

# Reinforcement learning II

Tomáš Svoboda

Vision for Robots and Autonomous Systems, Center for Machine Perception  
Department of Cybernetics  
Faculty of Electrical Engineering, Czech Technical University in Prague

April 22, 2020

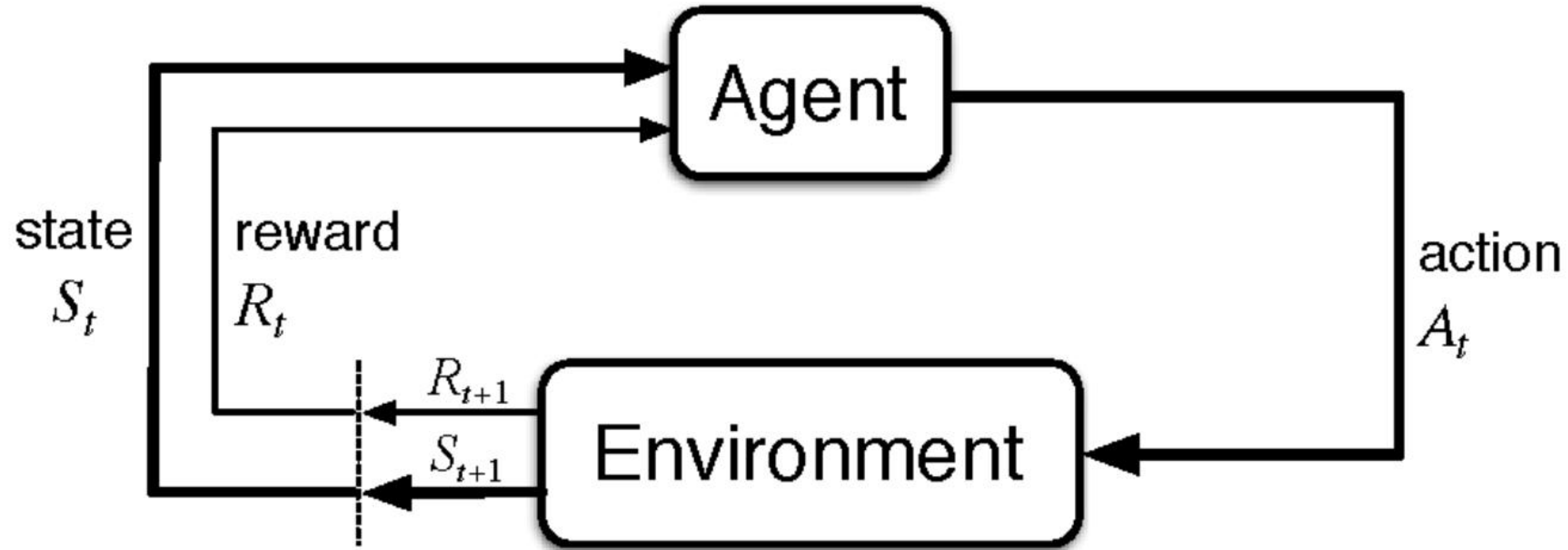
# Reinforcement learning II

Tomáš Svoboda

Vision for Robots and Autonomous Systems, Center for Machine Perception  
Department of Cybernetics  
Faculty of Electrical Engineering, Czech Technical University in Prague

April 22, 2020

## Recap: Reinforcement Learning



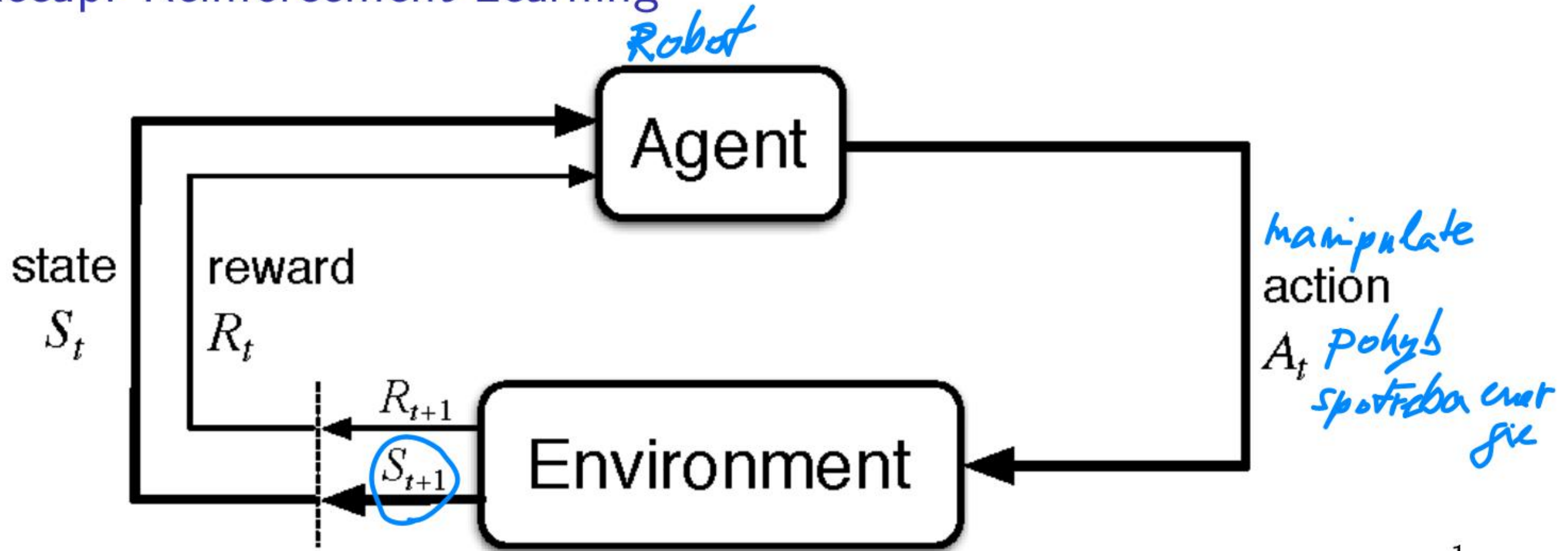
1

- ▶ Feedback in form of **Rewards**
- ▶ Learn to act so as to maximize sum of expected rewards.
- ▶ In `kuimaze` package, `env.step(action)` is the method.

---

<sup>1</sup>Scheme from [2]

# Recap: Reinforcement Learning



1

- ▶ Feedback in form of Rewards
- ▶ Learn to act so as to maximize sum of expected rewards.
- ▶ In `kuimaze` package, `env.step(action)` is the method.

<sup>1</sup>Scheme from [2]

# Learning to control flippers

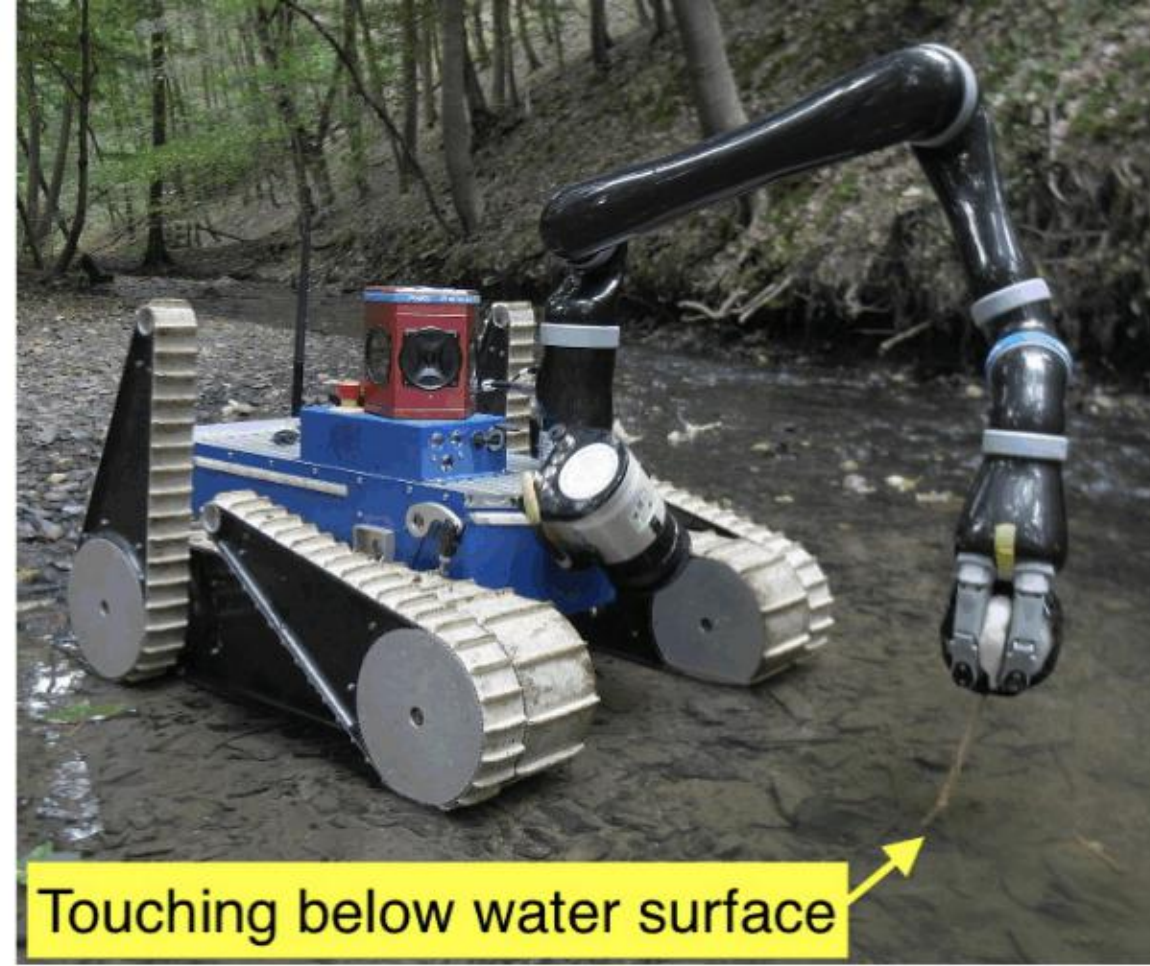
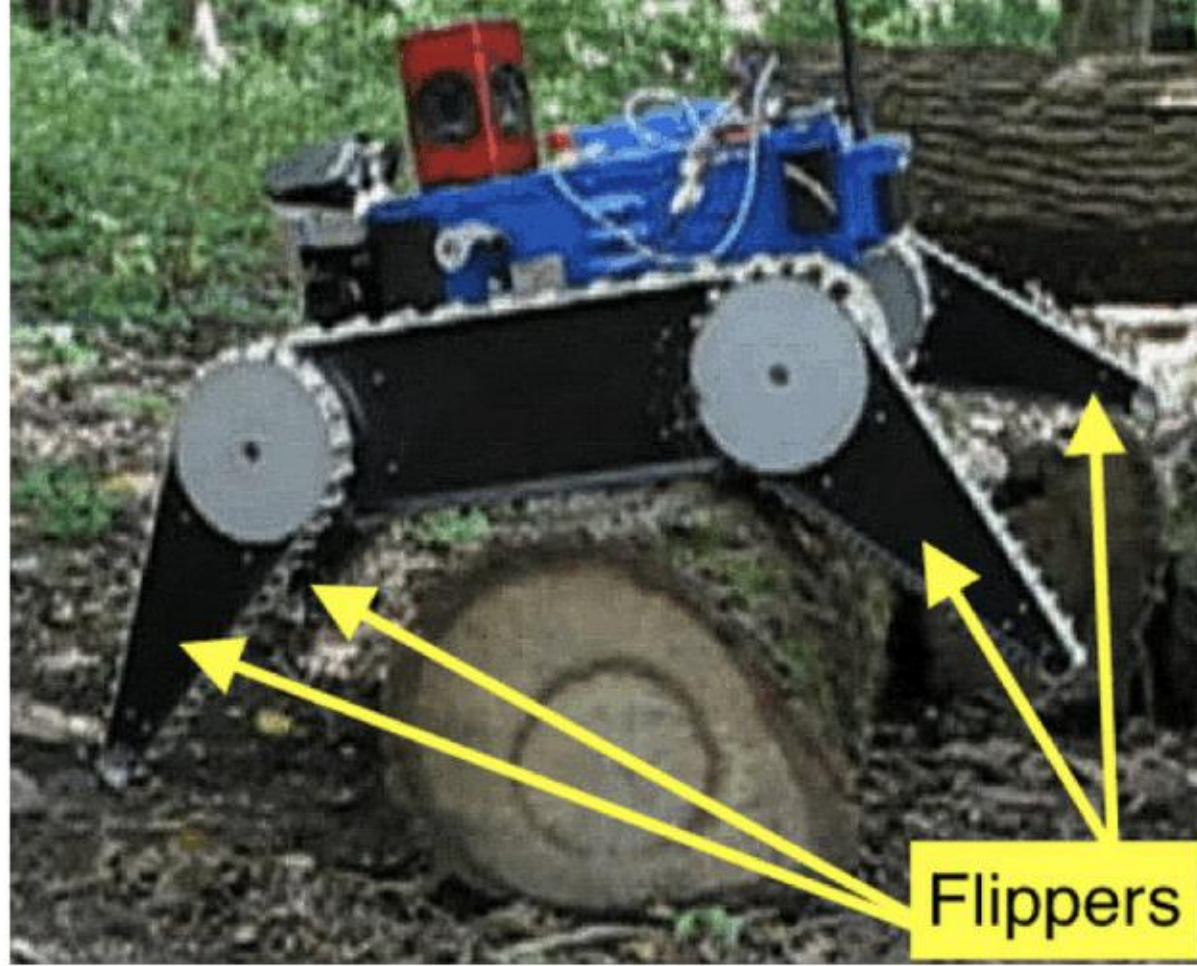


