



香港科技大學  
THE HONG KONG  
UNIVERSITY OF SCIENCE  
AND TECHNOLOGY

COMP 5212

Machine Learning

Lecture 23

# Reinforcement Learning

Junxian He  
May 3, 2024

# Reminder

1. Programming HW is due today
2. We will have HW4 released this week, all multi-choice questions, helping you review the contents in this semester

# Markov Decision Process

1. Start in some initial state  $s_0$
  2. For time step  $t$ :
    - a. Agent observes state  $s_t$
    - b. Agent takes action  $a_t = \pi(s_t)$       Deterministic policy
    - c. Agent receives reward  $r_t \sim p(r \mid s_t, a_t)$
    - d. Agent transitions to state  $s_{t+1} \sim p(s' \mid s_t, a_t)$
- MDPs make the **Markov assumption**: the reward and next state only depend on the current state and action.

# RL Example - gridworld

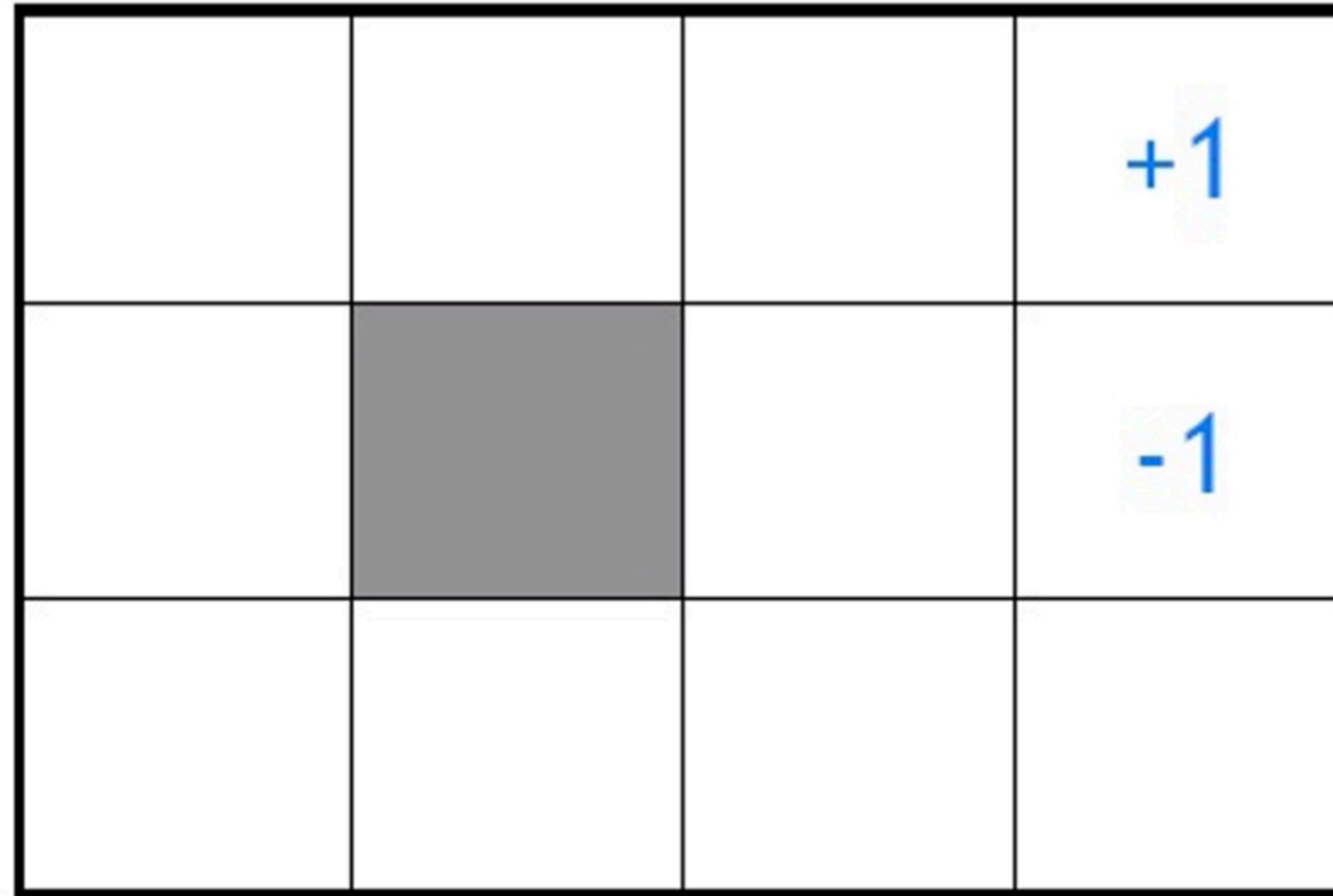
$\mathcal{S}$  = all empty squares in the grid

