

# Machine Learning I

## MATH60629

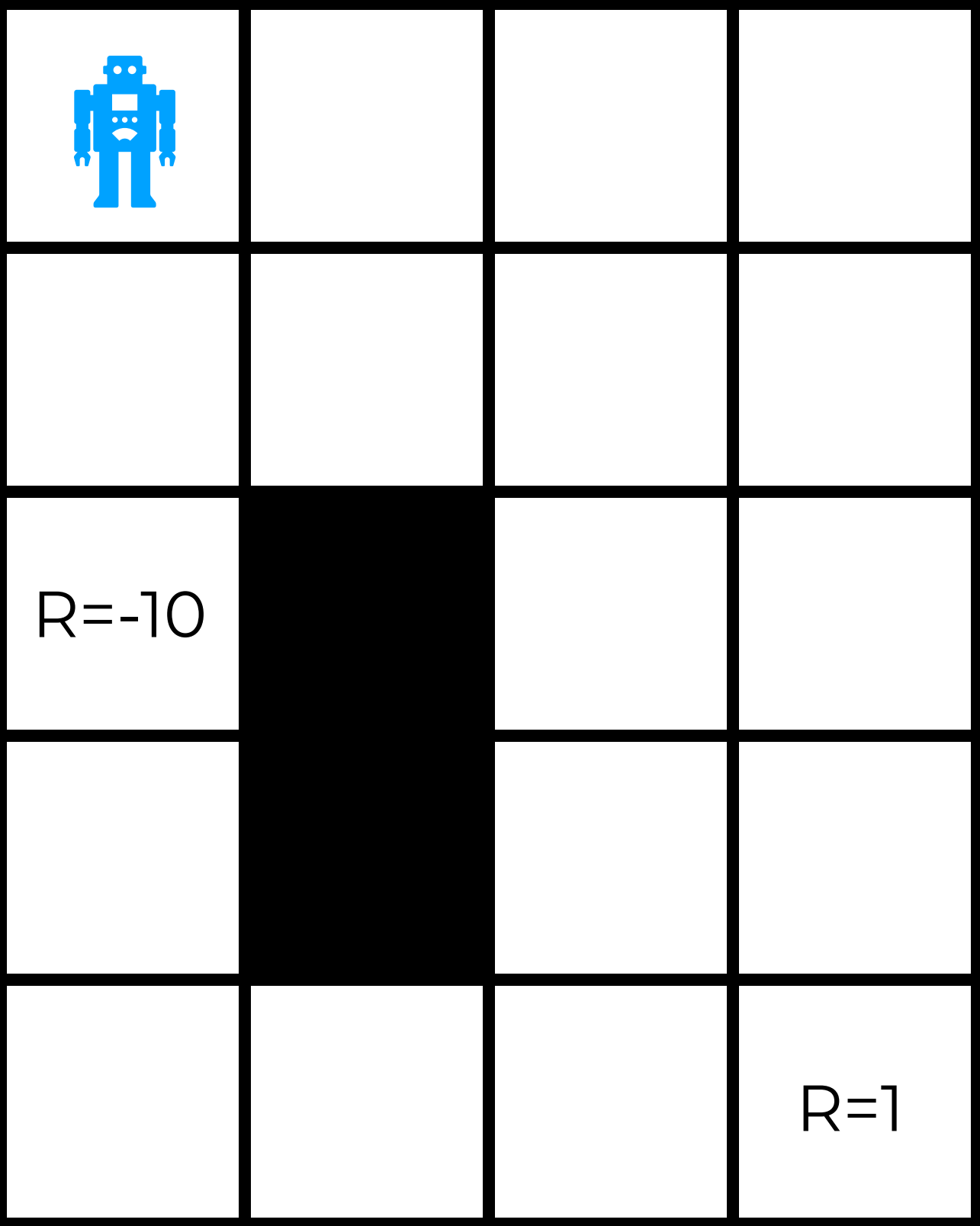
Sequential Decision Making I  
**Summary**  
— Week #12

