**Aalto University**
**School of Electrical**
**Engineering**

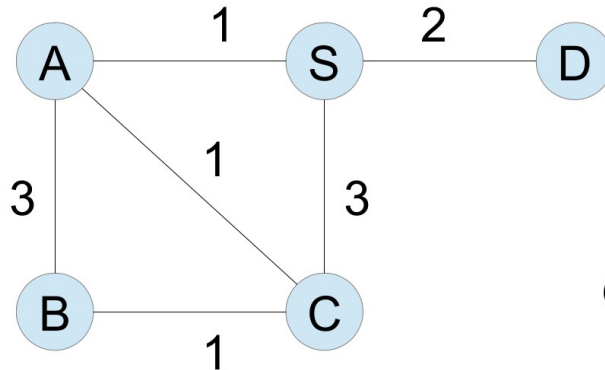# ELEC-E8125 Reinforcement Learning Solving discrete MDPs

Joni Pajarinen

12.9.2023

# Previous lecture: find shortest path exercise

- Use backward value iteration for

Initial start state $\longrightarrow$ $s_I = S$

$S_G = \{B\}$

Final cost

$$l_F(s) = \begin{array}{l} 0, s \in S_G \\ \infty, s \notin S_G \end{array}$$

Goal set



Reminder:

$$G^*(s) = min_a \{l(s,a) + G^*(f(s,a))\}$$

Hints: Moving from A to B costs $l(A, A \rightarrow B) = 3$. Arriving at final state B costs $l_F(B) = 0$. Arriving at final state C costs $l_F(C) = \infty$. When moving from C to B we arrive at B: $B = f(C, C \rightarrow B)$
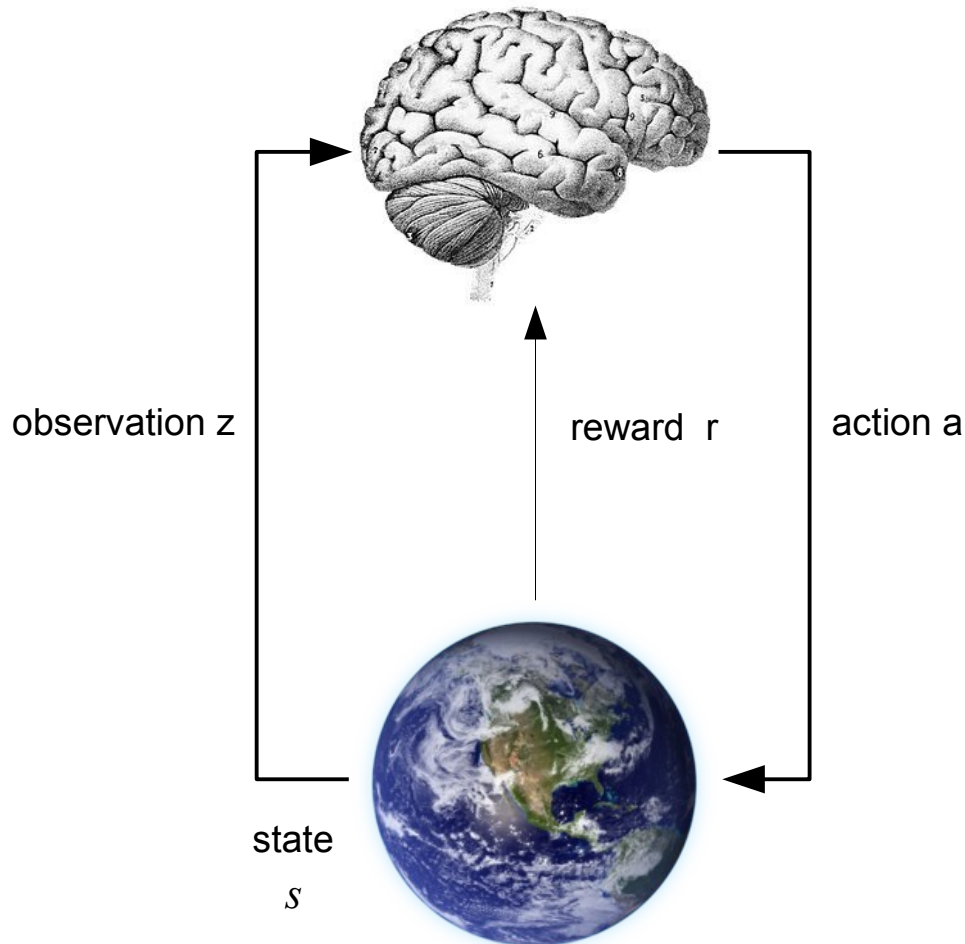
# Today