$\mathcal{A}$  = {up, down, left, right}

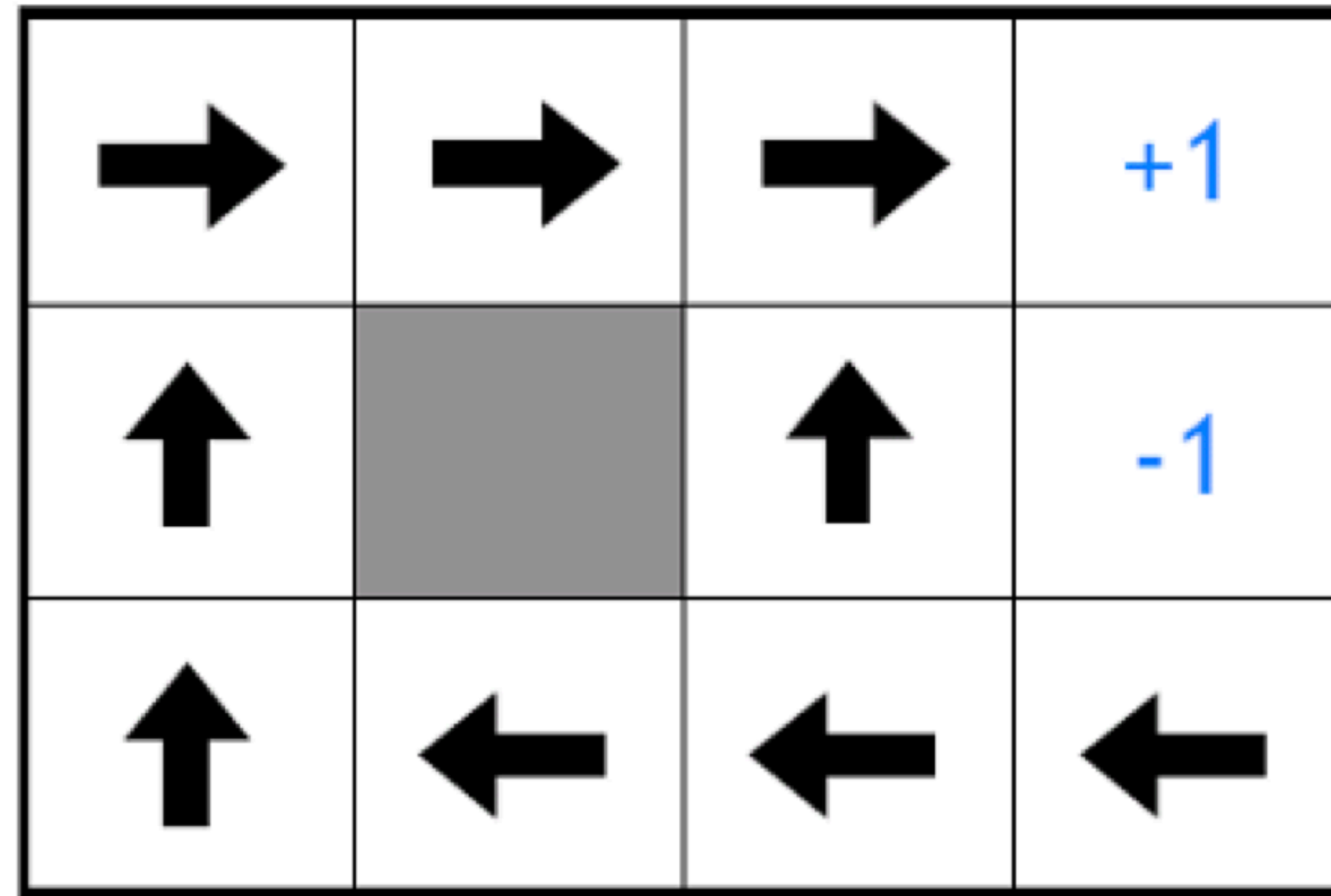
Deterministic transitions

Rewards of +1 and -1 for entering the labelled squares

Terminate after receiving either reward



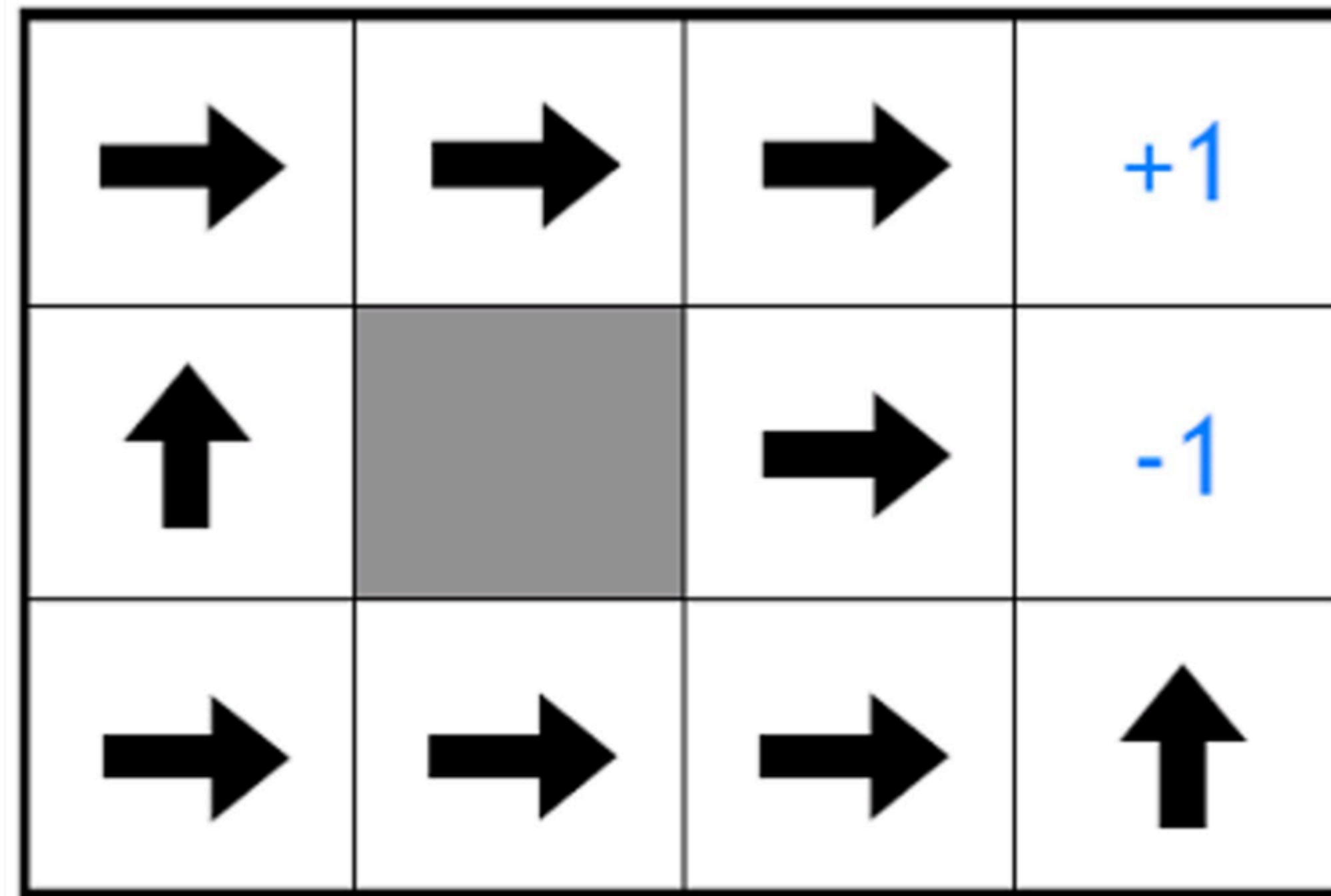
# RL Example - gridworld



Is this policy optimal?

