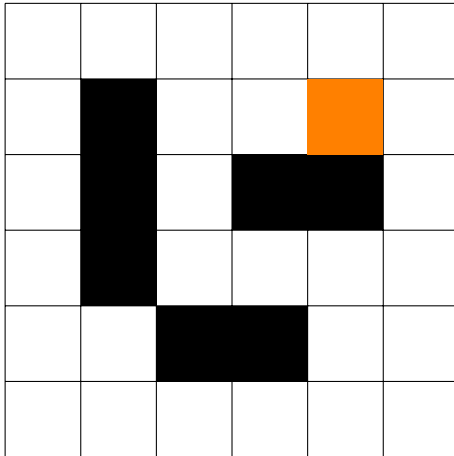
**Direct adaptive optimal control**

# 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

# Alternatives

If we had the optimal value function, then

`Val((3,6)) = Rew((3,6), ↓) + Val((3,5))`

But we have

`Val((3, 6)) = 0,  Rew((3,6), ↓) = -1, 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)))
```
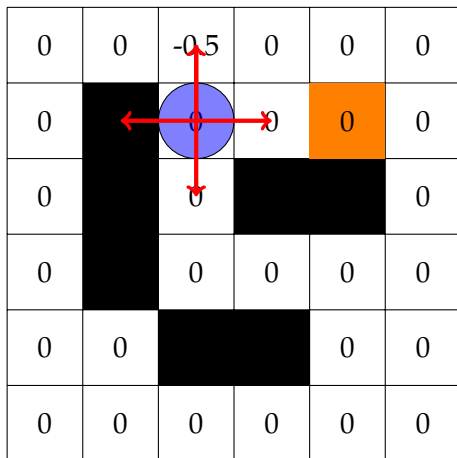
Therefore (say δ = 0.5)
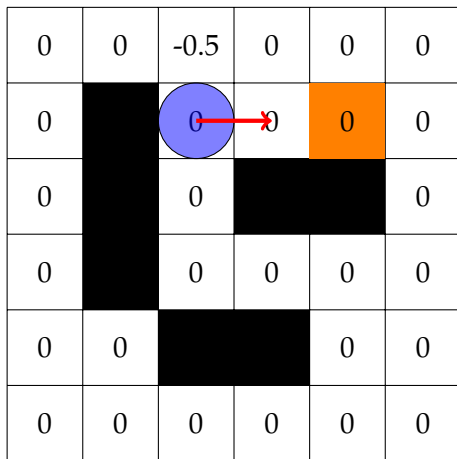
```
Val(3,6) <- 0 + 0.5 * (-1 + 0 - 0) = -0.5
```

# The new situation

## 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). Moreover, in many situations we need to operate with *infinite horizons*.

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

Instead of maximising

$Rew(s_0, pol(s_0)) + Rew(s_1, pol(s_1)) + Rew(s_2, pol(s_2)) + \ldots$

we maximise

$Rew(s_0, pol(s_0)) + \beta*Rew(s_1, pol(s_1)) + \beta^2*Rew(s_2, pol(s_2)) + \ldots$

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

# Value function

If we have an infinite horizon, the value function is *stationary*:

$V_i(s) = \max_{pol} \Sigma_i^\infty \ rew(s_t, \ pol(s_t))$ where

$s_i = s$

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

This is independent of i!

# Bellman's equation

We choose pol(s) that maximises

Rew(s, a) + β * Val(sys(s, a))

Therefore

Val(s) = max$_a$ (Rew(s, a) + β * Val(sys(s, a)))

# Value iteration and function approximation

```
Val(s) = maxₐ (Rew(s, a) + β * Val(sys(s, a)))
```

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

```
Val(s) <- Val(s) +
          δ * (maxₐ (Rew(s, a) + β * 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) <- Val(s) +
        δ * (maxₐ (Rew(s, a) + β * Val(sys(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**

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

`Val(s) = max`$_a$` (Rew(s, a) + β * Val(sys(s, a)))`

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

`Q(s, a) = Rew(s, a) + β * Val(sys(s, a))`

```
Q(s, a) = Rew(s, a) + β * Val(sys(s, a))
```

We have `Val(s) = max` $_a$ `Q(s, a)` therefore, knowing `Q` is sufficient for determining `Val`.

`Q(s, a) = Rew(s, a) + β * Val(sys(s, a))`

The optimal policy is the one that maximises

`Rew(s, a) + β * Val(sys(s, a))`

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

`aᵒᵖᵗ = arg maxₐ Q(s, a)`

# Recursive equation for Q-learning

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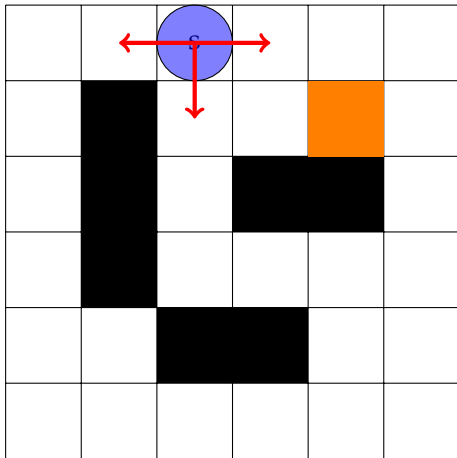
$Val(s) = \max_a Q(s, a)$

Therefore

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

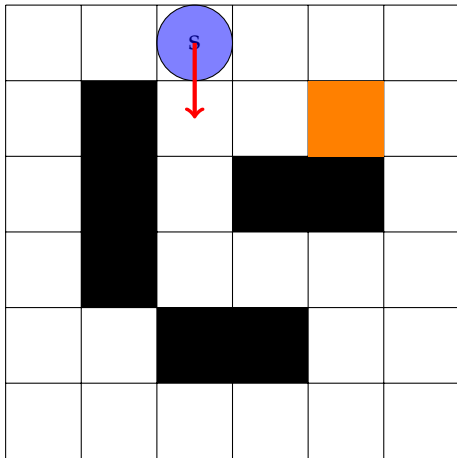We have a recursive equation for the `Q` function.

