

# Supervised Learning Attempt

---

- Suppose either returned 0 or 1 rewards....
- Prob of selection action 1 or action 2 (d=1 or d=2)
  - “Correct action” chosen:  
$$\pi_{t+1}(d_t) = \pi_t(d_t) + \alpha[1 - \pi_t(d_t)]$$
  - Other action adjusted so that both sum to 1
  
- What if both actions often succeed, e.g.,  $Q(d=1) = 0.98$  and  $Q(d=2) = 0.99$ ?

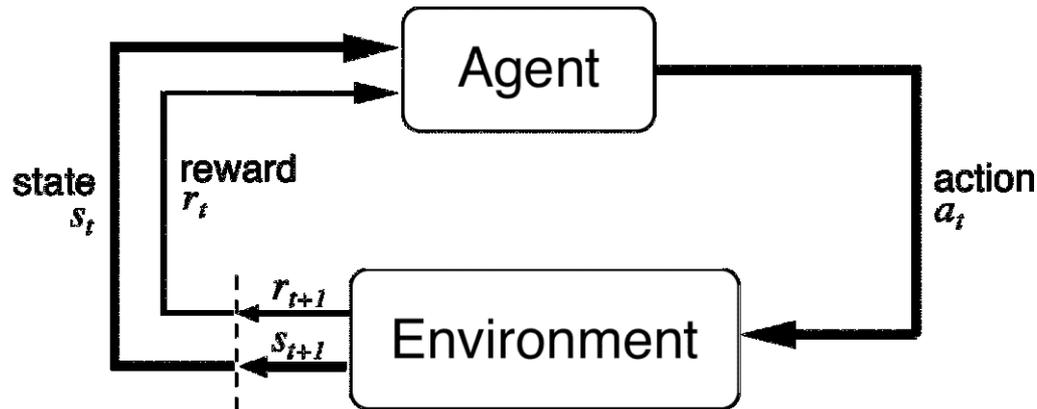
# Chapter 3: The Reinforcement Learning Problem

---

Objectives of this chapter:

- ❑ describe the RL problem we will be studying for the remainder of the course
- ❑ present idealized form of the RL problem for which we have precise theoretical results;
- ❑ introduce key components of the mathematics: value functions and Bellman equations;
- ❑ describe trade-offs between applicability and mathematical tractability.

# The Agent-Environment Interface



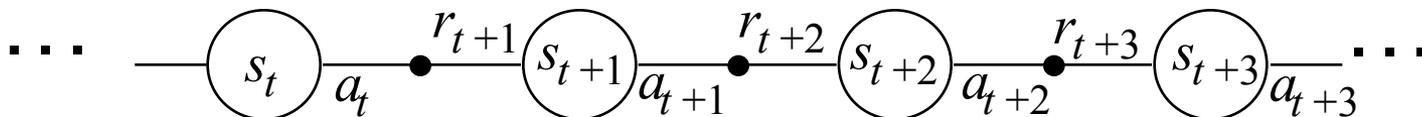
Agent and environment interact at discrete time steps:  $t = 0, 1, 2, \dots$

Agent observes state at step  $t$ :  $s_t \in \mathcal{S}$

produces action at step  $t$ :  $a_t \in A(s_t)$

gets resulting reward:  $r_{t+1} \in \mathcal{R}$

and resulting next state:  $s_{t+1}$



# The Agent Learns a Policy

---

**Policy** at step  $t$ ,  $\pi_t$  :

a mapping from states to action probabilities

$\pi_t(s, a) =$  probability that  $a_t = a$  when  $s_t = s$

- ❑ Reinforcement learning methods specify how the agent changes its policy as a result of experience.
- ❑ Roughly, the agent's goal is to get as much reward as it can over the long run.

# Getting the Degree of Abstraction Right

---

- ❑ Time steps need not refer to fixed intervals of real time.
- ❑ Actions can be low level (e.g., voltages to motors), or high level (e.g., accept a job offer), “mental” (e.g., shift in focus of attention), etc.
- ❑ States can low-level “sensations”, or they can be abstract, symbolic, based on memory, or subjective (e.g., the state of being “surprised” or “lost”).
- ❑ An RL agent is not like a *whole* animal or robot.
- ❑ Reward computation is in the agent’s environment because the agent cannot change it arbitrarily.
- ❑ The environment is not necessarily unknown to the agent, only incompletely controllable.

# Goals and Rewards

---

- ❑ Is a scalar reward signal an adequate notion of a goal?—maybe not, but it is surprisingly flexible.
- ❑ A goal should specify **what** we want to achieve, not **how** we want to achieve it.
- ❑ A goal must be outside the agent's direct control—thus outside the agent.
- ❑ The agent must be able to measure success:
  - explicitly;
  - frequently during its lifespan.

# The reward hypothesis

---

- That all of what we mean by goals and purposes can be well thought of as the maximization of the cumulative sum of a received scalar signal (reward)
- A sort of *null hypothesis*.
  - Probably ultimately wrong, but so simple we have to disprove it before considering anything more complicated

# Returns

---

Suppose the sequence of rewards after step  $t$  is :

$$r_{t+1}, r_{t+2}, r_{t+3}, \dots$$

What do we want to maximize?

In general,

we want to maximize the **expected return**,  $E\{R_t\}$ , for each step  $t$ .