- Markov decision processes (MDPs)

# Learning goals

- Understand MDPs and related concepts
- Understand value functions

- Be able to implement value iteration for determining an optimal policy (in exercise #2, you will get to do this for real)

# Markov decision process (MDP)



observation z

reward  r

action a

state

$s$

**MDP**
Environment observable
$z = s$

Defined by dynamics
$P(s_{t+1}|s_t, a_t)$

And reward function
$r_t = r(s_t, a_t)$

Solution, for example
$a^*_{1,...,T} = arg\ max_{a_1,...,a_T} \sum_{t=1}^{T} r_t$

Represented as policy
$a = \pi(s)$

Let's discuss MDPs in more detail

# Markov property


Andrey Markov

- "Future is independent of past given the present"

- State sequence $S$ is Markov iff ← "if and only if"
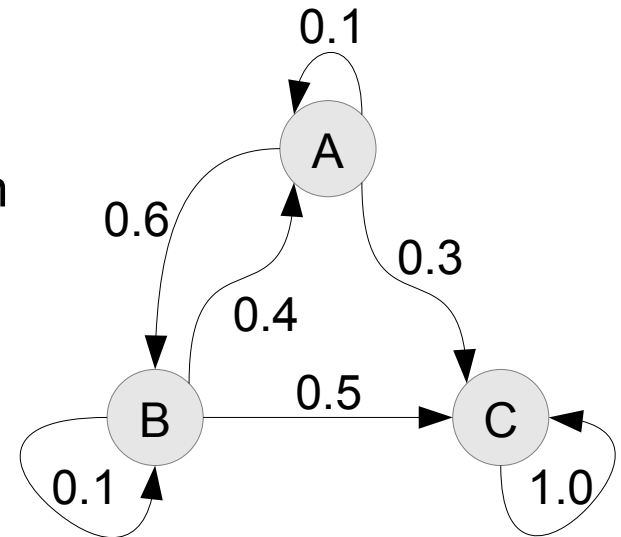
$$P(S_{t+1}|S_t) = P(S_{t+1}|S_1, \ldots, S_t)$$

- State captures all history
- Once state is known, history may be thrown away

**Aalto University**
**School of Electrical**
**Engineering**

# Markov process

No "decision" here!

- Markov process is a memoryless random process that generates a state sequence *S* with the Markov property

- Defined as *(S,T)*
  - *S:* set of states
  - *T*: S *x* S → *[0,1]* state transition function
    - $T_t(s,s') = P(s_{t+1} = s' | s_t = s)$
    - *P* can be represented as a transition probability matrix
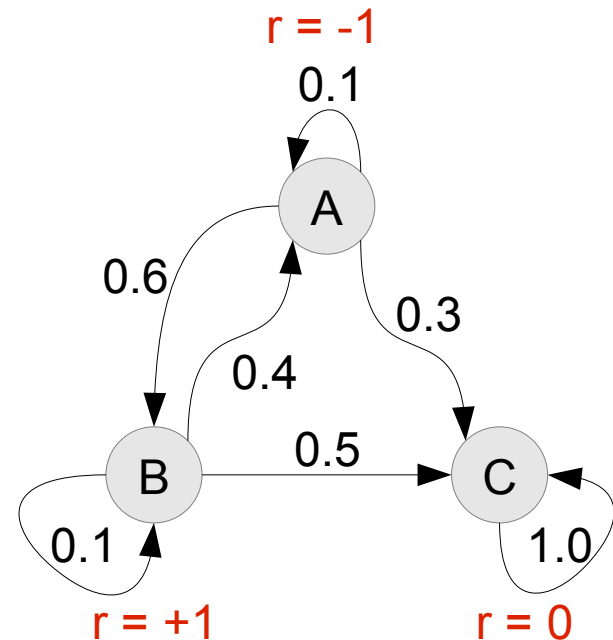- State sequences called *episodes*

0.1

A

0.6

0.3

0.4

B 0.5 C

0.1 1.0

How to calculate probability of a particular episode?
Starting from A, what is the probability of A,B,C?