# Q-learning in action
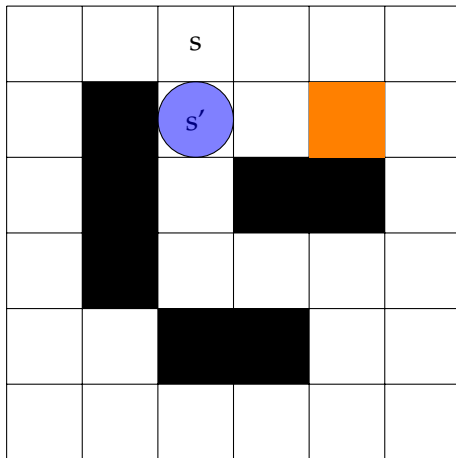


$Q(s, \leftarrow) = -10$
$Q(s, \downarrow) = -8$
$Q(s, \rightarrow) = -9$

$$Q(s, \downarrow) = -8$$
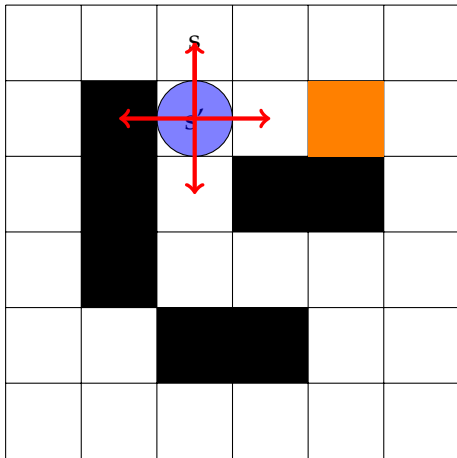
# Q-learning in action



$Q(s, \downarrow) = -8$
$Rew(s, \downarrow) = -1$

$$Q(s, \downarrow) = -8$$
$$Rew(s, \downarrow) = -1$$
$$Q(s', \leftarrow) = -\infty$$
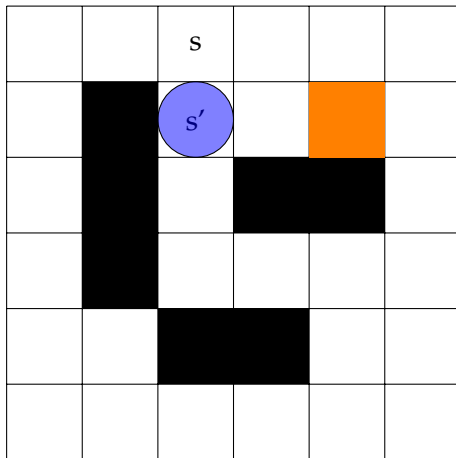$$Q(s', \downarrow) = -8$$
$$Q(s', \rightarrow) = -5$$
$$Q(s', \uparrow) = -10$$

$Q(s, \downarrow) = -8$
$Rew(s, \downarrow) = -1$
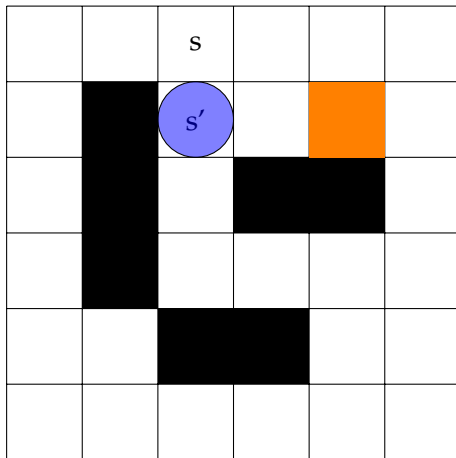$Val(s') =$
$max(-\infty, -8, -5, -10)$
$Rew(s, \downarrow) + \beta Val(s') =$
$-1 + 1 * (-5) = -6$

# Q-learning in action



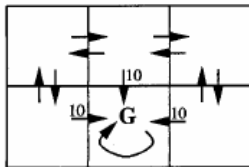$Q(s, \downarrow) = -8$
$Rew(s, \downarrow) + \beta Val(s') = -6$
$Q(s, \downarrow) \leftarrow$
$-8 + \delta(-6 - (-8)) = -6$

**13.2.** Consider the deterministic grid world shown below with the absorbing goal-state G. Here the immediate rewards are 10 for the labeled transitions and 0 for all unlabeled transitions.

(a) Give the $V^*$ value for every state in this grid world. Give the $Q(s, a)$ value for every transition. Finally, show an optimal policy. Use $\gamma = 0.8$.

(b) Suggest a change to the reward function $r(s, a)$ that alters the $Q(s, a)$ values, but does not alter the optimal policy. Suggest a change to $r(s, a)$ that alters $Q(s, a)$ but does not alter $V^*(s, a)$.

(c) Now consider applying the $Q$ learning algorithm to this grid world, assuming the table of $\hat{Q}$ values is initialized to zero. Assume the agent begins in the bottom left grid square and then travels clockwise around the perimeter of the grid until it reaches the absorbing goal state, completing the first training episode. Describe which $\hat{Q}$ values are modified as a result of this episode, and give their revised values. Answer the question again assuming the agent now performs a second identical episode. Answer it again for a third episode.

# Stochastic systems

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.*