**Episodic tasks:** interaction breaks naturally into episodes, e.g., plays of a game, trips through a maze.

$$R_t = r_{t+1} + r_{t+2} + \dots + r_T,$$

where  $T$  is a final time step at which a **terminal state** is reached, ending an episode.

# Returns for Continuing Tasks

---

**Continuing tasks:** interaction does not have natural episodes.

**Discounted return:**

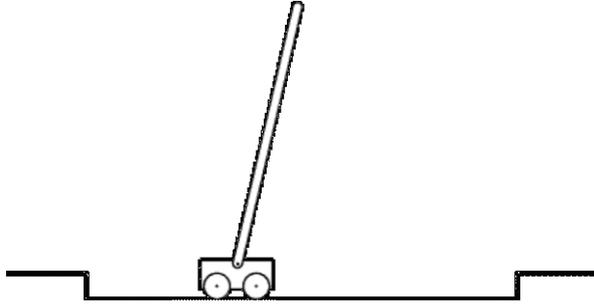
$$R_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1},$$

where  $\gamma$ ,  $0 \leq \gamma \leq 1$ , is the **discount rate**.

shortsighted  $0 \leftarrow \gamma \rightarrow 1$  farsighted

# An Example

---



Avoid **failure**: the pole falling beyond a critical angle or the cart hitting end of track.

As an **episodic task** where episode ends upon failure:

reward = +1 for each step before failure

$\Rightarrow$  return = number of steps before failure

As a **continuing task** with discounted return:

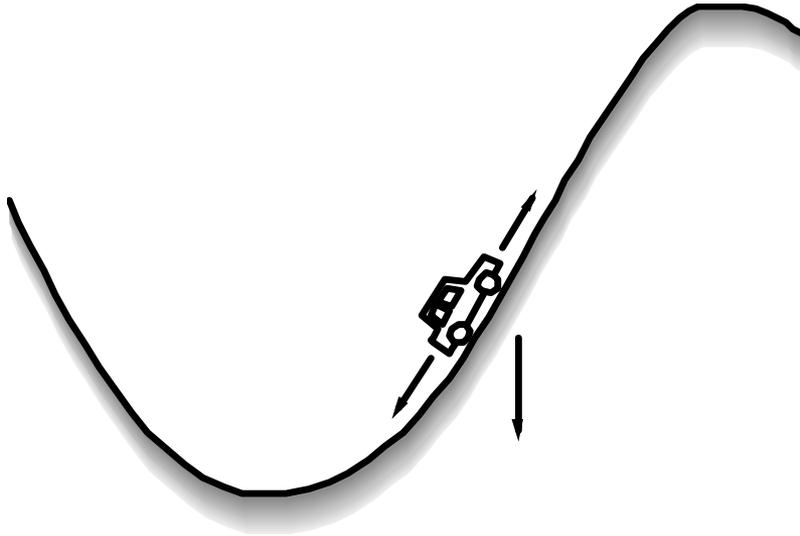
reward = -1 upon failure; 0 otherwise

$\Rightarrow$  return =  $-\gamma^k$ , for  $k$  steps before failure

In either case, return is maximized by avoiding failure for as long as possible.

# Another Example

---



Get to the top of the hill  
as quickly as possible.

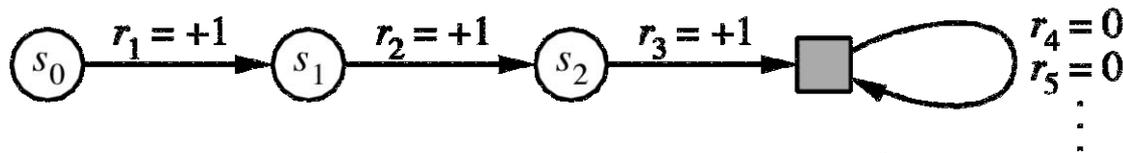
reward = -1 for each step where **not** at top of hill

⇒ return = - number of steps before reaching top of hill

Return is maximized by minimizing  
number of steps to reach the top of the hill.

# A Unified Notation

- In episodic tasks, we number the time steps of each episode starting from zero.
- We usually do not have to distinguish between episodes, so we write  $S_t$  instead of  $S_{t,j}$  for the state at step  $t$  of episode  $j$ .
- Think of each episode as ending in an absorbing state that always produces reward of zero:



- We can cover all cases by writing  $R_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$ ,

where  $\gamma$  can be 1 only if a zero reward absorbing state is always reached.

# The Markov Property

---

- By “the state” at step  $t$ , the book means whatever information is available to the agent at step  $t$  about its environment.
- The state can include immediate “sensations,” highly processed sensations, and structures built up over time from sequences of sensations.
- Ideally, a state should summarize past sensations so as to retain all “essential” information, i.e., it should have the **Markov Property**:

$$\Pr \{s_{t+1} = s', r_{t+1} = r \mid s_t, a_t, r_t, s_{t-1}, a_{t-1}, \dots, r_1, s_0, a_0\} = \Pr \{s_{t+1} = s', r_{t+1} = r \mid s_t, a_t\}$$

for all  $s', r$ , and histories  $s_t, a_t, r_t, s_{t-1}, a_{t-1}, \dots, r_1, s_0, a_0$ .