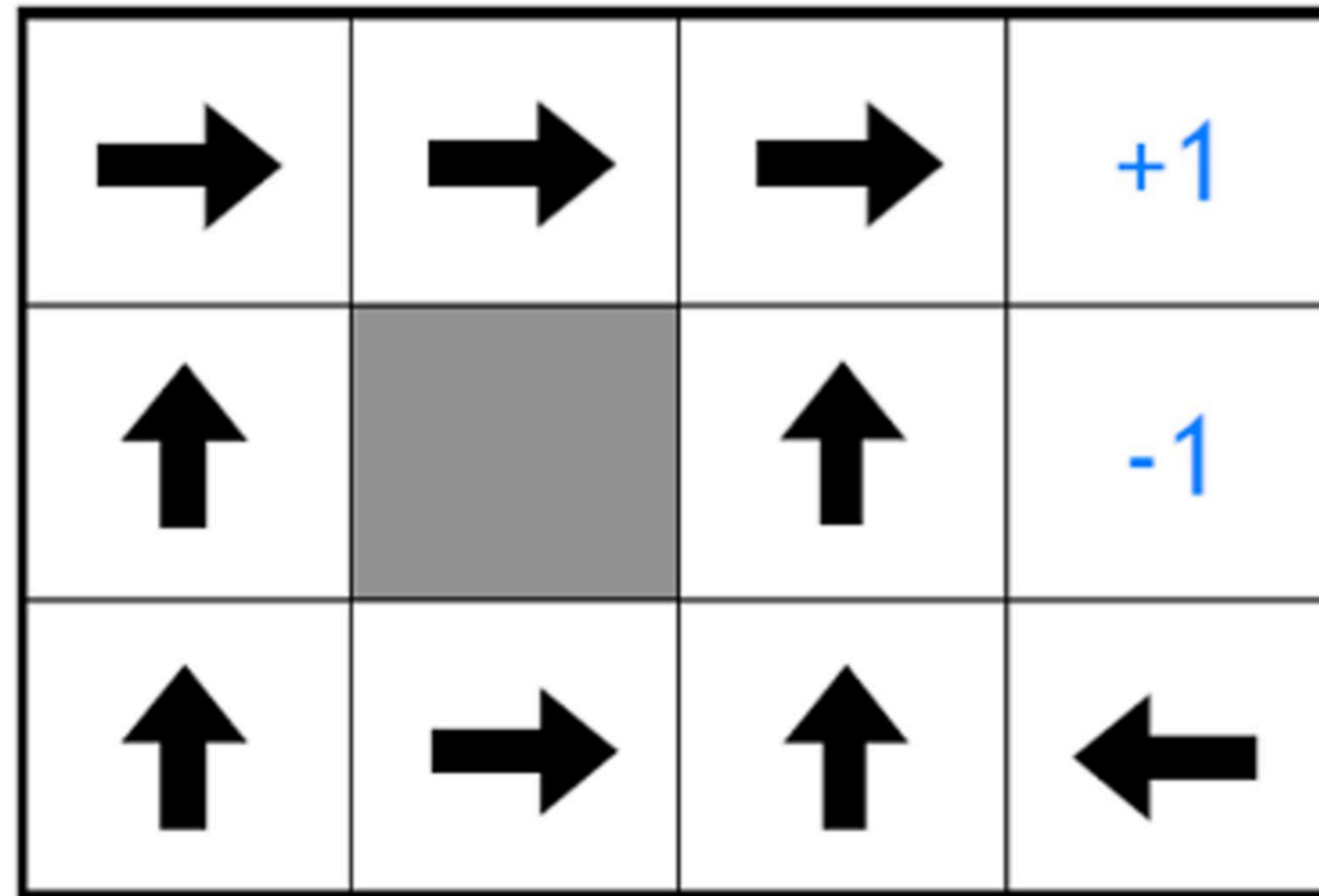
# RL Example - gridworld

Optimal policy given a reward of -2 per step



# RL Example - gridworld

Optimal policy given a reward of -0.5 per step



What would be the algorithm to find the optimal policy automatically?

# Discounted Reward

Total reward is  $\sum_{t=0}^{\infty} \gamma^t r_t = r_0 + \gamma r_1 + \gamma^2 r_2 + \gamma^3 r_3 + \dots$

where  $0 < \gamma < 1$  is some discount factor for future rewards

Why discount?

- Mathematically tractable – total reward doesn't explode

$$1 + 1 + 1 + \dots = \infty \text{ but } 1 + 0.8*1 + (0.8)^2*1 + \dots = 5$$

- Actions don't have lasting impact



# Key Challenges

- The algorithm has to gather its own training data
- The outcome of taking some action is often stochastic or unknown until after the fact
- Decisions can have a delayed effect on future outcomes (exploration-exploitation tradeoff)

**explore** decisions whose reward is uncertain

**exploit** decisions which give high reward

# RL: Objective function

- Find a policy  $\pi^* = \operatorname{argmax}_{\pi} V^{\pi}(s) \forall s \in \mathcal{S}$
- $V^{\pi}(s) = \mathbb{E}[\text{discounted total reward of starting in state } s \text{ and executing policy } \pi \text{ forever}]$

$$= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}[R(s_t, \pi(s_t))]$$

where  $0 < \gamma < 1$  is some discount factor for future rewards

# Value Function

## Bellman equations

$$V^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s_1 \in \mathcal{S}} p(s_1 | s, \pi(s)) V^\pi(s_1)$$

Recursive form

Immediate  
reward

Expected (Discounted)  
Future reward

# Solve Value Function

$$V^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s_1 \in \mathcal{S}} p(s_1 | s, \pi(s)) V^\pi(s_1)$$

Given  $R$ , transition function  $p$ , and policy  $\pi(s)$ , we can utilize this equation to solve  $V(s)$  for any  $s$

How?

Suppose the state size is finite  $|S|$ , you have  $|S|$  *linear* equations with  $|S|$  variables

# Optimal value function and policy

- Optimal value function:

$$V^*(s) = \max_{a \in \mathcal{A}} [R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) V^*(s')]$$

- System of  $|\mathcal{S}|$  equations and  $|\mathcal{S}|$  variables – nonlinear!

- Optimal policy:

$$\pi^*(s) = \operatorname{argmax}_{a \in \mathcal{A}} \underbrace{R(s, a)}_{\text{Immediate reward}} + \underbrace{\gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) V^*(s')}_{\text{Expected (Discounted) Future reward}}$$

- Insight: if you know the optimal value function, you can solve for the optimal policy!