# Markov reward process

- Markov reward process =
  Markov process with rewards
- Defined by (S, *T, r,* $\gamma$ )
  - *S, T* : as above
  - *r* : S $\rightarrow \mathcal{R}$   reward function
  - $\gamma \in [0, 1]$   discount factor
- Accumulated rewards in finite
  (*H* steps) or infinite horizon

$$\sum_{t=0}^{H} \gamma^t r_t \qquad \sum_{t=0}^{\infty} \gamma^t r_t$$

- *Return G*: accumulated rewards from time t

r = -1

0.1

A

0.6

0.3

0.4

B

0.5

C

0.1

1.0

r = +1

r = 0

A" Aalto University
School of Electrical
Engineering

$$G_t = \sum_{k=0}^{H} \gamma^k r_{t+k}$$

Why discount?

Return of (A,B,C), $\gamma$ = 0.9?

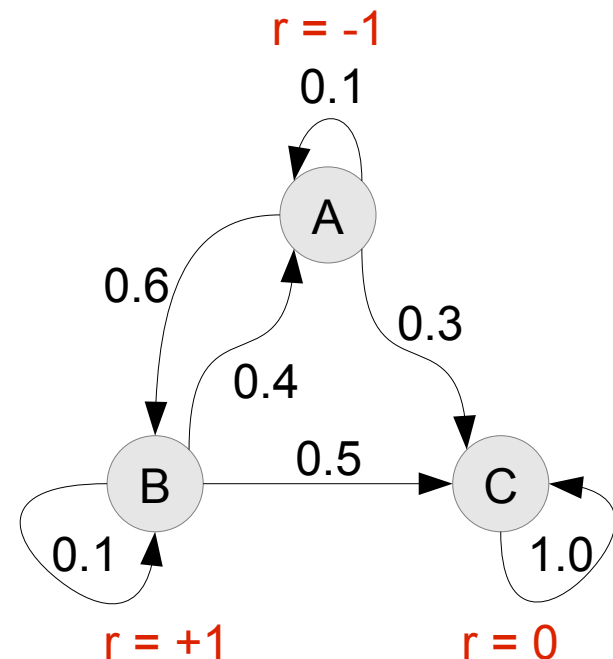# State value function for Markov reward processes

- State value function *V(s)* is expected cumulative reward starting from state *s*

$$V(s) = E[G_t | s_t = s]$$

- Value function can be defined by the Bellman equation

$$V(s) = E[G_t | s_t = s]$$
$$V(s) = E[r_t + \gamma V(s_{t+1}) | s_t = s]$$

r = -1

0.1

A

0.6

0.3

0.4

0.5

B

C

0.1

1.0

r = +1

r = 0

What is the value function for $\gamma$ = 0?

What is the value function for $\gamma$ = 0.5 after a single Bellman update when starting with zero values?

# Markov decision process (MDP)

Grid world



- Markov decision process
  defined by $(S, A, T, R, \gamma)$
  - $S, \gamma$ : as above
  - $A$: set of actions (inputs)
  - $T: S \times A \times S \rightarrow [0,1]$
    $$T_t(s, a, s') = P(s_{t+1} = s' | s_t = s, a_t = a)$$
  - $R: S \times A \rightarrow \Re$ reward function
    $$r_t(s, a) = r(s_t = s, a_t = a)$$

- Goal: Find policy $\pi(s)$ that maximizes
  expected cumulative reward

Agent tries to move forward:
P(success) = 0.8
P(left) = 0.1
P(right) = 0.1

Grid world example!

# Policy

- Deterministic policy $\pi$(S)*: S $\rightarrow$ A* is a mapping from states to actions

- Stochastic policy $\pi$(a|s): S,A $\rightarrow$ [0,1] is a distribution over actions given states

- Optimal policy $\pi_*$(s) is a policy that is better or equal than any other policy (in terms of cumulative rewards)
  - There always exists a deterministic optimal policy for an MDP

Grid world



Agent tries to move forward:
P(success) = 0.8
P(left) = 0.1
P(right) = 0.1

What is the optimal policy in the grid world?

# MDP value function

- *State-value function* of an MDP is the expected return starting from state $s$ and following policy $\pi$

$$V_\pi(s) = E_\pi[G_t | s_t = s]$$

- Can be decomposed into immediate and future components using Bellman expectation equation

$$V_\pi(s) = E_\pi[r_t + \gamma V_\pi(s_{t+1}) | s_t = s]$$

$$V_\pi(s) = r(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_\pi(s')$$

What is value function here?