- ▶ What are the states?
- ▶ How to design rewards?
- ▶ How to perform training episodes (roll-outs)?
- ▶ Simulator to reality gap.

# Learning to control flippers

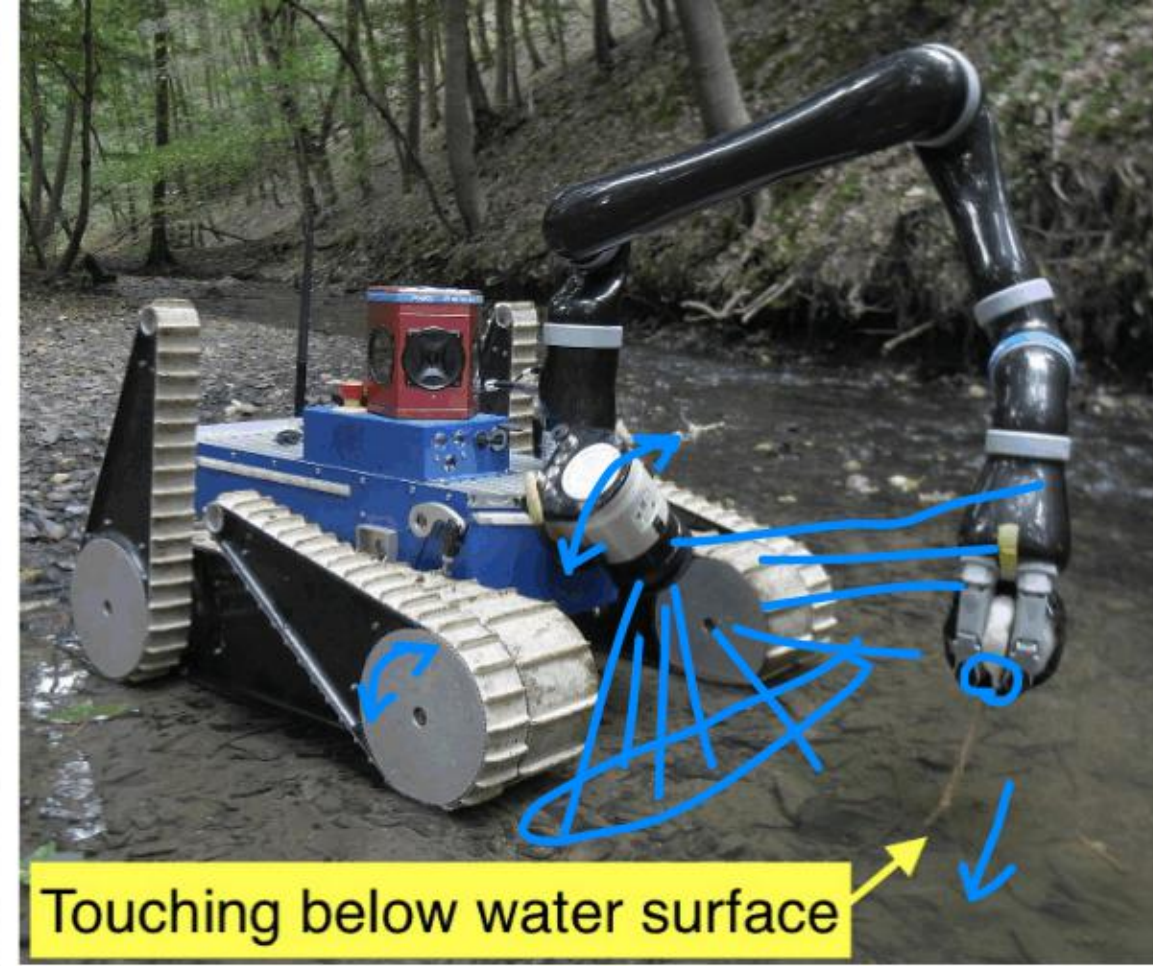
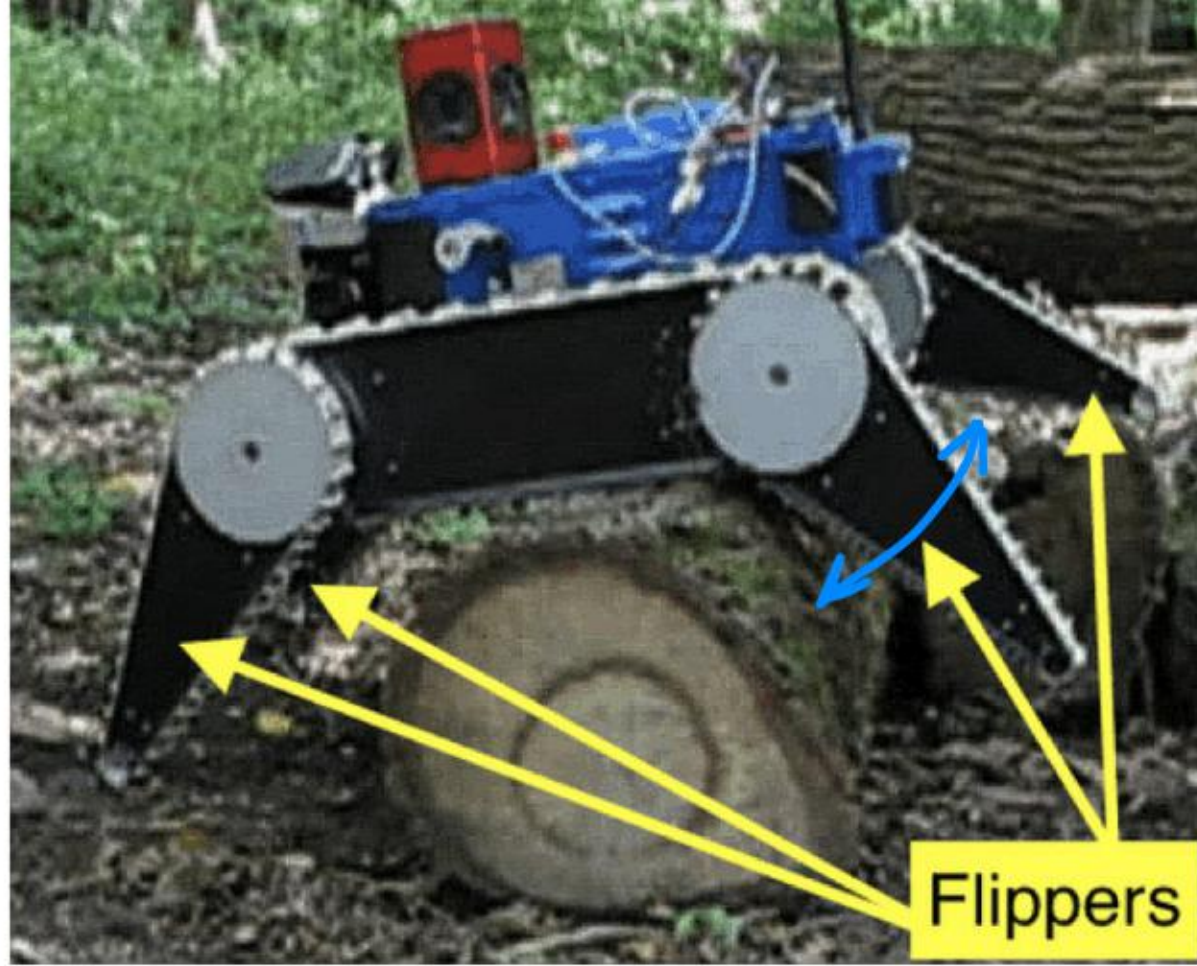


- ▶ What are the states?
- ▶ How to design rewards?
- ▶ How to perform training episodes (roll-outs)?
- ▶ Simulator to reality gap.



- ◆ **Construction:** 2× main tracks, 4× subtracks (flippers), differential break  
great stability and climbing capability
- ◆ **Sensor suite:** SICK LMS-151 range finder, Ladybug omnivision camera, Xsens MTi-G IMU  
3D sensing and localization
- ◆ **Control inputs:** Velocity vector, 4× flipper angle, 4× flipper stiffness,  
differential break (0/1)

difficult to control all of them manually!



- ◆ **Construction:** 2× main tracks, 4× subtracks (flippers), differential break  
great stability and climbing capability
- ◆ **Sensor suite:** SICK LMS-151 range finder, Ladybug omnivision camera, Xsens MTi-G IMU  
3D sensing and localization
- ◆ **Control inputs:** Velocity vector, 4× flipper angle, 4× flipper stiffness,  
differential break (0/1)

difficult to control all of them manually!



