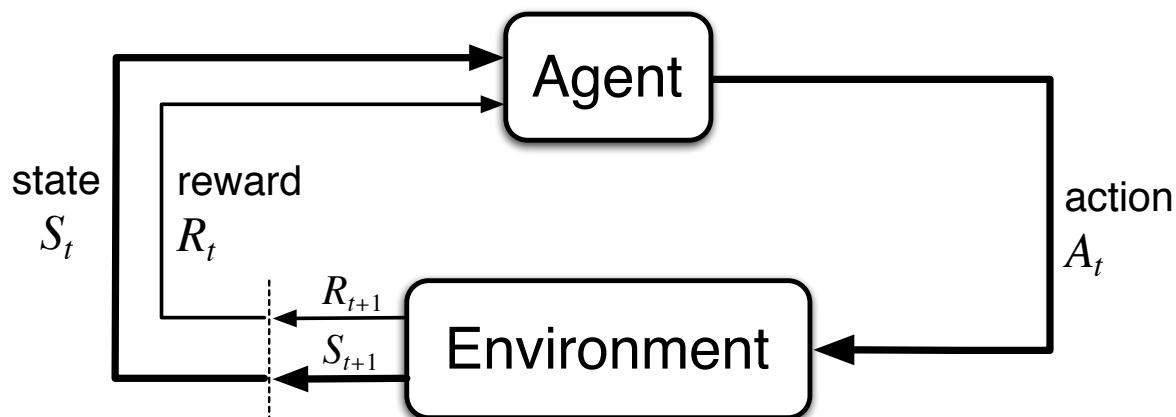# Chapter 3: The Reinforcement Learning Problem (Markov Decision Processes, or MDPs)

Objectives of this chapter:

- ❒ present Markov decision processes—an idealized form of the AI problem for which we have precise theoretical results

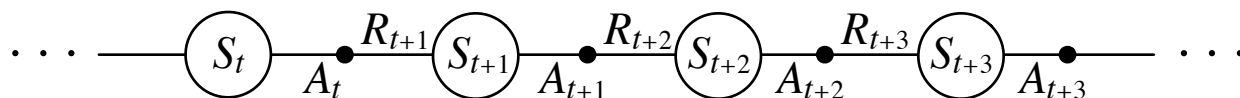- ❒ introduce key components of the mathematics: value functions and Bellman equations

Agent and environment interact at discrete time steps: $t = 0, 1, 2, 3, \ldots$

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

produces action at step $t$: $\quad A_t \in \mathcal{A}(S_t)$

gets resulting reward: $\quad R_{t+1} \in \mathcal{R} \subset \mathbb{R}$

and resulting next state: $S_{t+1} \in \mathcal{S}^+$

# Markov Decision Processes

❐ If a reinforcement learning task has the Markov Property, it is basically a **Markov Decision Process (MDP)**.

❐ If state and action sets are finite, it is a **finite MDP**.

❐ To define a finite MDP, you need to give:

- **state and action sets**

- one-step "dynamics"

$$p(s', r|s, a) = \mathbf{Pr}\{S_{t+1}=s', R_{t+1}=r \mid S_t=s, A_t=a\}$$

- there is also:

$$p(s'|s, a) \doteq \mathrm{Pr}\{S_{t+1}=s' \mid S_t=s, A_t=a\} = \sum_{r \in \mathcal{R}} p(s', r|s, a)$$

$$r(s, a) \doteq \mathbb{E}[R_{t+1} \mid S_t=s, A_t=a] = \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s', r|s, a)$$

# The Agent Learns a Policy

**Policy** at step $t$ $= \pi_t =$

a mapping from states to action probabilities

$\pi_t(a \mid s) =$ probability that $A_t = a$ when $S_t = s$

Special case - *deterministic policies*:

$\pi_t(s) =$ the action taken with prob=1 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.

# The Meaning of Life
## (goals, rewards, and returns)

# Return

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

$$R_{t+1}, R_{t+2}, R_{t+3}, \ldots$$

What do we want to maximize?

At least three cases, but in all of them,

we seek to maximize the **expected return**, $E\{G_t\}$, on each step $t$.

- <u>Total reward</u>, $G_t$ = sum of all future reward in the episode
- <u>Discounted reward</u>, $G_t$ = sum of all future *discounted* reward
- <u>Average reward</u>, $G_t$ = average reward per time step

# Episodic Tasks

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

In episodic tasks, we almost always use simple *total reward*:

$$G_t = R_{t+1} + R_{t+2} + \cdots + R_T,$$

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

# Continuing Tasks

**Continuing tasks**: interaction does not have natural episodes, but just goes on and on...

For continuing tasks we would use *discounted return*:

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1},$$

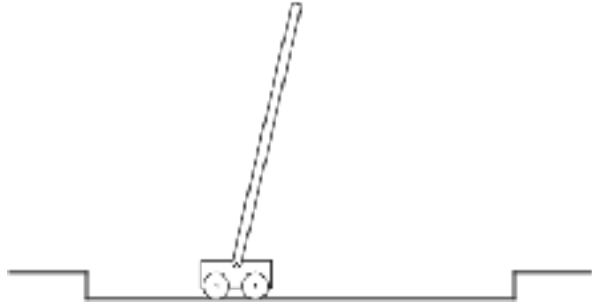where $\gamma$, $0 \leq \gamma \leq 1$, is the **discount rate**.