# Value Iteration

- Inputs:  $R(s, a)$ ,  $p(s' | s, a)$ ,  $0 < \gamma < 1$
- Initialize  $V^{(0)}(s) = 0 \forall s \in \mathcal{S}$  (or randomly) and set  $t = 0$
- While not converged, do:

- For  $s \in \mathcal{S}$

$$V^{(t+1)}(s) \leftarrow \max_{a \in \mathcal{A}} \underbrace{\left[ R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) V^{(t)}(s') \right]}_{Q(s, a)}$$

- $t = t + 1$

$$\pi^*(s) \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} \left[ R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) V^{(t)}(s') \right]$$

After finding the optimal value function, we can find the optimal policy

# Value Iteration: Convergence

**Theorem 1:** Value function convergence

$V$  will converge to  $V^*$  if each state is “visited” infinitely often (Bertsekas, 1989)

# Policy Iteration

- Inputs:  $R(s, a)$ ,  $p(s' | s, a)$ ,  $0 < \gamma < 1$
- Initialize  $\pi$  randomly
- While not converged, do:
  - Solve the Bellman equations defined by policy  $\pi$  Linear equation system

$$V^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, \pi(s)) V^\pi(s')$$

- Update  $\pi$

$$\pi(s) \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) V^\pi(s')$$

Both value iteration and policy iteration are standard algorithms for solving MDPs, there isn't universal agreement over which is better



# Learning a Model for an MDP

State transition  $p(s' | s, a)$  and reward function  $R(s, a)$  are unknown in practice

Suppose we have a number of trials:

$$\begin{aligned} s_0^{(1)} &\xrightarrow{a_0^{(1)}} s_1^{(1)} \xrightarrow{a_1^{(1)}} s_2^{(1)} \xrightarrow{a_2^{(1)}} s_3^{(1)} \xrightarrow{a_3^{(1)}} \dots \\ s_0^{(2)} &\xrightarrow{a_0^{(2)}} s_1^{(2)} \xrightarrow{a_1^{(2)}} s_2^{(2)} \xrightarrow{a_2^{(2)}} s_3^{(2)} \xrightarrow{a_3^{(2)}} \dots \\ &\dots \end{aligned}$$

$$P_{sa}(s') = \frac{\text{\#times took we action } a \text{ in state } s \text{ and got to } s'}{\text{\#times we took action } a \text{ in state } s}$$

Similarly we can estimate  $R(s, a)$  to be the the average reward observed at state  $s$  with action  $a$