# Action-value function

- *Action-value function Q* is expected return starting from state *s*, taking action *a*, and then following policy $\pi$

$$Q_\pi(s,a) = E_\pi[G_t | s_t = s, a_t = a]$$

- Using Bellman expectation equation

$$Q_\pi(s,a) = E_\pi[r_t + \gamma Q_\pi(s_{t+1}, a_{t+1} | s_t = s, a_t = a)]$$

$$Q_\pi(s,a) = r(s,a) + \gamma \sum_{s'} T(s,a,s') Q_\pi(s', \pi(s'))$$

# Optimal value function

- Optimal state-value function is maximum value function over all policies

$$V^*(s) = max_\pi V_\pi(s)$$

- Optimal action-value function is maximum action-value function over all policies

$$Q^*(s,a) = max_\pi Q_\pi(s,a)$$

- All optimal policies achieve optimal state- and action-value functions

What is the optimal action if we know $Q$*?
What about $V$*?

# Optimal policy vs optimal value function

- Optimal policy for optimal action-value function

$$\pi^*(s) = arg\,max_a\, Q^*(s, a)$$

- Optimal action for optimal state-value function

$$\pi^*(s) = arg\,max_a\, E_{s'}[r(s, a) + \gamma\, V^*(s')]$$

$$\pi^*(s) = arg\,max_a\left(r(s, a) + \gamma \sum_{s'} T(s, a, s')\, V^*(s')\right)$$

# Value iteration

Do you notice that this is an expectation?

- Starting from $V_0^*(s) = 0 \quad \forall s$
  iterate

$$V_{i+1}^*(s) = max_a \left( r(s,a) + \gamma \sum_{s'} T(s,a,s') V_i^*(s') \right)$$

until convergence

- Value iteration converges to *V\*(s)*

Compare to

$$G^*(s) = min_a \left[ l(s,a) + G^*(f(s,a)) \right]$$

from last week!
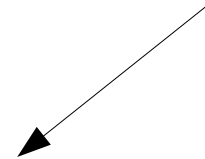
**A''** **Aalto University**
**School of Electrical**
**Engineering**

# Iterative policy evaluation

- Problem: Evaluate value of policy $\pi$

- Solution: Iterate Bellman expectation back-ups

- $V_1 \rightarrow V_2 \rightarrow \ldots \rightarrow V_\pi$

- Using synchronous back-ups:
  - For all states $s$
  - Update $V_{k+1}(s)$ from $V_k(s')$
  - Repeat

From slide 12

$$V_{k+1}(s) = r(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_k(s')$$

$$V_{k+1}(s) = \sum_a \pi(a|s) \left( r(s,a) + \gamma \sum_{s'} T(s,a,s') V_k(s') \right)$$

Note: Starting point can be random policy

# V

# Greedy policy

## k = 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 | 0.0 | 0.0 |

|  |   | 1 | 2 | 3 |
|  | 4 | 5 | 6 | 7 |
|  | 8 | 9 | 10 | 11 |
|  | 12 | 13 | 14 |  |

## k = 1

| 0.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | 0.0 |

r = -1 for all actions

## k = 2

| 0.0 | -1.7 | -2.0 | -2.0 |
| -1.7 | -2.0 | -2.0 | -2.0 |
| -2.0 | -2.0 | -2.0 | -1.7 |
| -2.0 | -2.0 | -1.7 | 0.0 |

# Policy improvement and policy iteration

- Given a policy $\pi$, it can be improved by
  - Evaluating $V_\pi$
  - Forming a new policy by acting greedily with respect to $V_\pi$

- This always improves the policy

- Iterating multiple times called *policy* iteration
  - Converges to optimal policy

# Computational limits – Value iteration

- Complexity $O(|A||S|^2)$ per iteration
- Effective up to medium size problems (millions of states)

- Complexity when applied to action-value function $O(|A|^2|S|^2)$ per iteration

# Summary

- Markov decision processes represent environments with uncertain dynamics

- Deterministic optimal policies can be found using state-value or action-value functions

- Dynamic programming is used in value iteration and policy iteration algorithms

# Next week: From MDPs to RL

- Readings
  - Sutton & Barto Ch. 5-5.4, 5.6, 6-6.5