# Three main components

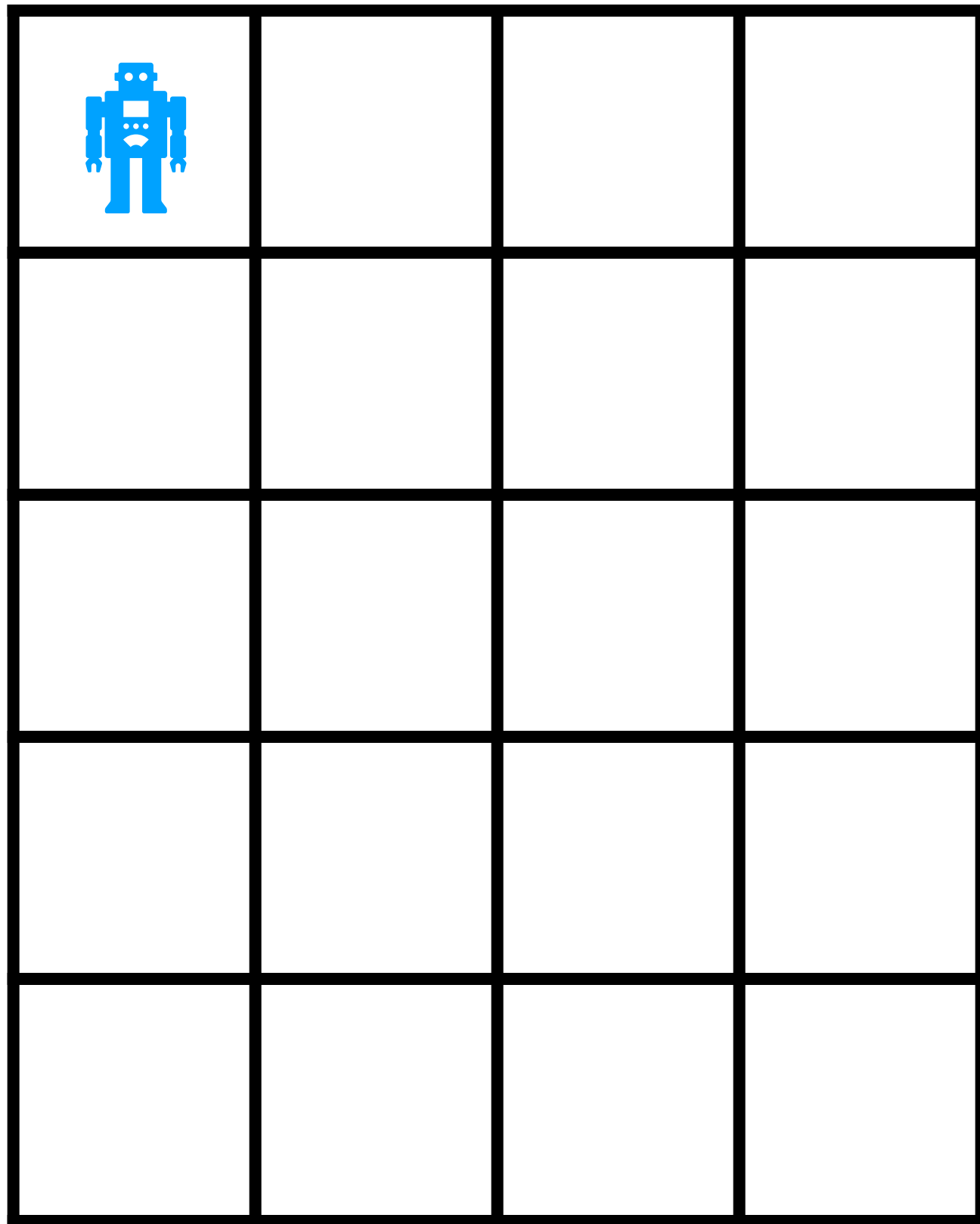
- **Task (T)**
- **Performance measure (P)**
- **Experience (E)**

# Supervised learning

- Experience a fixed data set
- Fit a model using this data
- Use the model to make predictions about unseen data
- Eventually: Predictions may be used downstream