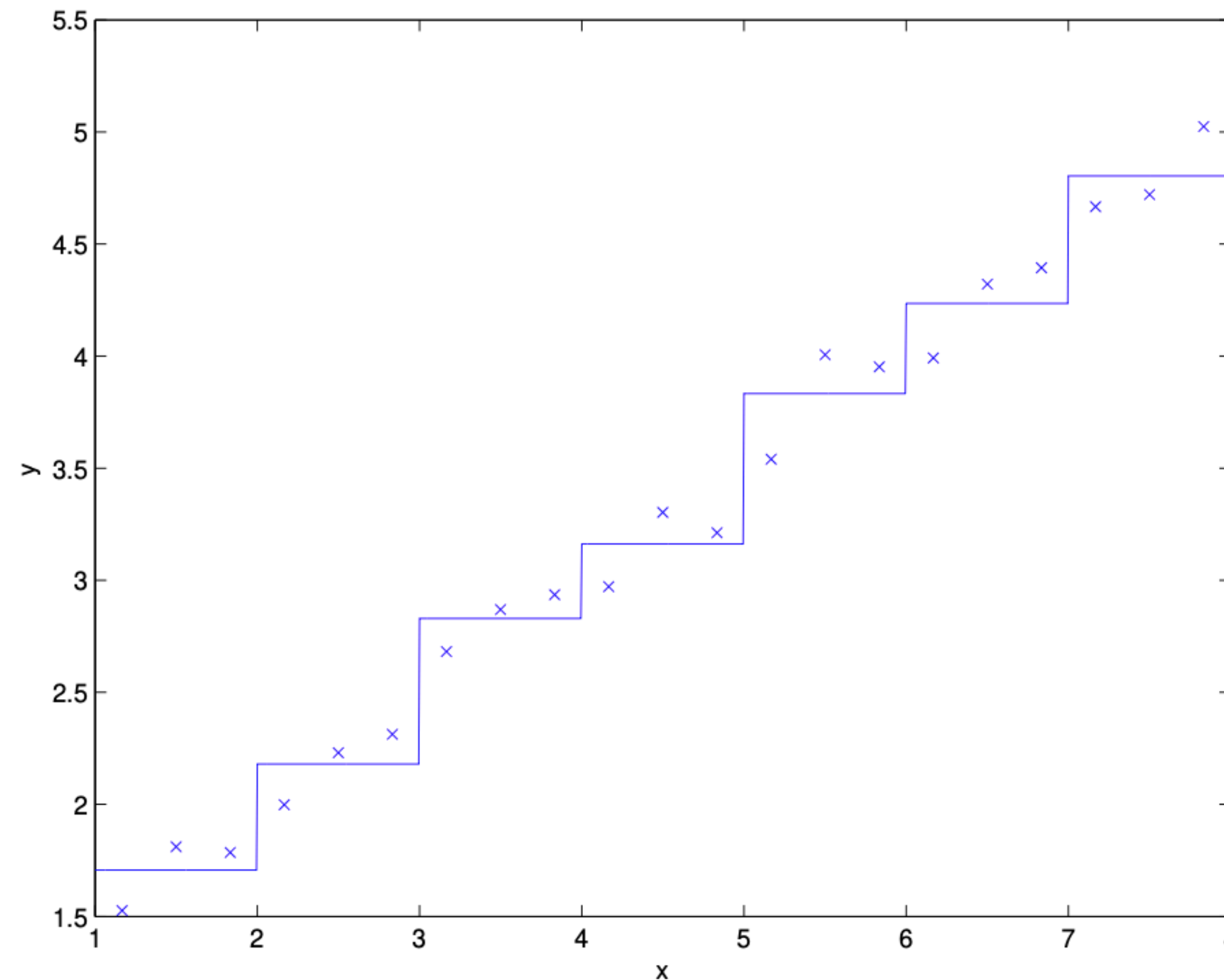
# Finding the Optimal Policy

1. Initialize  $\pi$  randomly.
2. Repeat {
  - (a) Execute  $\pi$  in the MDP for some number of trials.
  - (b) Using the accumulated experience in the MDP, update our estimates for  $P_{sa}$  (and  $R$ , if applicable).
  - (c) Apply value iteration with the estimated state transition probabilities and rewards to get a new estimated value function  $V$ .
  - (d) Update  $\pi$  to be the greedy policy with respect to  $V$ .}

# Continuous State

Continuous states are common, e.g., using  $(x, y)$  coordinates and velocity to express the state of a car

## Discretization



Dimensions are high when we have several states

# Value Function Approximation

Learn a function  $f_{\theta}(s) : s \rightarrow V(s)$

Similar to supervised learning,  $f_{\theta}(s)$  could be a neural network

1. Randomly sample  $n$  states  $s^{(1)}, s^{(2)}, s^{(3)} \dots$   
Want to compose a training dataset, next is to estimate  $V(s)$  for these  $n$  points
2. For every state, repeat value iteration to approximate  $V(s^{(i)})$
3. Now we have a supervised dataset to train  $f_{\theta}(s)$

# Learning a Proxy Model

Policy is a function parameterized by  $\theta$ :  $\pi_\theta$

$\tau = (s_0, a_0, s_1, a_1, \dots, s_{T-1}, a_{T-1}, s_T)$  is the trajectory from  $\pi_\theta$

$\tau$  is a random variable

Total payoff for the policy is:

$$\eta_\theta = \mathbb{E}_{\tau \sim p_\theta} \left[ \sum_{t=0}^{T-1} \gamma^t R(s_t, a_t) \right]$$

$p_\theta$  contains the policy and the transition  $p(s' | s, a)$

We want to optimize  $\theta$  to maximize  $\eta$

# Policy Gradient

$$\eta_{\theta} = \mathbb{E}_{\tau \sim p_{\theta}} \left[ \sum_{t=0}^{T-1} \gamma^t R(s_t, a_t) \right]$$