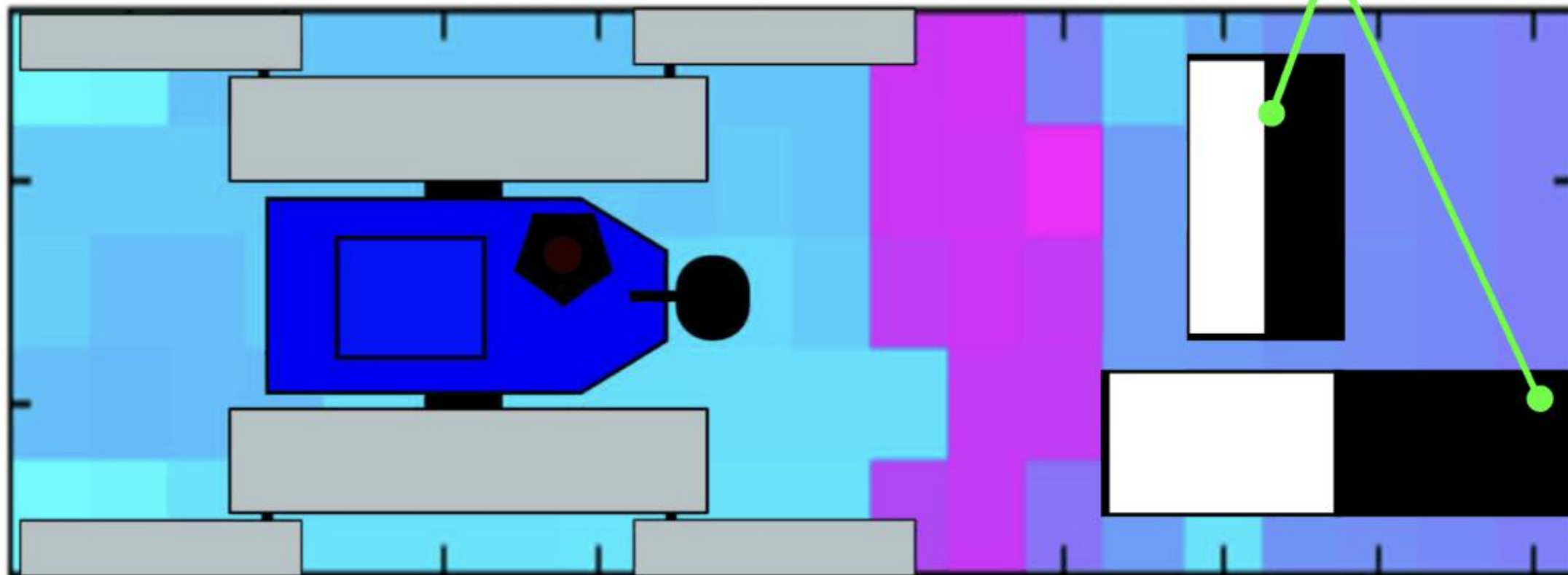
State  $s \in \mathcal{S} \subset \mathbb{R}^n$  concatenates:

◆ Proprioceptive measurements: roll, pitch, torques, velocity, acceleration.

◆ Local exteroceptive measurements.

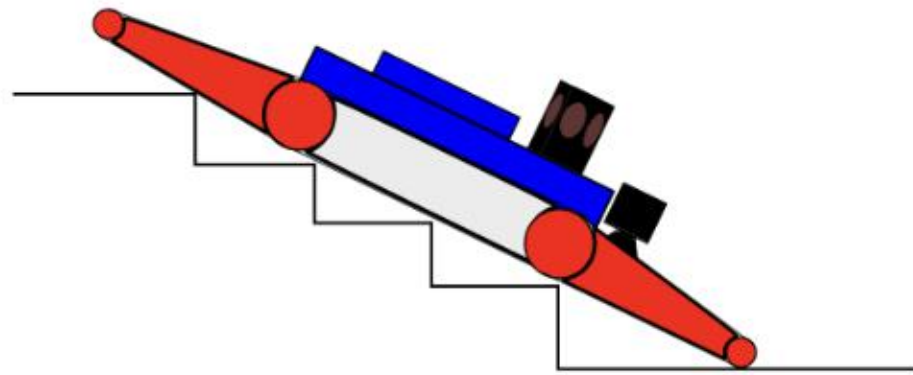
features on digital elevation map  
with fixed size.

Features

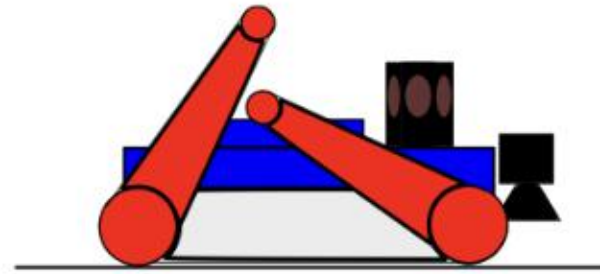


Instead of  $\mathbf{a} \in \mathcal{A} \subset \mathbb{R}^8$  we consider only 5 configurations<sup>2</sup>:

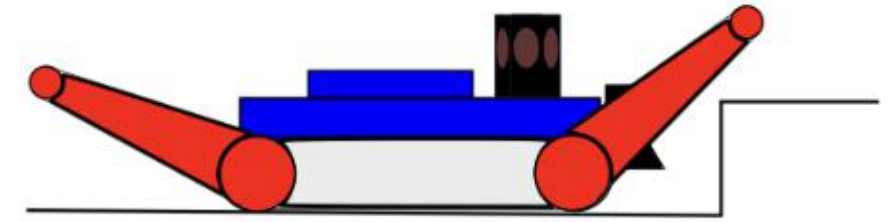
$$\mathcal{A} = \{\text{I-shape, V-shape, L-shape, U-shape soft, U-shape hard}\}$$



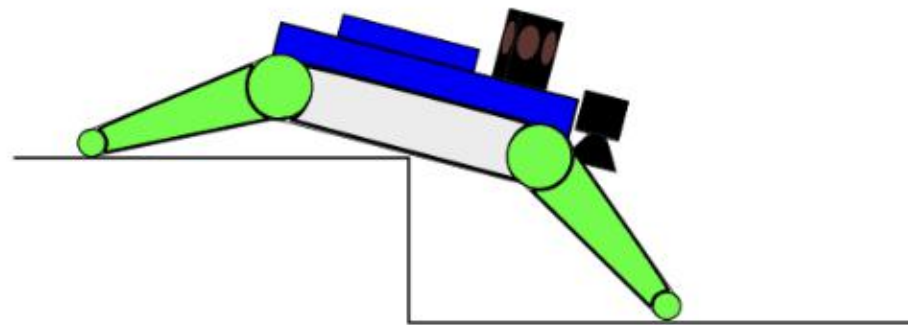
I-shape  
(Maximizes traction)



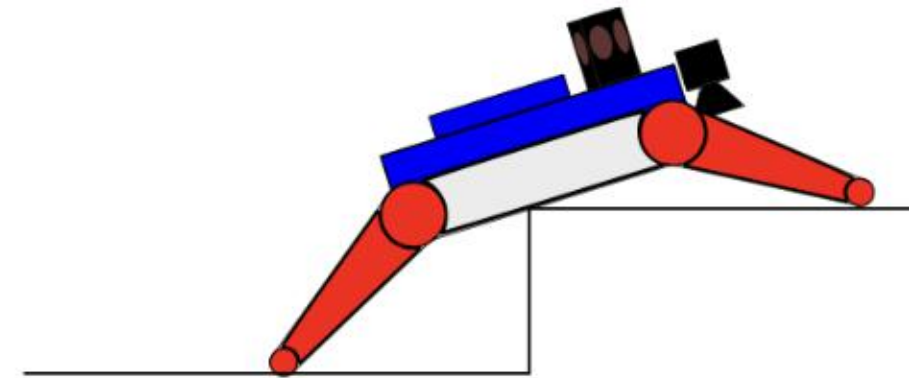
V-shape  
(Provides observability)



L-shape  
(Forward approach)



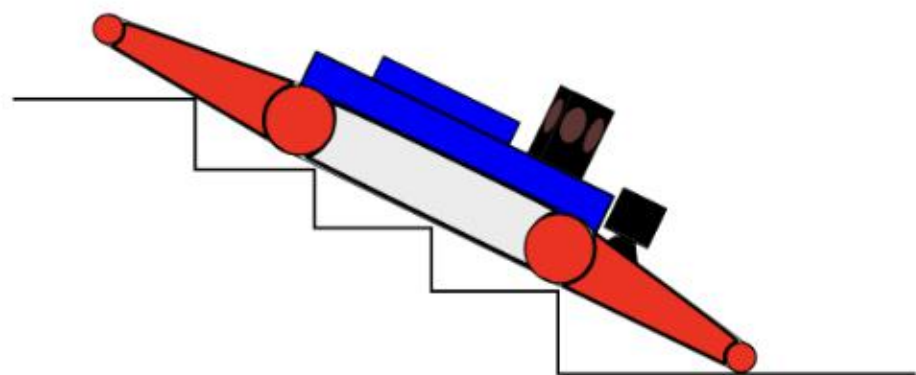
U-shape soft  
(Smooth climbing down)



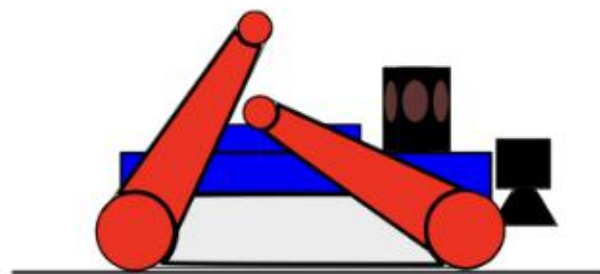
U-shape hard  
(Lifts the body up)

Instead of  $\underline{a} \in \mathcal{A} \subset \mathbb{R}^8$  we consider only 5 configurations<sup>2</sup>:

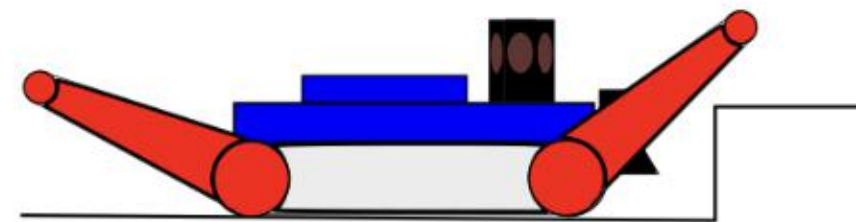
$$\mathcal{A} = \{\text{I-shape, V-shape, L-shape, U-shape soft, U-shape hard}\}$$