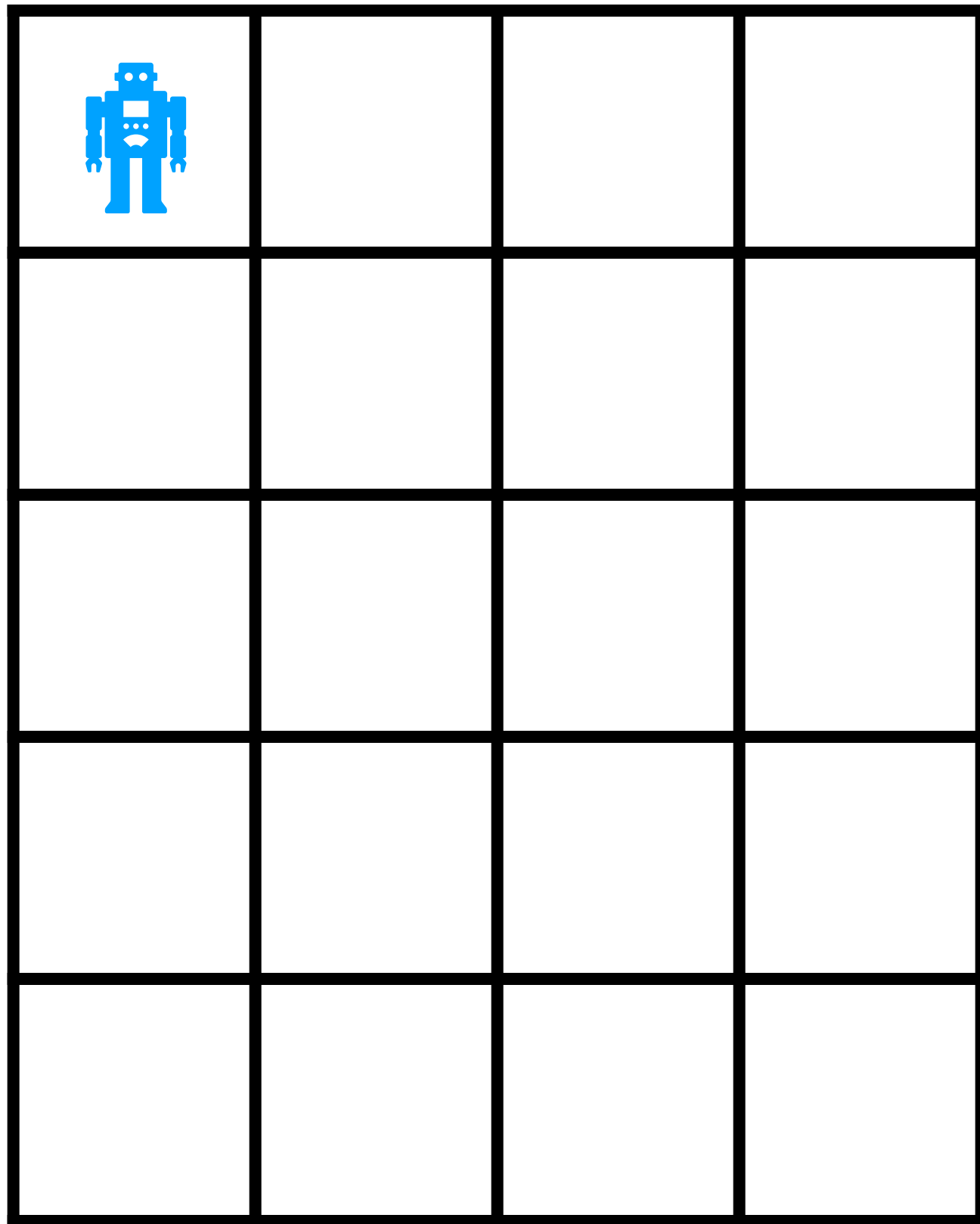


# Initial example with grid world



- Each cell is a state ( $S$ )
- Actions indicate which movements are possible:  $A := \{L, R, U, D\}$
- Rewards encode the task:  $R(s)$
- Transition probabilities encode the outcome of an action:  $P(s' | s, a)$

# Initial example with grid world



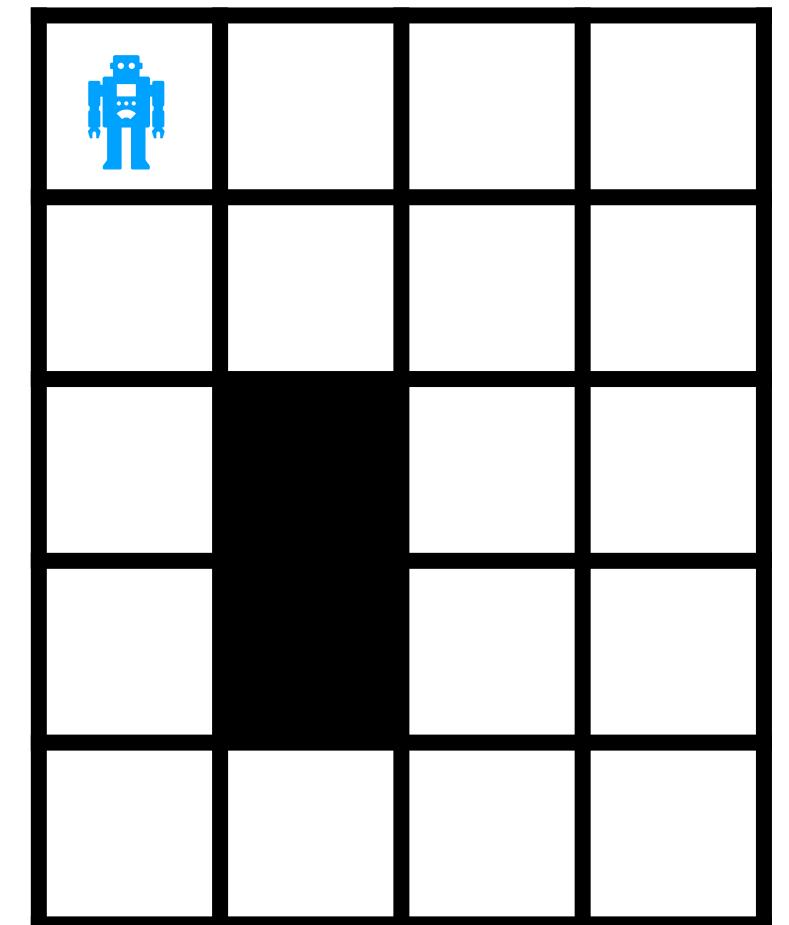
- Each cell is a state ( $S$ )
- Actions indicate which movements are possible:  $A := \{L, R, U, D\}$
- Rewards encode the task:  $R(s)$
- Transition probabilities encode the outcome of an action:  $P(s' | s, a)$