We define  $f(\tau) = \sum_{t=0}^{T-1} \gamma^t R(s_t, a_t)$ .

$$\eta(\theta) = \mathbb{E}_{\tau \sim P_{\theta}} [f(\tau)]$$

Connection to VAE? Reparameterization trick?

# Policy Gradient

$$\begin{aligned}\nabla_{\theta} \mathbf{E}_{\tau \sim P_{\theta}} [f(\tau)] &= \nabla_{\theta} \int P_{\theta}(\tau) f(\tau) d\tau \\ &= \int \nabla_{\theta} (P_{\theta}(\tau) f(\tau)) d\tau \quad (\text{swap integration with gradient}) \\ &= \int (\nabla_{\theta} P_{\theta}(\tau)) f(\tau) d\tau \quad (\text{because } f \text{ does not depend on } \theta) \\ &= \int P_{\theta}(\tau) (\nabla_{\theta} \log P_{\theta}(\tau)) f(\tau) d\tau \\ &= \mathbf{E}_{\tau \sim P_{\theta}} [(\nabla_{\theta} \log P_{\theta}(\tau)) f(\tau)]\end{aligned}$$

Can be approximated using MC sampling

Policy gradient is a commonly used method to propagate gradients through discrete variables

# Policy Gradient

$$\mathbf{E}_{\tau \sim P_\theta} [(\nabla_\theta \log P_\theta(\tau)) f(\tau)]$$

What is  $\nabla_\theta \log P_\theta(\tau)$

$$P_\theta(\tau) = \mu(s_0) \pi_\theta(a_0|s_0) P_{s_0 a_0}(s_1) \pi_\theta(a_1|s_1) P_{s_1 a_1}(s_2) \cdots P_{s_{T-1} a_{T-1}}(s_T)$$

$$\nabla_\theta \log P_\theta(\tau) = \nabla_\theta \log \pi_\theta(a_0|s_0) + \nabla_\theta \log \pi_\theta(a_1|s_1) + \cdots + \nabla_\theta \log \pi_\theta(a_{T-1}|s_{T-1})$$

$$\begin{aligned} \nabla_\theta \eta(\theta) &= \nabla_\theta \mathbf{E}_{\tau \sim P_\theta} [f(\tau)] = \mathbf{E}_{\tau \sim P_\theta} \left[ \left( \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t|s_t) \right) \cdot f(\tau) \right] \\ &= \mathbf{E}_{\tau \sim P_\theta} \left[ \left( \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t|s_t) \right) \cdot \left( \sum_{t=0}^{T-1} \gamma^t R(s_t, a_t) \right) \right] \end{aligned}$$