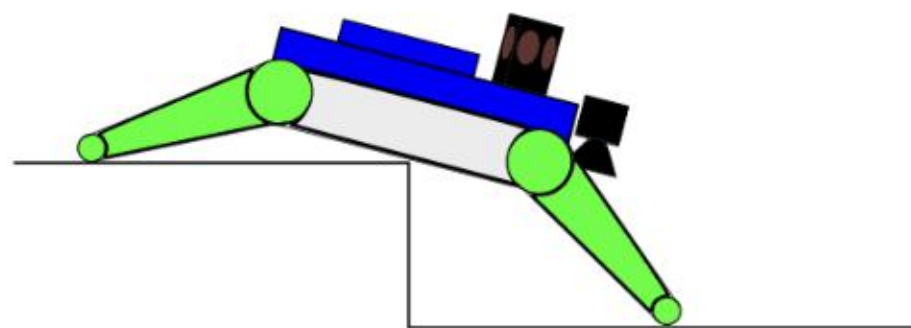
I-shape  
(Maximizes traction)



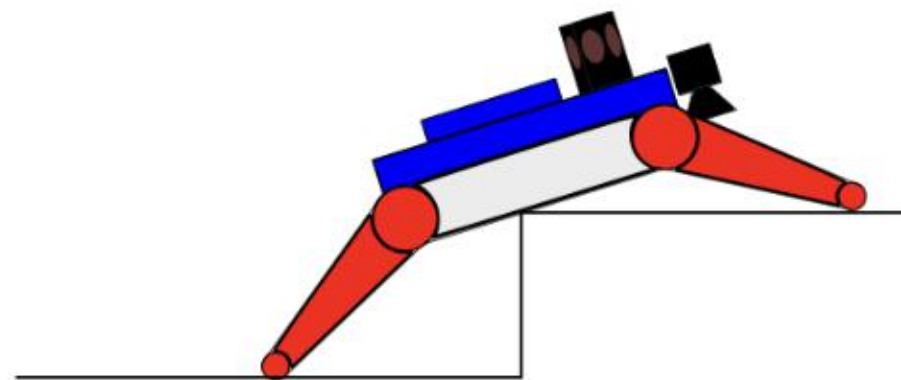
V-shape  
(Provides observability)



L-shape  
(Forward approach)



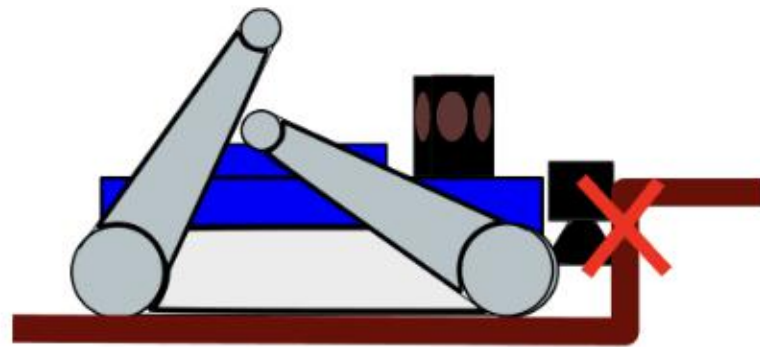
U-shape soft  
(Smooth climbing down)



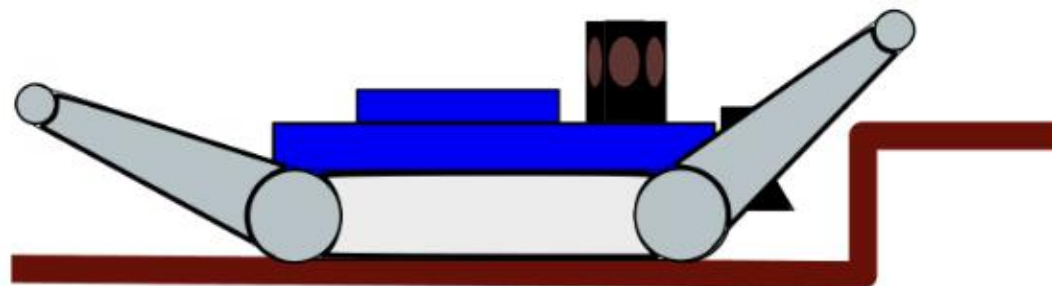
U-shape hard  
(Lifts the body up)

Reward  $r(a, s) : \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$  is a weighted sum of following contributions:

1. Safe pitch and roll reward, avoiding tipping over
2. Smoothness reward, suppresses body hits
3. Speed reward, drives robot forward
4. User denoted reward (penalty) indicating the success (failure) of the particular maneuver indicates failure/possible damages



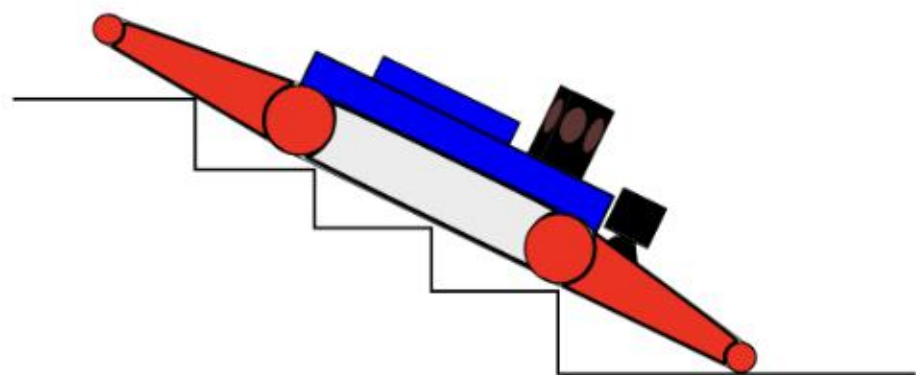
$$r(\text{V-shape}, s) = -1$$



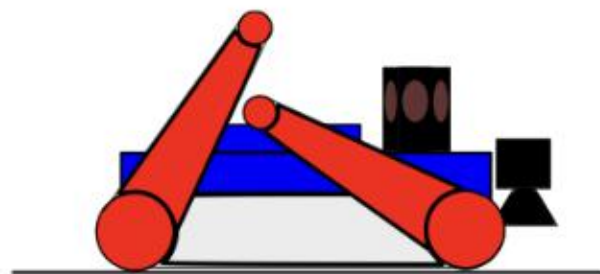
$$r(\text{L-shape}, s) = 1$$

Instead of  $\mathbf{a} \in \mathcal{A} \subset \mathbb{R}^8$  we consider only 5 configurations<sup>2</sup>:

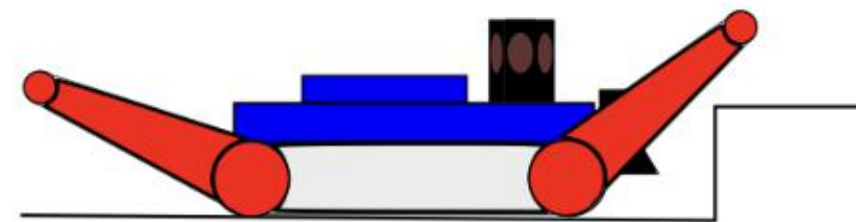
$$\mathcal{A} = \{\text{I-shape}, \text{V-shape}, \text{L-shape}, \text{U-shape soft}, \text{U-shape hard}\}$$



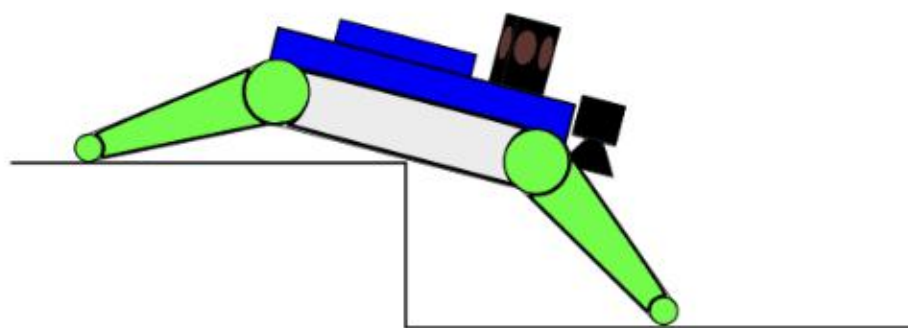
I-shape  
(Maximizes traction)



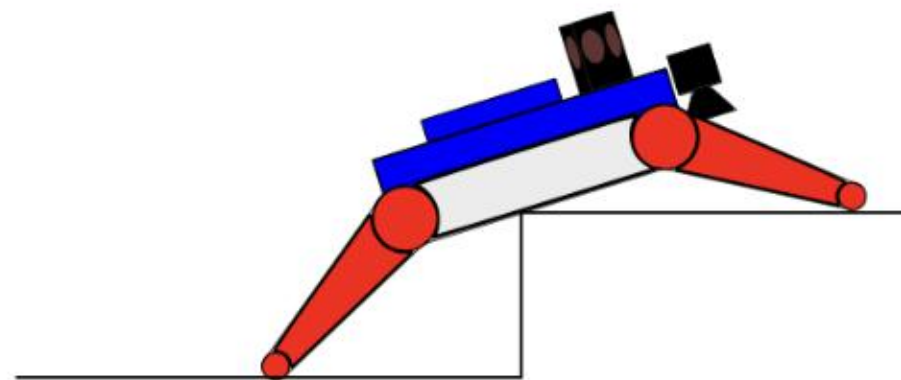
V-shape  
(Provides observability)



L-shape  
(Forward approach)



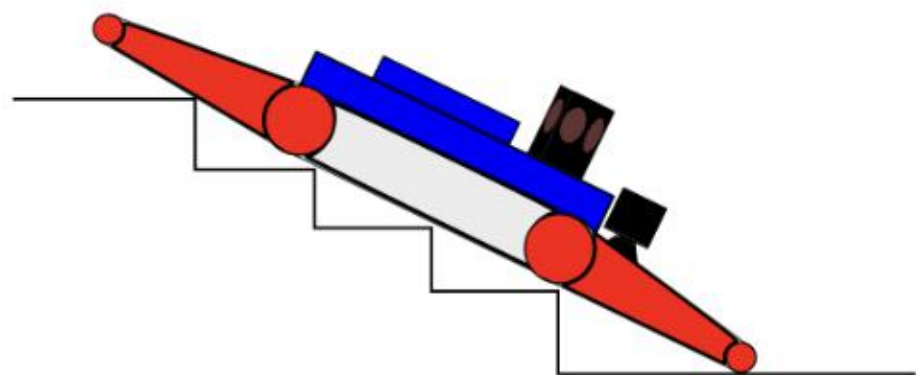
U-shape soft  
(Smooth climbing down)



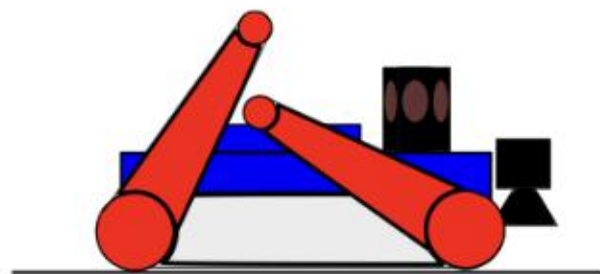
U-shape hard  
(Lifts the body up)