## Planning

This week we discuss a version of RL where these are observed

# Markov Decision Process (MDP)

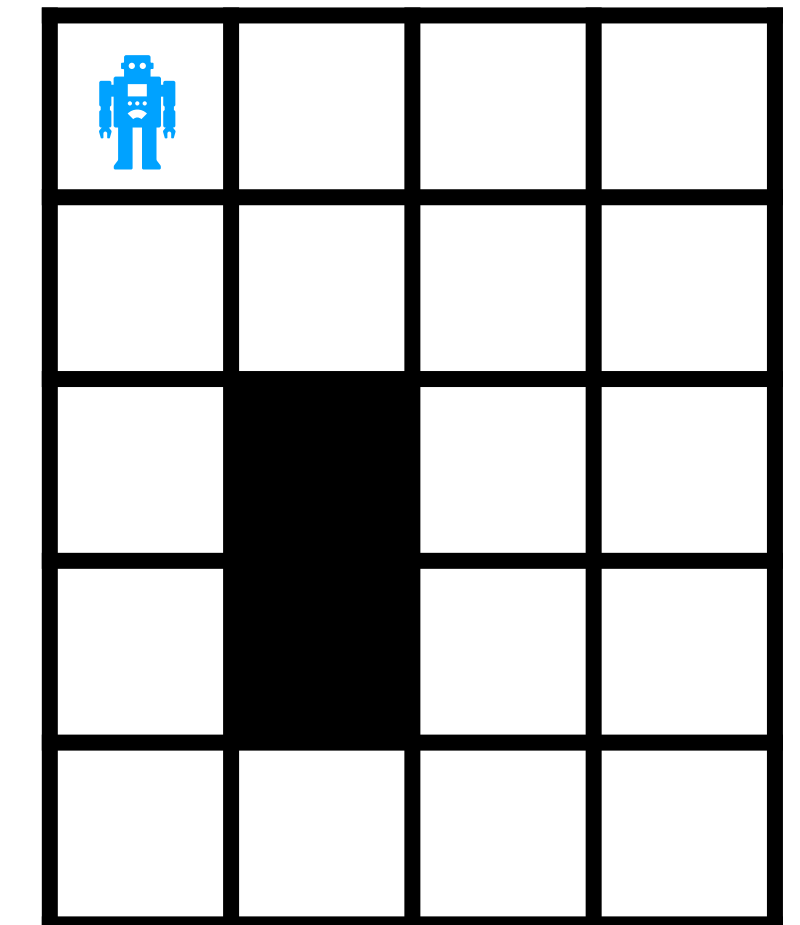
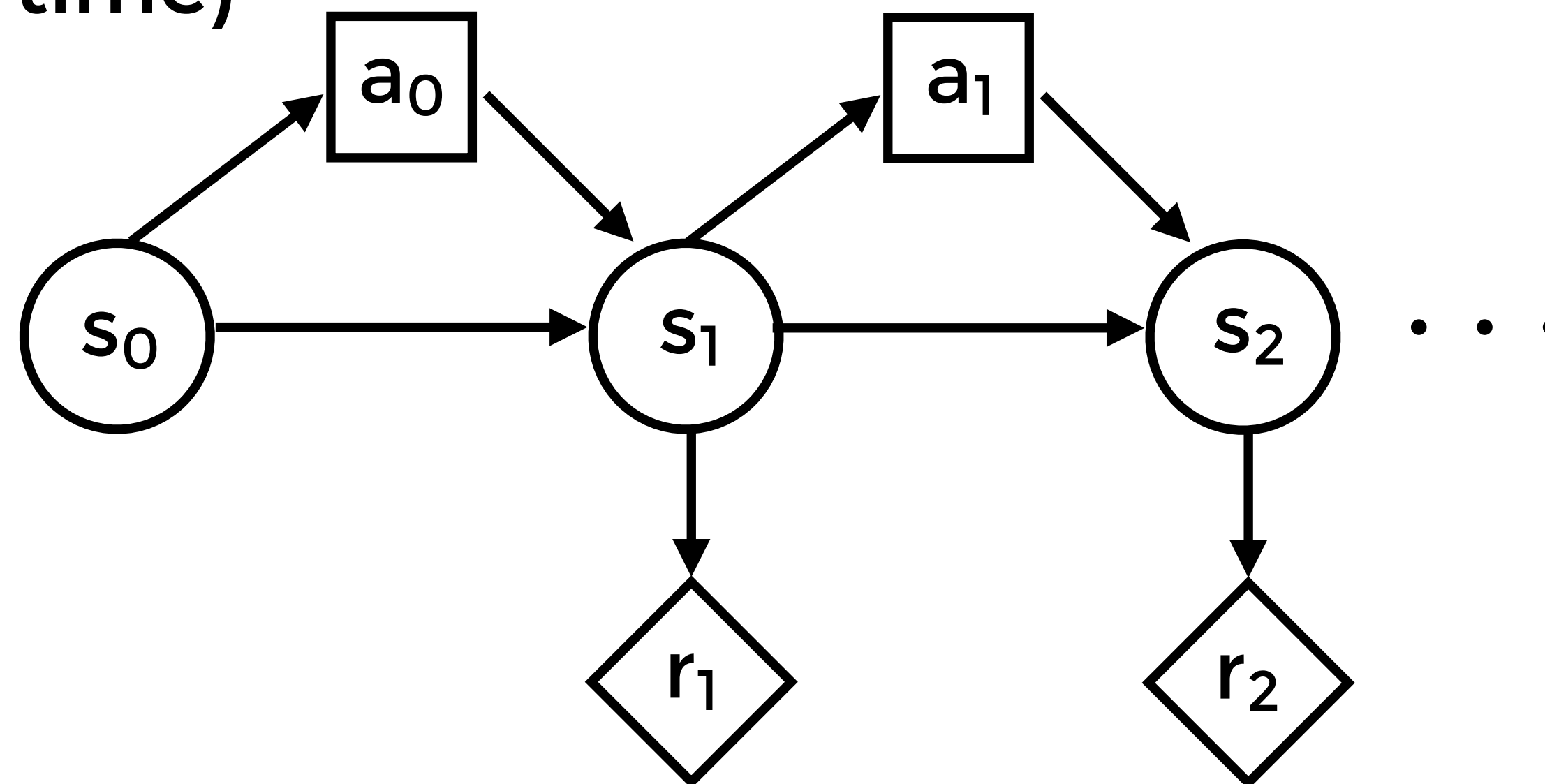
- Provide a framework for decision-making under uncertainty
- Markov process with decisions and utilities
- Assumes stationarity (i.e., transitions are fixed across time)