# A Lemma

$$\mathbb{E}_{x \sim P_\theta} [\nabla_\theta \log P_\theta(x)] = 0.$$

$$\mathbb{E}_{a_t \sim \pi_\theta(\cdot | s_t)} \nabla_\theta \log \pi_\theta(a_t | s_t) = 0$$

$$\mathbb{E}_{\tau \sim P_\theta} \left[ \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t) \right] = 0$$

# Policy Gradient

$$\begin{aligned}\nabla_{\theta}\eta(\theta) &= \sum_{t=0}^{T-1} \mathbb{E}_{\tau \sim P_{\theta}} \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \cdot \left( \sum_{j=0}^{T-1} \gamma^j R(s_j, a_j) \right) \right] \\ &= \sum_{t=0}^{T-1} \mathbb{E}_{\tau \sim P_{\theta}} \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \cdot \left( \sum_{j \geq t} \gamma^j R(s_j, a_j) \right) \right]\end{aligned}$$

Loss does not mean much, and you should only care about the return

# Learning Policy and Value Function Together

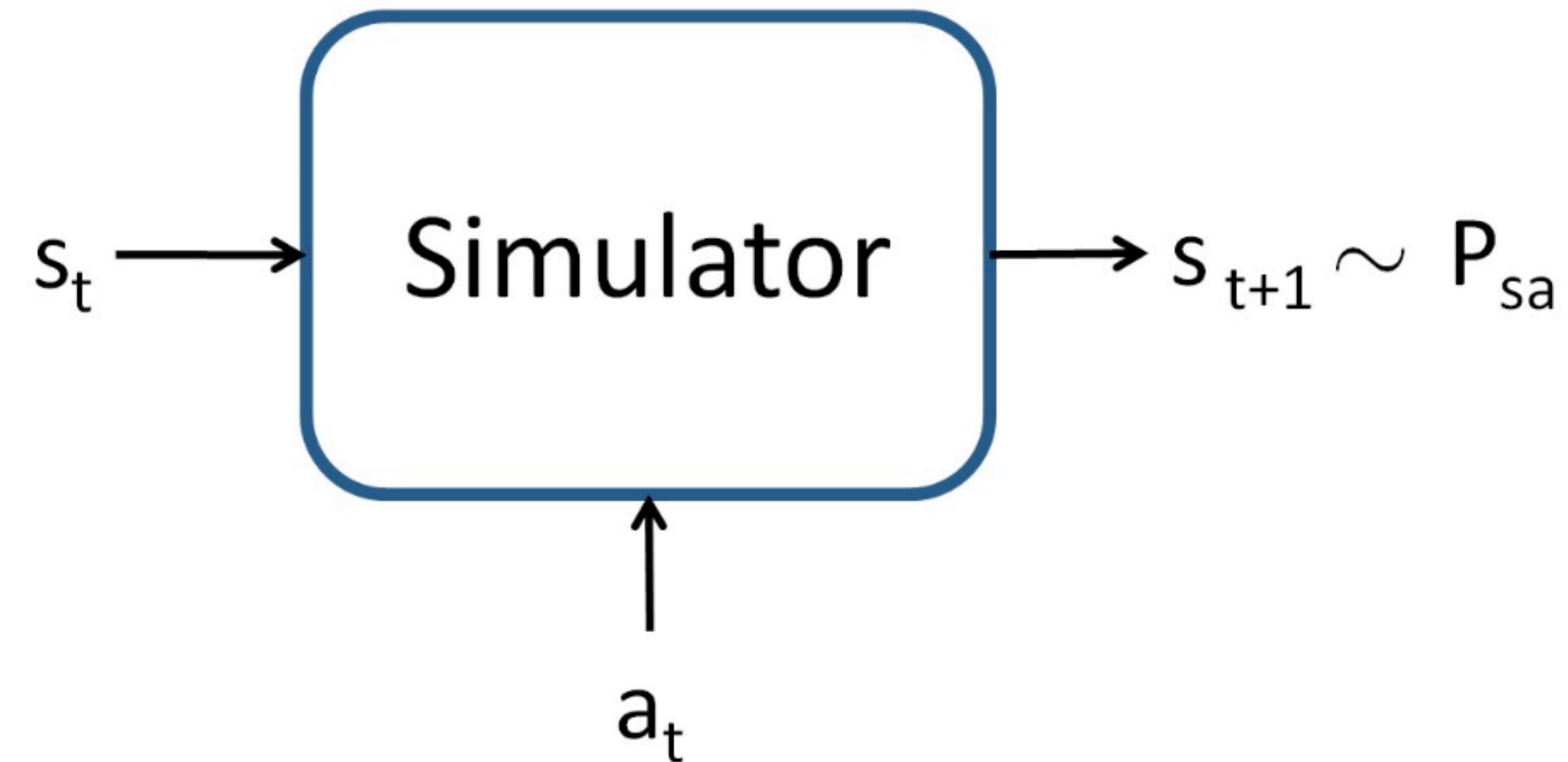
Repeat{

1. Perform a number of trials from policy  $\pi_\theta$  to get all the trajectory
2. Update the policy with the current value function
3. Compute the expected reward for each state in the trajectories
4. Supervised training to train the value function