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

Typically, $\gamma = 0.9$

# An Example: Pole Balancing

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

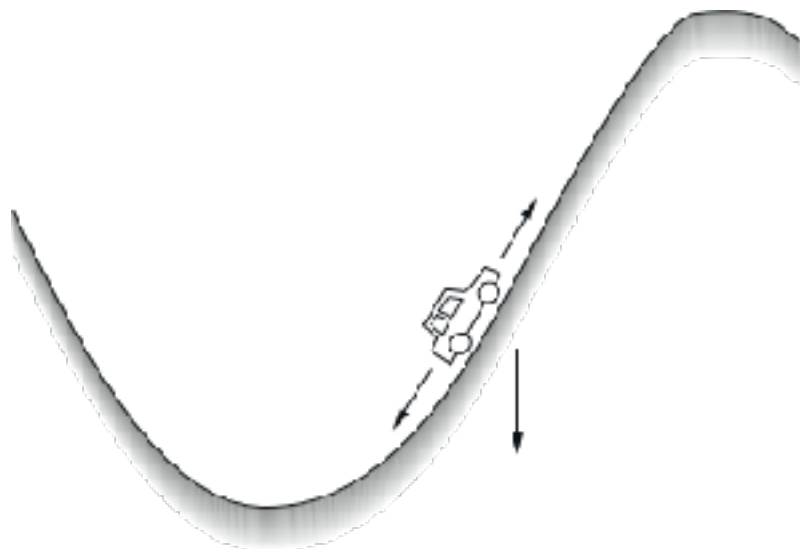As a **continuing task** with discounted return:

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: Mountain Car

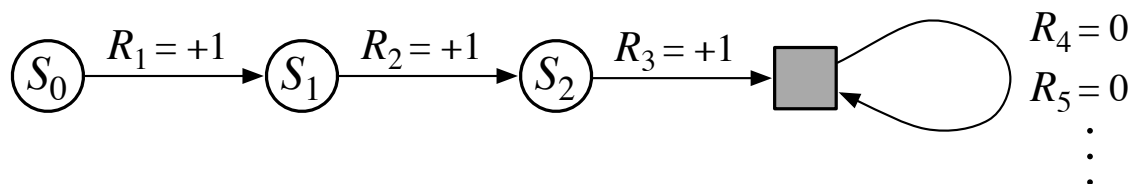Get to the top of the hill
as quickly as possible.

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

$\Rightarrow$ 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 Trick to Unify Notation for Returns

❏ In episodic tasks, we number the time steps of each episode starting from zero.

❏ We usually do not have to distinguish between episodes, so instead of writing $S_{t,j}$ for states in episode $j$, we write just $S_t$

❏ Think of each episode as ending in an absorbing state that always produces reward of zero:

$$S_0 \xrightarrow{R_1 = +1} S_1 \xrightarrow{R_2 = +1} S_2 \xrightarrow{R_3 = +1} \blacksquare \circlearrowright \quad \begin{matrix} R_4 = 0 \\ R_5 = 0 \\ \vdots \end{matrix}$$

❏ We can cover <u>all</u> cases by writing $G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$,

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

# Rewards and returns

- The objective in RL is to maximize long-term future reward

- That is, to choose $A_t$ so as to maximize $R_{t+1}, R_{t+2}, R_{t+3}, \ldots$

- But what exactly should be maximized?

- The <u>discounted *return* at time t</u>:

  the *discount rate*

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \cdots \qquad \gamma \in [0, 1)$$

| $\gamma$ | Reward sequence | Return |
|---|---|---|
| 0.5(or any) | 1 0 0 0… | 1 |
| 0.5 | 0 0 2 0 0 0… | 0.5 |
| 0.9 | 0 0 2 0 0 0… | 1.62 |
| 0.5 | -1 2 6 3 2 0 0 0… | 2 |

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \cdots \qquad \gamma \in [0, 1)$$

- Suppose $\gamma = 0.5$ and the reward sequence is

  $R_1 = 1, R_2 = 6, R_3 = -12, R_4 = 16,$ then zeros for $R_5$ and later

- What are the following returns?

  $G_4 = 0 \quad G_3 = 16 \quad G_2 = -4 \quad G_1 = 4 \quad G_0 = 3$

- Suppose $\gamma = 0.5$ and the reward sequence is all 1s.

  $$G = \frac{1}{1 - \gamma} = 2$$