Instead of  $\underline{a} \in \mathcal{A} \subset \mathbb{R}^8$  we consider only 5 configurations<sup>2</sup>:

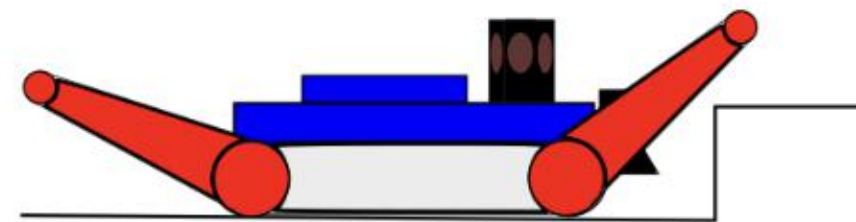
$$\mathcal{A} = \{\text{I-shape, V-shape, L-shape, U-shape soft, U-shape hard}\}$$



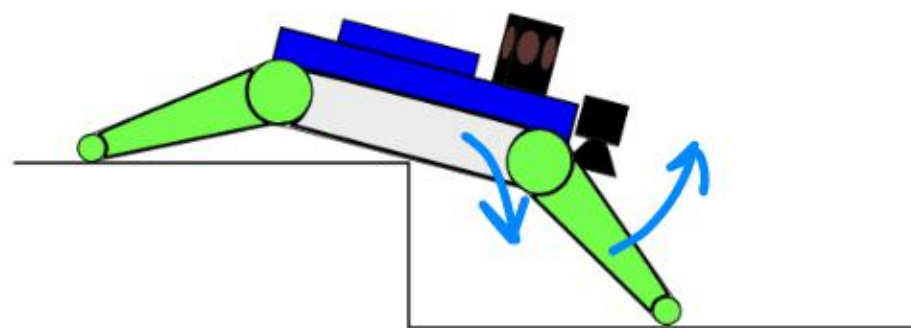
I-shape  
(Maximizes traction)



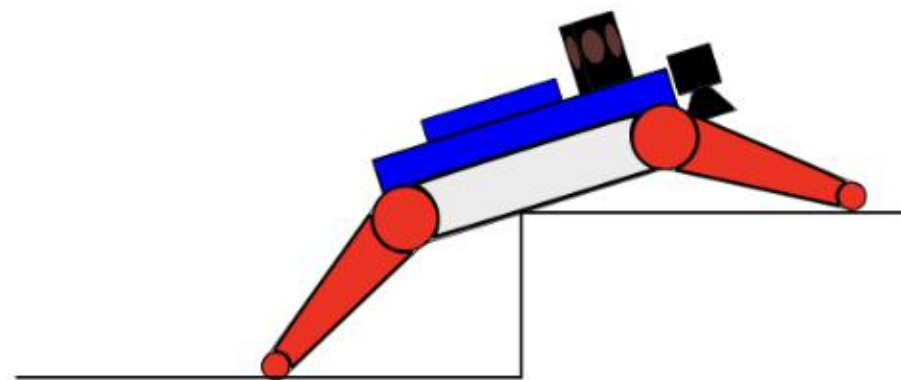
V-shape  
(Provides observability)



L-shape  
(Forward approach)



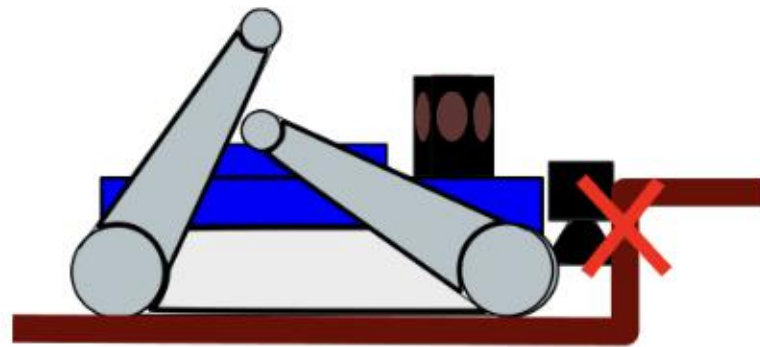
U-shape soft  
(Smooth climbing down)



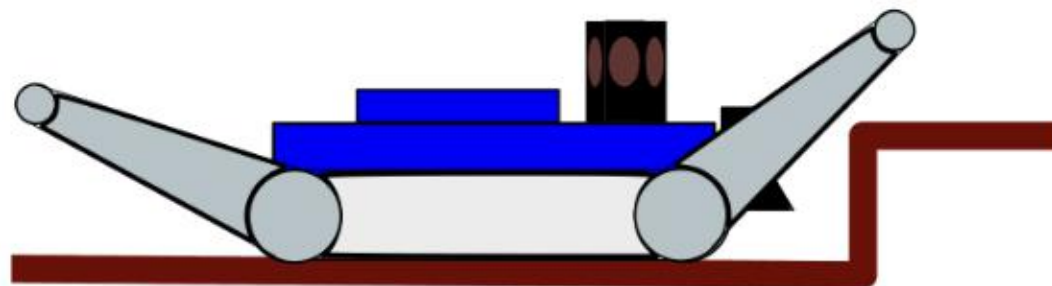
U-shape hard  
(Lifts the body up)

Reward  $r(a, s) : \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$  is a weighted sum of following contributions:

1. Safe pitch and roll reward, avoiding tipping over
2. Smoothness reward, suppresses body hits
3. Speed reward, drives robot forward
4. User denoted reward (penalty) indicating the success (failure) of the particular maneuver indicates failure/possible damages



$$r(\text{V-shape}, s) = -1$$

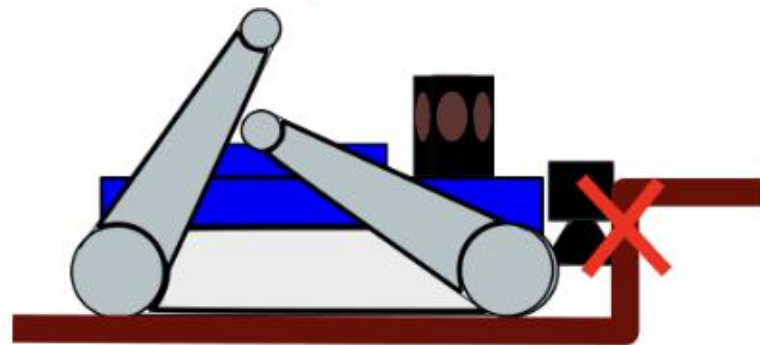


$$r(\text{L-shape}, s) = 1$$

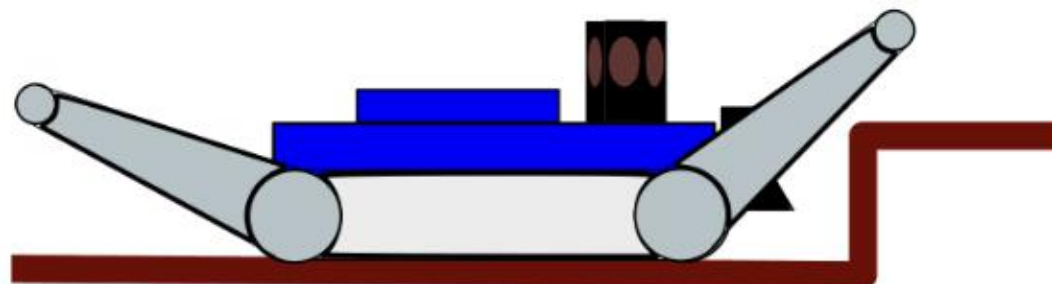
Reward  $r(a, s) : \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$  is a weighted sum of following contributions:

1. Safe pitch and roll reward, avoiding tipping over
2. Smoothness reward, suppresses body hits
3. Speed reward, drives robot forward
4. User denoted reward (penalty) indicating the success (failure) of the particular maneuver indicates failure/possible damages

*Reward shaping*



$$r(\text{V-shape}, s) = -1$$



$$r(\text{L-shape}, s) = 1$$



# From off-line (MDPs) to on-line (RL)

Markov decision process – MDPs. Off-line search, we know:

- ▶ A set of states  $s \in \mathcal{S}$  (map)
- ▶ A set of actions per state.  $a \in \mathcal{A}$
- ▶ A transition model  $p(s'|s, a)$  (robot)
- ▶ A reward function  $r(s, a, s')$  (map, robot)

Looking for the optimal policy  $\pi(s)$ . We can plan/search before the robot enters the environment.

## On-line problem:

- ▶ Transition  $p$  and reward  $r$  functions not known.
- ▶ Agent/robot must act and learn from experience.

# (Transition) Model-based learning

The main idea: Do something and:

- ▶ Learn an approximate model from experiences.
- ▶ Solve as if the model were correct.

Learning MDP model:

- ▶ Try  $s, a$ , observe  $s'$ , count  $s, a, s'$ .
- ▶ Normalize to get an estimate of  $p(s'|s, a)$
- ▶ Discover each  $r(s, a, s')$  when experienced.

Solve the learned MDP.

# (Transition) Model-based learning

The main idea: Do something and:

- ▶ Learn an approximate model from experiences.
- ▶ Solve as if the model were correct.

Learning MDP model:

- ▶ Try  $s, a$ , observe  $s'$ , count  $s, a, s'$ .
- ▶ Normalize to get an estimate of  $p(s'|s, a)$
- ▶ Discover each  $r(s, a, s')$  when experienced.

Solve the learned MDP.

# Model-free learning

- ▶  $r, p$  not known.
- ▶ Move around, observe
- ▶ And learn on the way.
- ▶ **Goal:** learn the state value  $v(s)$  or (better) q-value  $q(s, a)$  functions.