}

Reward models are often trained in advance

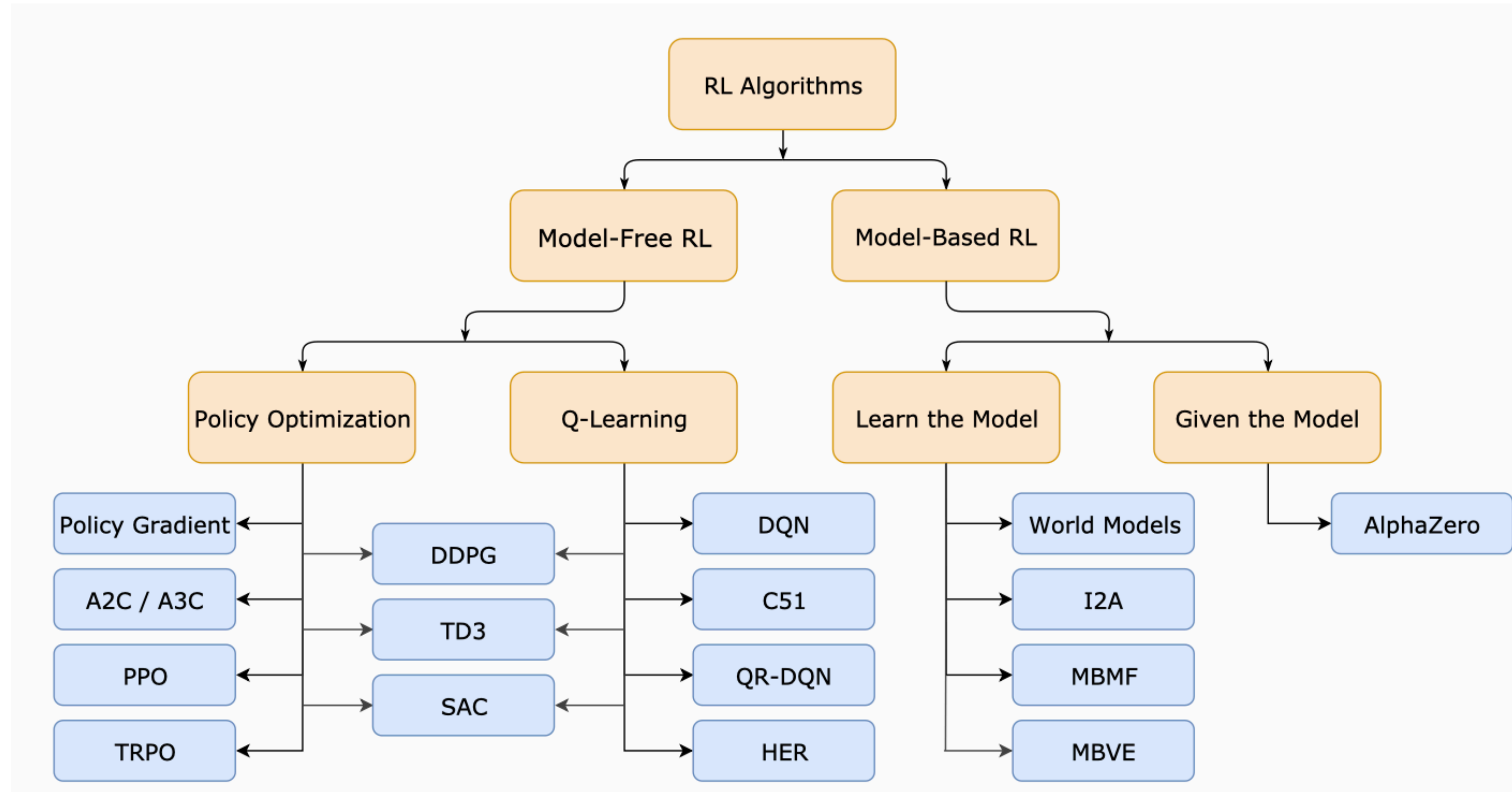
# Model for the Environment



## Model-based Reinforcement Learning

1. Interaction with real environment can be slow
2. Interaction with real environment can be risky

# Taxonomy of RL



[https://spinningup.openai.com/en/latest/spinningup/rl\\_intro2.html#part-2-kinds-of-rl-algorithms](https://spinningup.openai.com/en/latest/spinningup/rl_intro2.html#part-2-kinds-of-rl-algorithms)