# Markov Decision Process (MDP)

- Provide a framework for decision-making under uncertainty
- Markov process with decisions and utilities
- Assumes stationarity (i.e., transitions are fixed across time)

- Square nodes: decisions
- Circle nodes: States
- Diamond nodes: utility

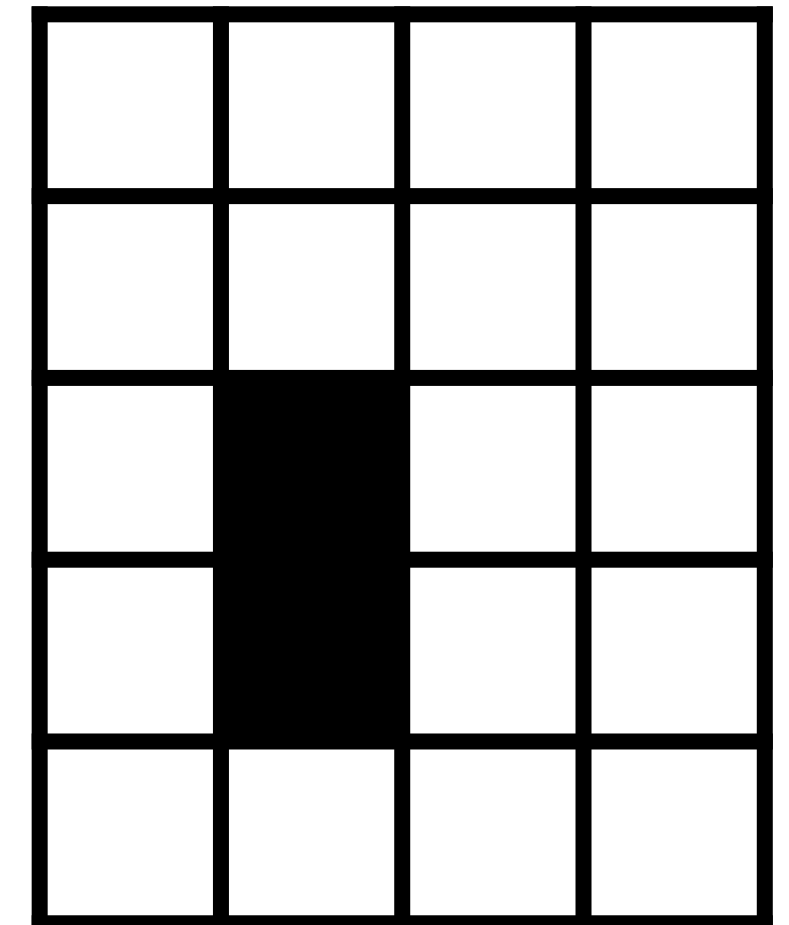




# Markov Decision Process (MDP)

$\langle A, S, P, R, \gamma \rangle$

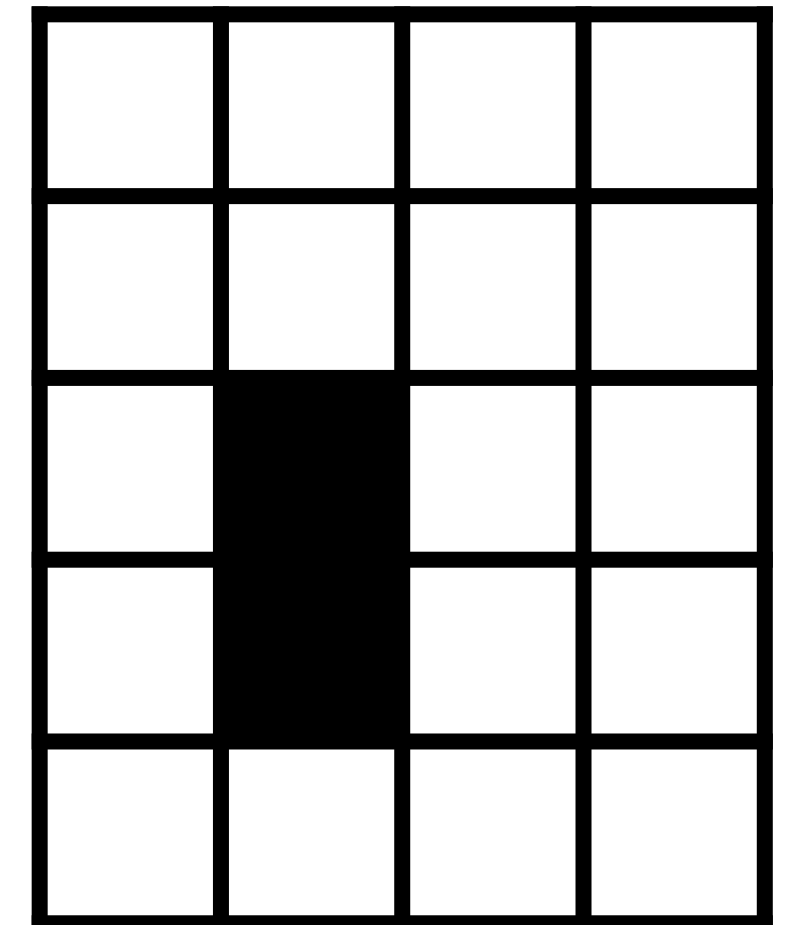
- **A**: set of actions
- $P(S' | S, A)$ : transition probabilities
- $R(S)$ : reward function
- $\gamma$  : discount factor  $\in [0, 1]$



# Markov Decision Process (MDP)

$\langle A, S, P, R, \gamma \rangle$

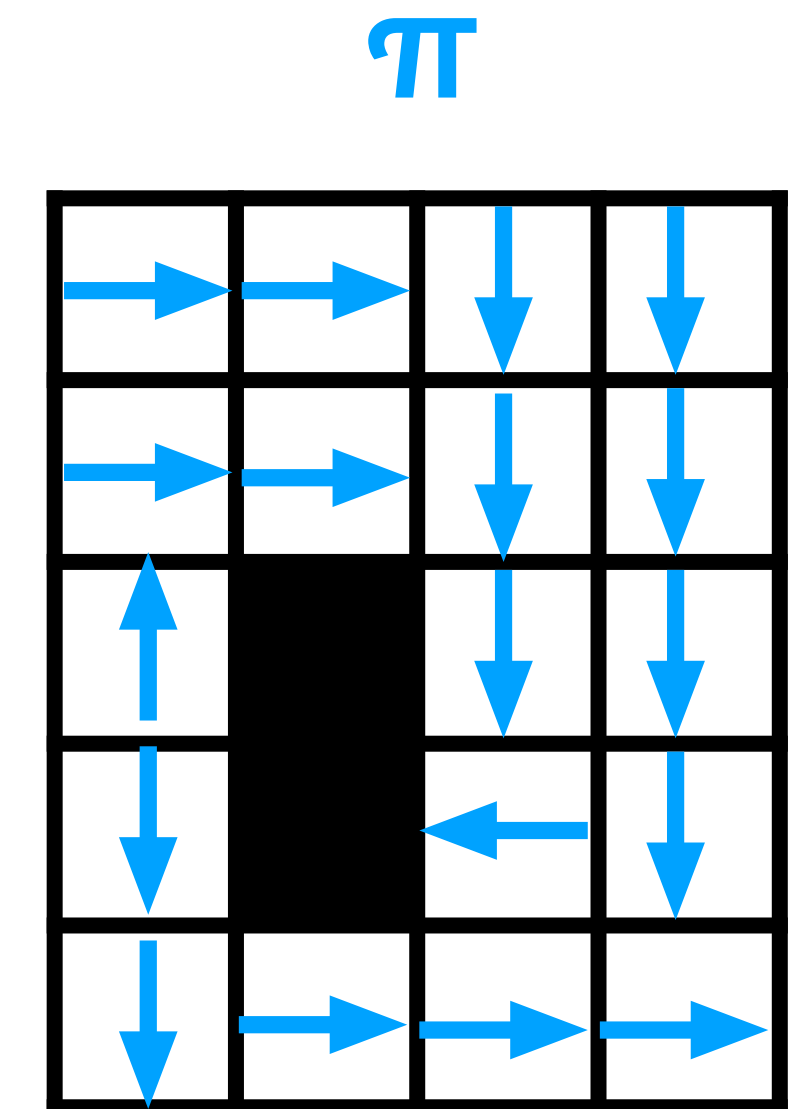
- **A**: set of actions
- $P(S' | S, A)$ : transition probabilities
- $R(S)$ : reward function
- $\gamma$  : discount factor  $\in [0, 1]$
- A policy:  $\pi : S \rightarrow A$