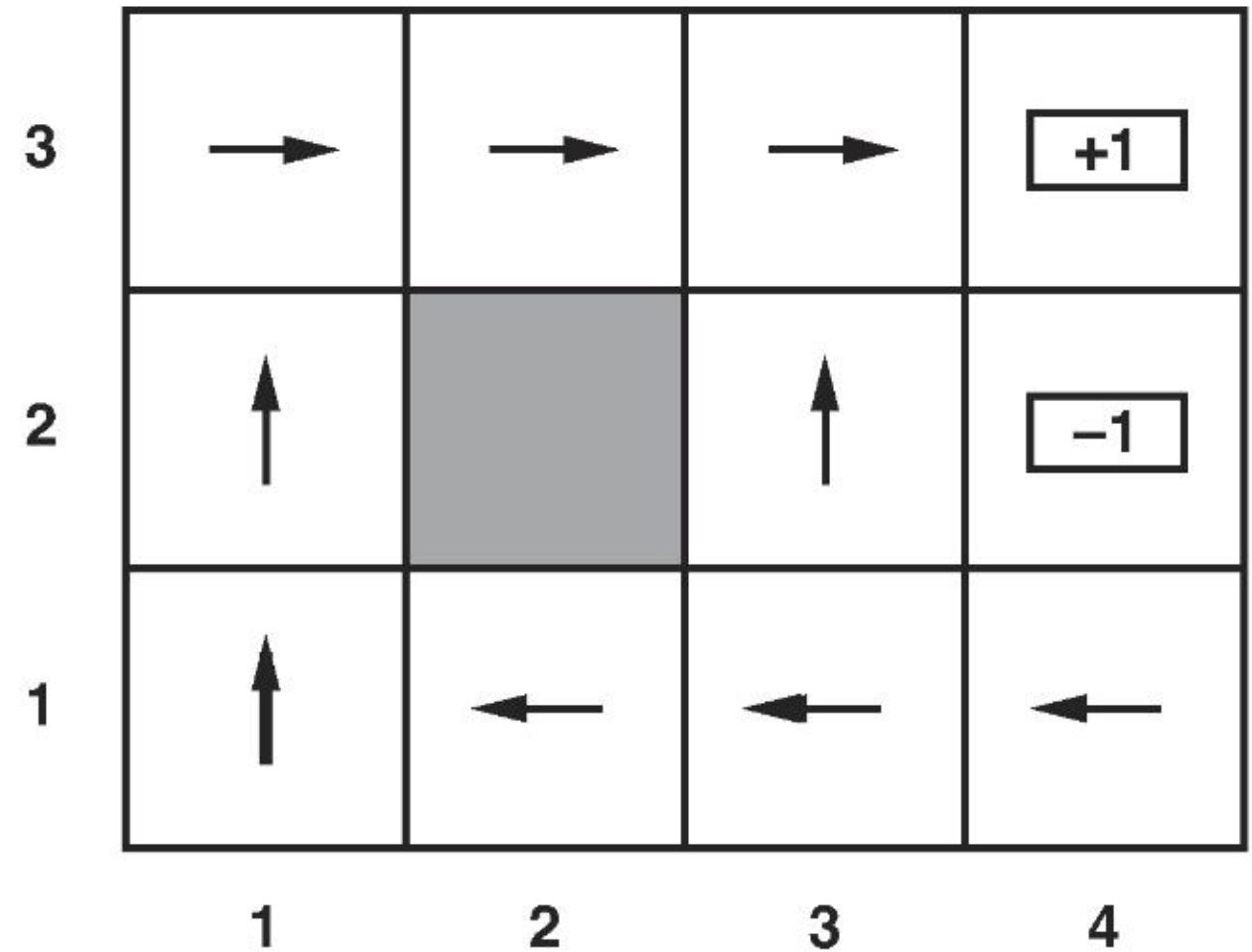


Image from [1]

# Model-free learning

- ▶  $r, p$  not known.
- ▶ Move around, observe
- ▶ And learn on the way.
- ▶ **Goal:** learn the state value  $v(s)$  or (better) q-value  $q(s, a)$  functions.

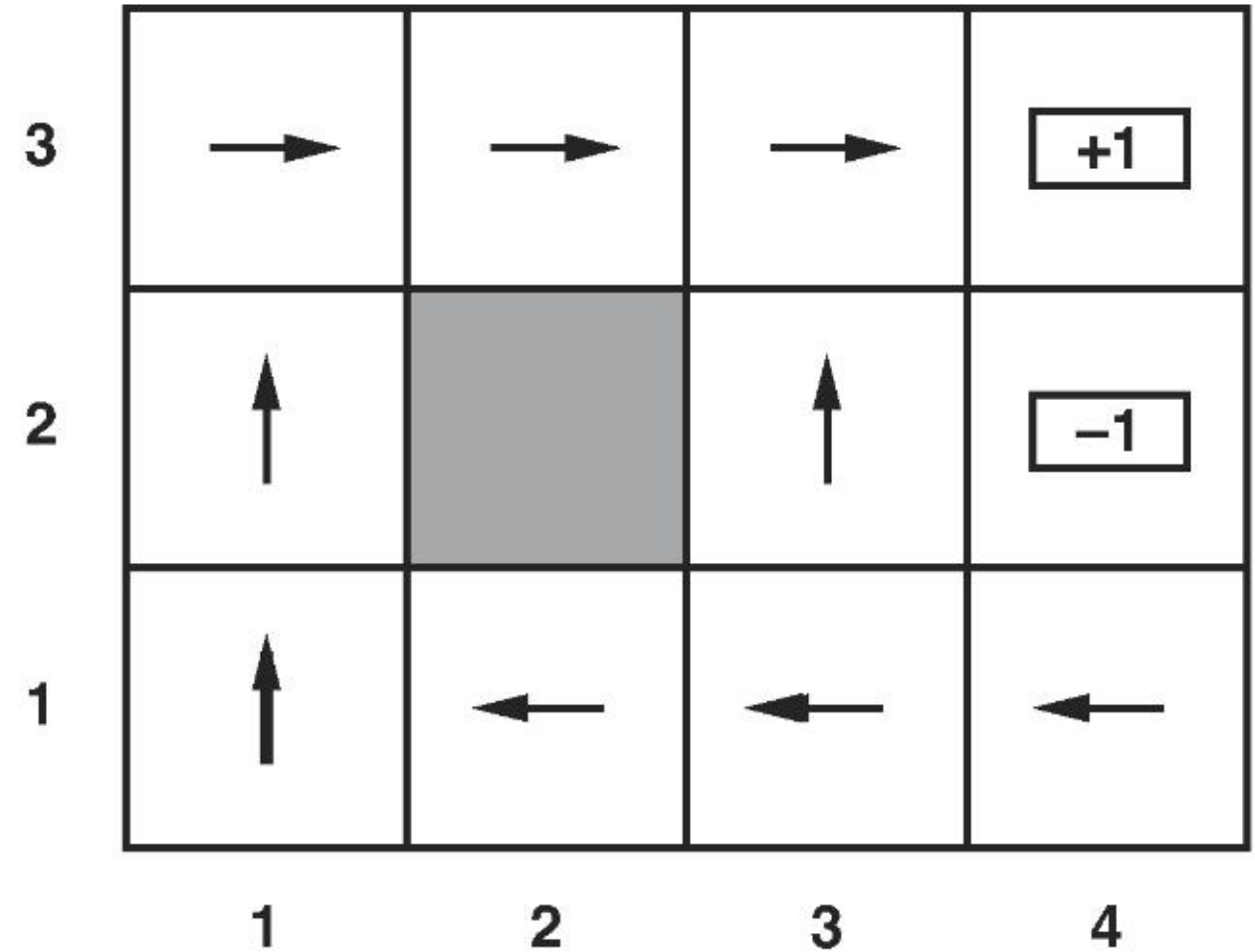
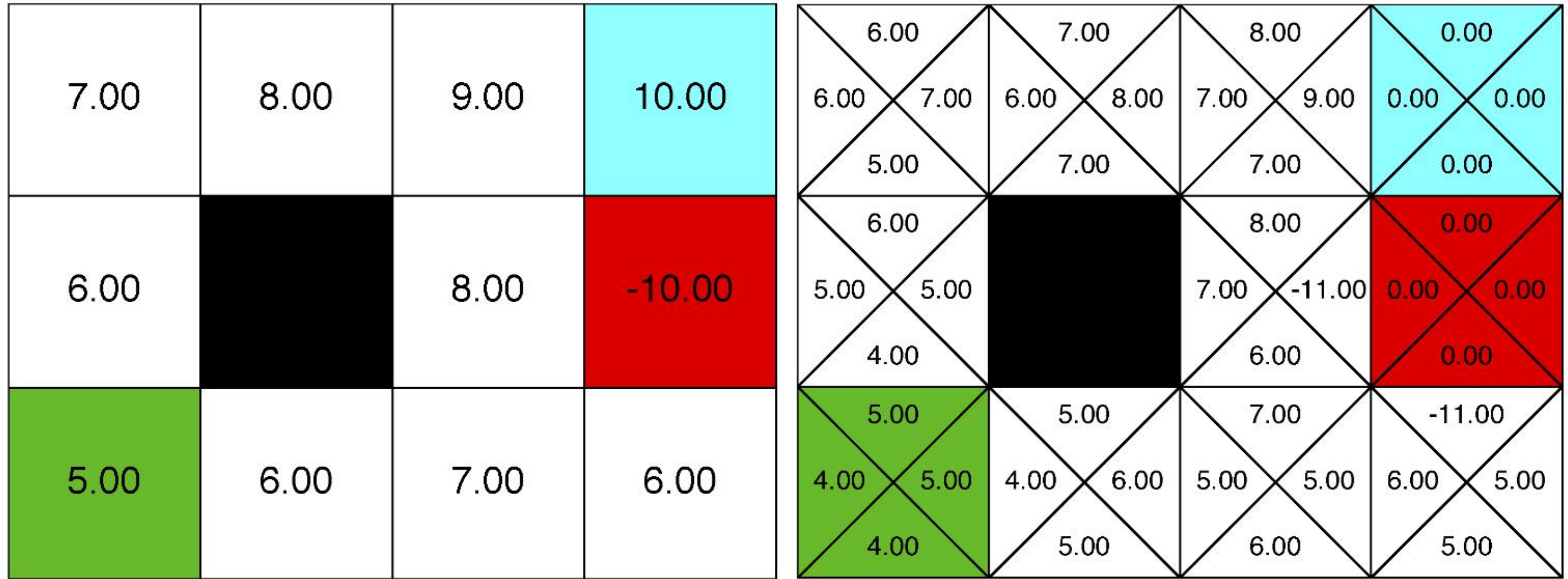


Image from [1]

# Recap: $V$ – and $Q$ – values, converged ...

$\gamma = 1$ , rewards  $-1, +10, -10$ , and no confusion - deterministic robot



$$V(S_t) = R_{t+1} + V(S_{t+1})$$

$$Q(S_t, A_t) = R_{t+1} + \max_a Q(S_{t+1}, a)$$

# Model-free TD learning, updating after each transition

- ▶ Observe, experience environment through learning episodes, collecting:

$$S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}, R_{t+2}, \dots$$

- ▶ Update by mimicking Bellman updates after each transition  $(S_t, A_t, R_{t+1}, S_{t+1})$

# Model-free TD learning, updating after each transition

- ▶ Observe, experience environment through learning episodes, collecting:

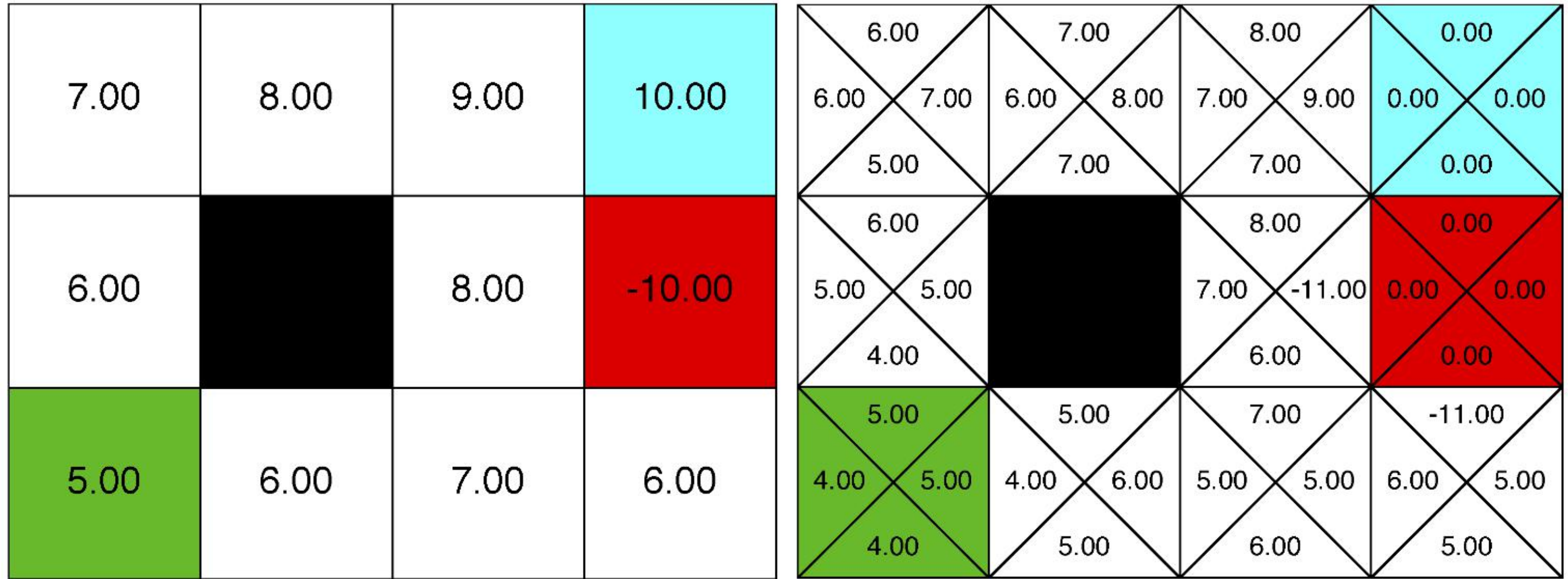
$$S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}, R_{t+2}, \dots$$

- ▶ Update by mimicking Bellman updates after each transition  $(S_t, A_t, R_{t+1}, S_{t+1})$



# Recap: $V$ – and $Q$ – values, converged ...

$\gamma = 1$ , rewards  $-1, +10, -10$ , and no confusion - deterministic robot



$$V(S_t) = R_{t+1} + V(S_{t+1})$$

$$Q(S_t, A_t) = R_{t+1} + \max_a Q(S_{t+1}, a)$$

# Model-free TD learning, updating after each transition

- ▶ Observe, experience environment through learning episodes, collecting:

$$S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}, R_{t+2}, \dots$$

- ▶ Update by mimicking Bellman updates after each transition  $(S_t, A_t, R_{t+1}, S_{t+1})$

# Model-free TD learning, updating after each transition

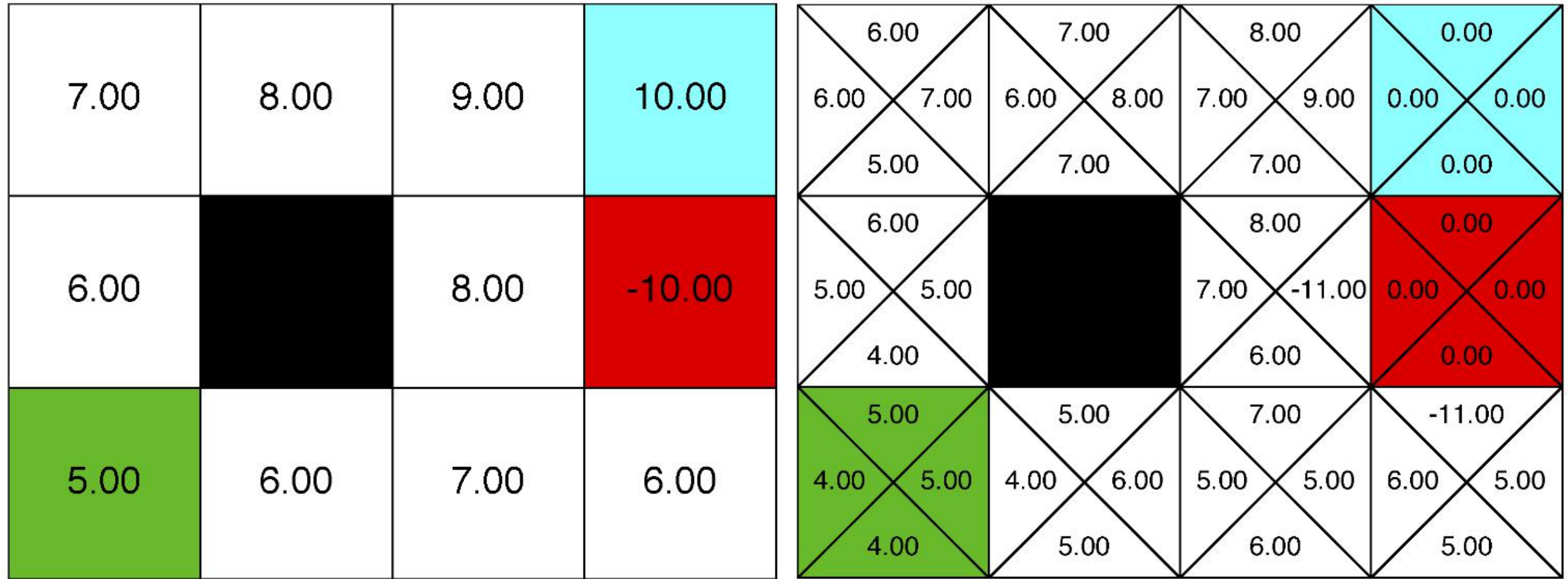
- ▶ Observe, experience environment through learning episodes, collecting:

$$S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}, R_{t+2}, \dots$$

- ▶ Update by mimicking Bellman updates after each transition  $(S_t, A_t, R_{t+1}, S_{t+1})$

# Recap: $V$ – and $Q$ – values, converged ...

$\gamma = 1$ , rewards  $-1, +10, -10$ , and no confusion - deterministic robot



$$V(S_t) = R_{t+1} + V(S_{t+1})$$

$$Q(S_t, A_t) = R_{t+1} + \max_a Q(S_{t+1}, a)$$

# Model-free TD learning, updating after each transition

- ▶ Observe, experience environment through learning episodes, collecting:

$$S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}, R_{t+2}, \dots$$

- ▶ Update by mimicking Bellman updates after each transition  $(S_t, A_t, R_{t+1}, S_{t+1})$

# Model-free TD learning, updating after each transition

*temporal difference*

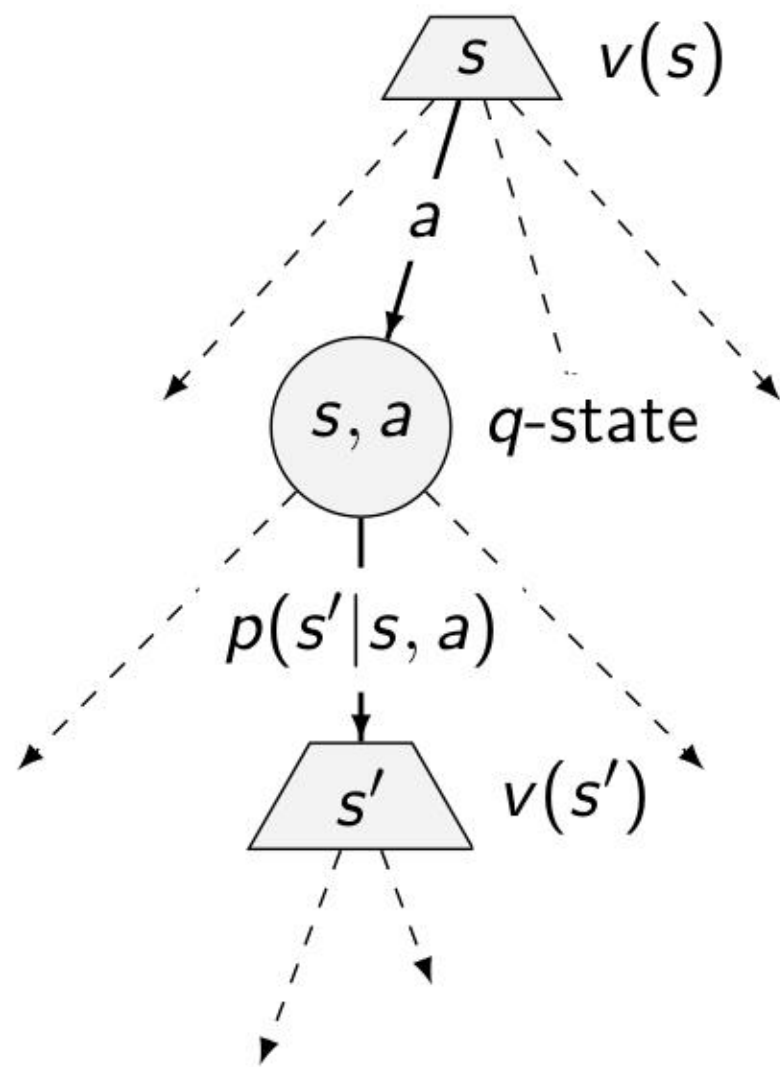
- ▶ Observe, experience environment through learning episodes, collecting:

$$S_t, A_t, \underbrace{R_{t+1}, S_{t+1}}_{\text{transition}}, A_{t+1}, R_{t+2}, \dots$$

- ▶ Update by mimicking Bellman updates after each transition  $(S_t, A_t, R_{t+1}, S_{t+1})$

## Recap: Bellman optimality equations for $v(s)$ and $q(s, a)$

$$\begin{aligned}v(s) &= \max_a \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma v(s')] \\ &= \max_a q(s, a)\end{aligned}$$



The value of a q-state  $(s, a)$ :

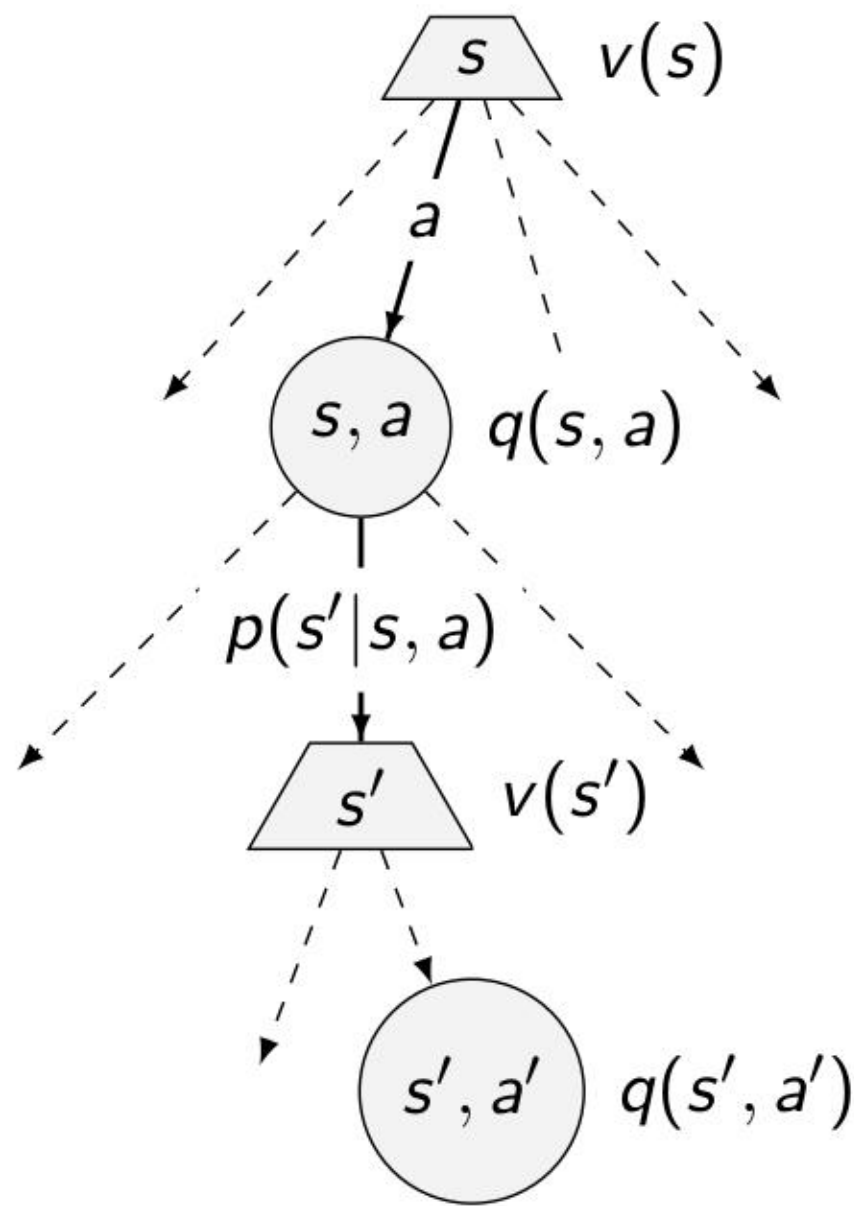
$$\begin{aligned}q(s, a) &= \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma v(s')] \\ &= \sum_{s'} p(s'|s, a) \left[ r(s, a, s') + \gamma \max_{a'} q(s', a') \right]\end{aligned}$$

## Recap: Bellman optimality equations for $v(s)$ and $q(s, a)$

$$\begin{aligned}v(s) &= \max_a \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma v(s')] \\ &= \max_a q(s, a)\end{aligned}$$

The value of a  $q$ -state  $(s, a)$ :

$$\begin{aligned}q(s, a) &= \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma v(s')] \\ &= \sum_{s'} p(s'|s, a) \left[ r(s, a, s') + \gamma \max_{a'} q(s', a') \right]\end{aligned}$$



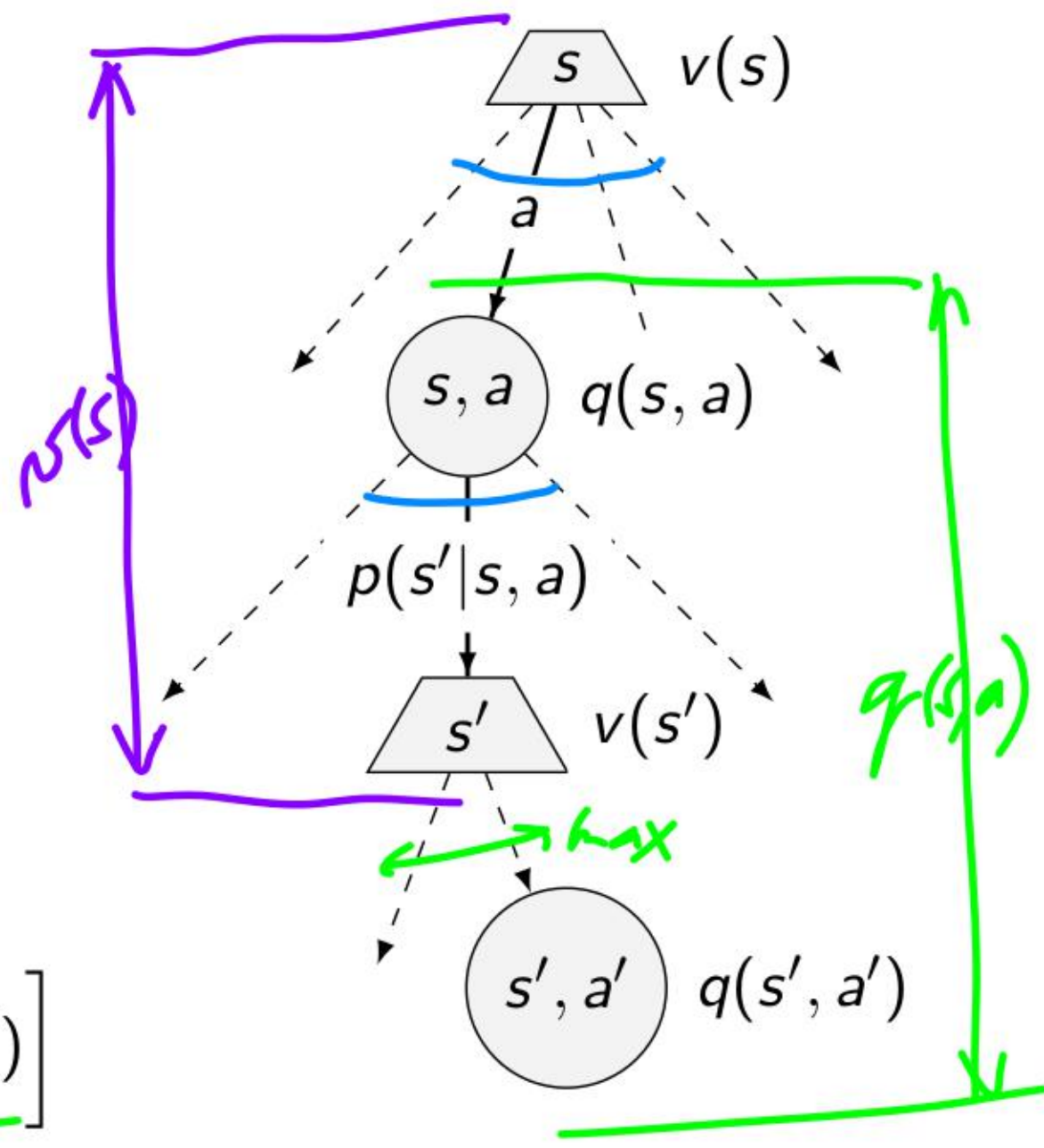


# Recap: Bellman optimality equations for $v(s)$ and $q(s, a)$

$$\begin{aligned}v(s) &= \max_a \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma v(s')] \\ &= \max_a q(s, a)\end{aligned}$$

The value of a  $q$ -state  $(s, a)$ :

$$\begin{aligned}q(s, a) &= \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma v(s')] \\ &= \sum_{s'} p(s'|s, a) \left[ r(s, a, s') + \gamma \max_{a'} q(s', a') \right]\end{aligned}$$



# Q-learning

Learn  $Q$  values as the robot/agent goes (temporal difference). If some  $Q$  quantity not known, initialize.

- ▶ time  $t$ , at  $S_t$
- ▶ take  $A_t \in \mathcal{A}(S_t)$ , observe  $R_{t+1}, S_{t+1}$
- ▶ compute trial/sample estimate at time  $t$   
trial =  $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶  $\alpha$  temporal difference update  
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{trial} - Q(S_t, A_t))$
- ▶  $S_t \leftarrow S_{t+1}$  and repeat (unless  $S_t$  is terminal)

In each step  $Q$  approximates the optimal  $q^*$  function.

# Q-learning

Learn  $Q$  values as the robot/agent goes (temporal difference). If some  $Q$  quantity not known, initialize.

- ▶ time  $t$ , at  $S_t$
- ▶ take  $A_t \in \mathcal{A}(S_t)$ , observe  $R_{t+1}, S_{t+1}$
- ▶ compute trial/sample estimate at time  $t$   
trial =  $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶  $\alpha$  temporal difference update  
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{trial} - Q(S_t, A_t))$
- ▶  $S_t \leftarrow S_{t+1}$  and repeat (unless  $S_t$  is terminal)

In each step  $Q$  approximates the optimal  $q^*$  function.

# Q-learning

Learn  $Q$  values as the robot/agent goes (temporal difference). If some  $Q$  quantity not known, initialize.

- ▶ time  $t$ , at  $S_t$
- ▶ take  $A_t \in \mathcal{A}(S_t)$ , observe  $R_{t+1}, S_{t+1}$
- ▶ compute trial/sample estimate at time  $t$   
trial =  $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶  $\alpha$  temporal difference update  
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{trial} - Q(S_t, A_t))$
- ▶  $S_t \leftarrow S_{t+1}$  and repeat (unless  $S_t$  is terminal)

In each step  $Q$  approximates the optimal  $q^*$  function.

# Q-learning

Learn  $Q$  values as the robot/agent goes (temporal difference). If some  $Q$  quantity not known, initialize.

- ▶ time  $t$ , at  $S_t$
- ▶ take  $A_t \in \mathcal{A}(S_t)$ , observe  $R_{t+1}, S_{t+1}$
- ▶ compute trial/sample estimate at time  $t$   
trial =  $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶  $\alpha$  temporal difference update  
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{trial} - Q(S_t, A_t))$
- ▶  $S_t \leftarrow S_{t+1}$  and repeat (unless  $S_t$  is terminal)

In each step  $Q$  approximates the optimal  $q^*$  function.

# Q-learning

Learn  $Q$  values as the robot/agent goes (temporal difference). If some  $Q$  quantity not known, initialize.

- ▶ time  $t$ , at  $S_t$
- ▶ take  $A_t \in \mathcal{A}(S_t)$ , observe  $R_{t+1}, S_{t+1}$
- ▶ compute trial/sample estimate at time  $t$   
trial =  $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶  $\alpha$  temporal difference update  
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{trial} - Q(S_t, A_t))$
- ▶  $S_t \leftarrow S_{t+1}$  and repeat (unless  $S_t$  is terminal)

In each step  $Q$  approximates the optimal  $q^*$  function.

# Q-learning

Learn  $Q$  values as the robot/agent goes (temporal difference). If some  $Q$  quantity not known, initialize.

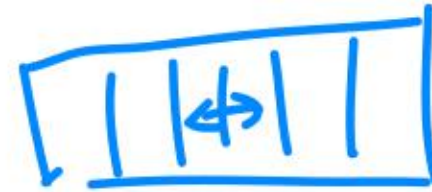
- ▶ time  $t$ , at  $S_t$
- ▶ take  $A_t \in \mathcal{A}(S_t)$ , observe  $R_{t+1}, S_{t+1}$
- ▶ compute trial/sample estimate at time  $t$   
trial =  $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶  $\alpha$  temporal difference update  
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{trial} - Q(S_t, A_t))$
- ▶  $S_t \leftarrow S_{t+1}$  and repeat (unless  $S_t$  is terminal)

In each step  $Q$  approximates the optimal  $q^*$  function.

# