## **Reinforcement Learning**

Basic Concepts in RL

Akka Zemmari

### Introduction

Agent

Environment

State

Action

Reward

Episode

## Introduction

### RL in a nutshell



Figure 1: RL: An agent interacting with an environment

### Key Concepts:

- Agent
- Environment
- Action
- Reward
- Episode

### Toy Example



Figure 2: Toy Example: Grid World

- The agent navigates from a starting position to a goal while avoiding obstacles.
- Hands-on: see files Grid.py and starter\_grid.py.

## Agent



**Agent:** The learner or decision maker that interacts with the environment.

### **Components:**

- Policy
- Value Function
- Model

- In the toy example, the agent is the robot navigating the grid world.

## Environment

# **Environment:** The external system with which the agent interacts.

### **Components:**

- State
- Action
- Reward

## State

State: A representation of the environment.



**Figure 3:** Toy Example: the state is the position of the robot in the Grid World.

## Action

Action

Action: The set of possible moves the agent can make.



**Figure 4:** Toy Example: the agent can move up  $(a_1)$ , right  $(a_2)$ , down  $(a_3)$ , left  $(a_4)$ , or stay in its place  $(a_5)$ .

### Reward



**Reward:** A scalar feedback signal from the environment.



**Figure 5:** Toy Example: the agent receives a reward of +1 when reaching the goal, -1 when hitting an obstacle or get out of the boundary, and 0 otherwise.

Reward

**Reward:** Tabular representation (suitable for programming).

	a <sub>1</sub>	a <sub>2</sub>	a <sub>3</sub>	a <sub>4</sub>	a <sub>5</sub>
<i>s</i> <sub>1</sub>	-1	0	0	-1	0
<i>s</i> <sub>2</sub>	-1	0	0	0	0
<i>S</i> 3	-1	-1	-1	0	0
<i>S</i> <sub>4</sub>	0	0	-1	-1	0
<i>S</i> 5	0	-1	0	0	0
<i>s</i> <sub>6</sub>	0	-1	+1	0	-1
<i>S</i> <sub>7</sub>	0	0	-1	-1	-1
<i>S</i> <sub>8</sub>	0	+1	-1	-1	0
<b>S</b> 9	-1	-1	-1	0	+1

**Return:** The sum of rewards over time steps.

A trajectory is a sequence of states, actions, and rewards:

$$s_1 \xrightarrow{r_1}_{a_1} s_2 \xrightarrow{r_2}_{a_2} s_3 \xrightarrow{r_3}_{a_3} \dots \xrightarrow{r_{T-1}}_{a_{T-1}} s_T.$$
(1)

The return is the sum of rewards:

$$return = r_1 + r_2 + \ldots + r_{T-1}.$$
 (2)



### **Return:** The sum of rewards over time steps.



### Figure 6: Toy Example: trajectories and returns.

**Discounted Return:** The sum of rewards over time steps with a discount factor  $\gamma$ .

Definition:

discounted return =  $r_1 + \gamma r_2 + \gamma^2 r_3 + \ldots + \gamma^{T-1} r_{T-1}$ . (3)

- $\gamma \in [0,1]$  is the discount factor:
  - γ = 0: the agent is myopic (only cares about the immediate reward).
  - $\gamma = 1$ : the agent is far-sighted (cares about all future rewards).
- Toy Example: see whiteboard.

# Episode

### Episode

**Episode:** A sequence of time steps where the agent interacts with the environment.



Figure 7: Toy Example: an episode in the grid world.

- An episode is usually assumed to terminate in a finite number of time steps. Tasks with episodes are called episodic tasks.