# Markov Decision Process (MDP)

$$\langle \mathbf{A}, \mathbf{S}, \mathbf{P}, \mathbf{R}, \gamma \rangle$$

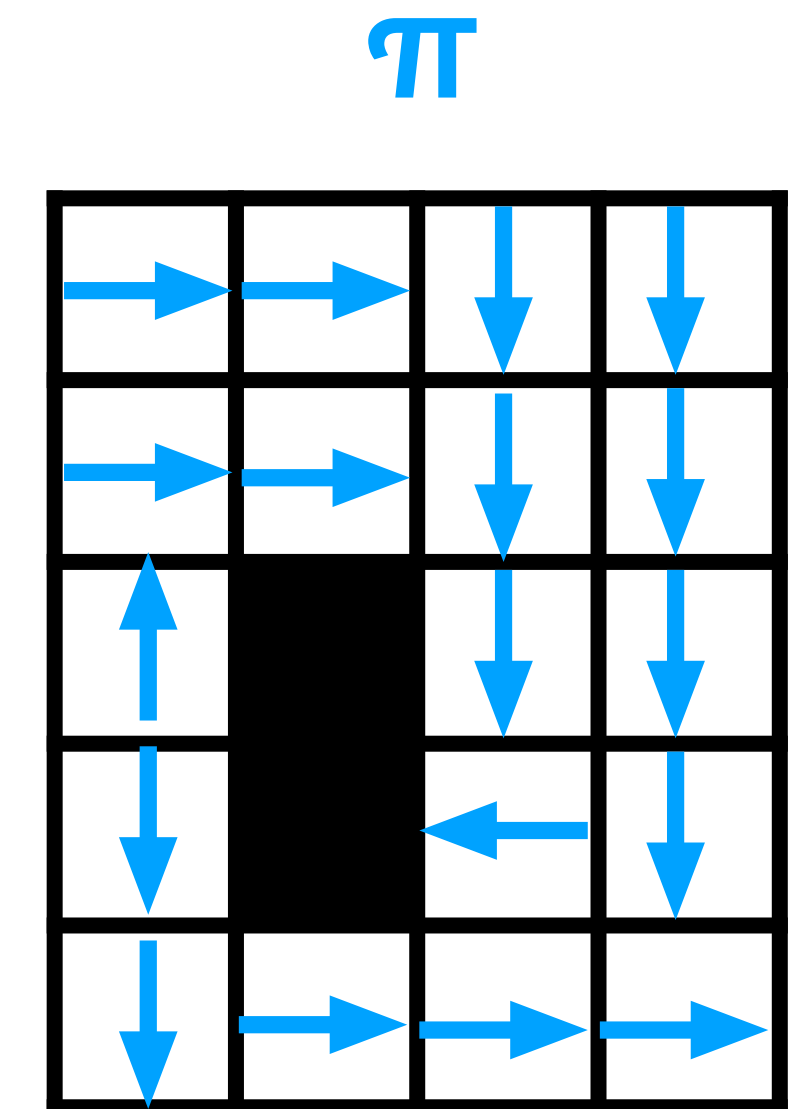
- $\mathbf{A}$ : set of actions
- $P(S' | S, A)$ : transition probabilities
- $R(S)$ : reward function
- $\gamma$  : discount factor  $\in [0, 1]$
- A policy:  $\pi : \mathbf{S} \rightarrow \mathbf{A}$



# Markov Decision Process (MDP)

$$\langle \mathbf{A}, \mathbf{S}, \mathbf{P}, \mathbf{R}, \gamma \rangle$$

- $\mathbf{A}$ : set of actions
- $P(S' | S, A)$ : transition probabilities
- $R(S)$ : reward function
- $\gamma$  : discount factor  $\in [0, 1]$
- A policy:  $\pi : S \rightarrow A$
- **Goal**: find the optimal policy



# Solving an MDP

- Find the optimal policy of an MDP

$$\pi^*(s) \quad \forall s$$

# Solving an MDP

- Find the optimal policy of an MDP

$$\pi^*(s) \quad \forall s$$

- Policies are evaluated using their expected utility:

$$EU(\pi) = \sum_{t=0}^{\infty} \gamma^t \sum_{s_{t+1}} P(s_{t+1} | s_t, \pi(s_t)) R(s_{t+1})$$

# Solving an MDP

- Find the optimal policy of an MDP

$$\pi^*(s) \quad \forall s$$

- Policies are evaluated using their expected utility:

$$EU(\pi) = \sum_{t=0}^{\infty} \gamma^t \sum_{s_{t+1}} P(s_{t+1} | s_t, \pi(s_t)) R(s_{t+1})$$

# Solving an MDP

- Find the optimal policy of an MDP

$$\pi^*(s) \quad \forall s$$

- Policies are evaluated using their expected utility:

$$EU(\pi) = \sum_{t=0}^{\infty} \gamma^t \sum_{s_{t+1}} P(s_{t+1} | s_t, \pi(s_t)) R(s_{t+1})$$

- The optimal policy is the one with highest expected utility:  $EU(\pi^*) \geq EU(\pi) \quad \forall \pi$



# Solving an MDP

- Three well-known techniques:
  1. Value iteration
  2. Policy Iteration
  3. Linear Programming

# Value Function

- $V(s_t)$ : The value of being in state  $s$  at time  $t$

# Value Function

- $V(s_t)$ : The value of being in state  $s$  at time  $t$

$V(s_t) :=$  expected sum of rewards of being in  $s$

# Finite horizon

- Assume that the process has  $T$  steps

# Finite horizon

- Assume that the process has  $T$  steps
- The value at step  $T$  is

# Finite horizon

- Assume that the process has  $T$  steps
- The value at step  $T$  is  $V(s_T) = R(s_T)$

# Finite horizon

- Assume that the process has  $T$  steps
- The value at step  $T$  is  $V(\mathbf{s}_T) = R(\mathbf{s}_T)$
- The value at step  $T-1$  is

$$V(\mathbf{s}_{T-1}) = \max_{\mathbf{a}_{T-1}} \left\{ R(\mathbf{s}_{T-1}) + \gamma \sum_{\mathbf{s}_T} P(\mathbf{s}_T \mid \mathbf{s}_{T-1}, \mathbf{a}_{T-1}) R(\mathbf{s}_T) \right\}$$

# Finite horizon

- Assume that the process has  $T$  steps
- The value at step  $T$  is  $V(\mathbf{s}_T) = R(\mathbf{s}_T)$

- The value at step  $T-1$  is

$$V(\mathbf{s}_{T-1}) = \max_{\mathbf{a}_{T-1}} \left\{ R(\mathbf{s}_{T-1}) + \gamma \sum_{\mathbf{s}_T} P(\mathbf{s}_T | \mathbf{s}_{T-1}, \mathbf{a}_{T-1}) R(\mathbf{s}_T) \right\}$$

- The value at step  $t$  is ( $0 \leq t \leq T$ )

$$V(\mathbf{s}_t) = \max_{\mathbf{a}_t} \left\{ R(\mathbf{s}_t) + \gamma \sum_{\mathbf{s}_{t+1}} P(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) V(\mathbf{s}_{t+1}) \right\}$$



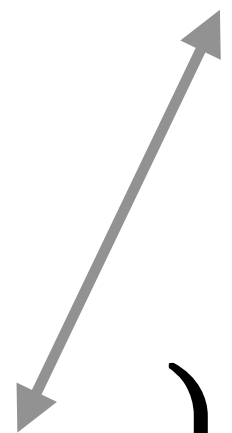
# Finite horizon

- Assume that the process has  $T$  steps
- The value at step  $T$  is  $V(\mathbf{s}_T) = R(\mathbf{s}_T)$

- The value at step  $T-1$  is

$$V(\mathbf{s}_{T-1}) = \max_{\mathbf{a}_{T-1}} \left\{ R(\mathbf{s}_{T-1}) + \gamma \sum_{\mathbf{s}_T} P(\mathbf{s}_T | \mathbf{s}_{T-1}, \mathbf{a}_{T-1}) R(\mathbf{s}_T) \right\}$$

- The value at step  $t$  is ( $0 \leq t \leq T$ )

$$V(\mathbf{s}_t) = \max_{\mathbf{a}_t} \left\{ R(\mathbf{s}_t) + \gamma \sum_{\mathbf{s}_{t+1}} P(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) V(\mathbf{s}_{t+1}) \right\}$$


# Value iteration (VI)

- Iteratively update  $V(s)$  for each state until convergence

# Value iteration (VI)

- Iteratively update  $V(s)$  for each state until convergence
- (Initialize  $V(s)$  for every state)

# Value iteration (VI)

- Iteratively update  $V(s)$  for each state until convergence
- (Initialize  $V(s)$  for every state)

- For  $i=1,2,3,\dots$

- For  $s=1,\dots,S$

$$V(s) = \max_a \left\{ R(s) + \gamma \sum_{s'} P(s' | s, a) V(s') \right\}$$

# Value iteration (VI)

- Iteratively update  $V(s)$  for each state until convergence
- (Initialize  $V(s)$  for every state)

- For  $i=1,2,3,\dots$

- For  $s=1,\dots,S$

$$V(s) = \max_a \left\{ R(s) + \gamma \sum_{s'} P(s' | s, a) V(s') \right\}$$

- The policy is implicit

- Once converged:  $\pi^*(s) = \arg \max_a \left\{ R(s) + \gamma \sum_{s'} P(s' | s, a) V^*(s') \right\} \forall s$

# Policy Iteration (PI)

- Improve policy explicitly.

# Policy Iteration (PI)

- Improve policy explicitly.

Start with any (e.g., random) policy  $\pi$

# Policy Iteration (PI)

- Improve policy explicitly.

Start with any (e.g., random) policy  $\pi$

Iterate until convergence:

1. Given current policy get the value of each state

$$\mathbf{V}^\pi(\mathbf{s}) = \mathbf{R}(\mathbf{s}) + \gamma \sum_{s'} \mathbf{P}(s' | \mathbf{s}, \pi(\mathbf{s})) \mathbf{V}^\pi(s') \quad \forall \mathbf{s}$$



# Policy Iteration (PI)

- Improve policy explicitly.

Start with any (e.g., random) policy  $\pi$

Iterate until convergence:

1. Given current policy get the value of each state

$$\mathbf{V}^\pi(\mathbf{s}) = \mathbf{R}(\mathbf{s}) + \gamma \sum_{s'} \mathbf{P}(s' | \mathbf{s}, \pi(\mathbf{s})) \mathbf{V}^\pi(s') \quad \forall \mathbf{s}$$

2. Update the current policy

$$\pi'(\mathbf{s}) = \arg \max_{\mathbf{a}} \left\{ \mathbf{R}(\mathbf{s}) + \gamma \sum_{s'} \mathbf{P}(s' | \mathbf{s}, \mathbf{a}) \mathbf{V}^\pi(s') \right\} \quad \forall \mathbf{s}$$

# Policy Iteration (PI)

- Improve policy explicitly.

Start with any (e.g., random) policy  $\pi$

Iterate until convergence:

1. Given current policy get the value of each state

$$\mathbf{V}^\pi(\mathbf{s}) = \mathbf{R}(\mathbf{s}) + \gamma \sum_{s'} \mathbf{P}(s' | \mathbf{s}, \pi(\mathbf{s})) \mathbf{V}^\pi(s') \quad \forall \mathbf{s}$$

2. Update the current policy

$$\pi'(\mathbf{s}) = \arg \max_{\mathbf{a}} \left\{ \mathbf{R}(\mathbf{s}) + \gamma \sum_{s'} \mathbf{P}(s' | \mathbf{s}, \mathbf{a}) \mathbf{V}^\pi(s') \right\} \quad \forall \mathbf{s}$$

Policy  
Evaluation

Policy  
Update

- Demo of the PI algorithm in a deterministic environment

[http://www.cs.toronto.edu/~lcharlin/courses/60629/reinforcejs/gridworld\\_dp.html](http://www.cs.toronto.edu/~lcharlin/courses/60629/reinforcejs/gridworld_dp.html)

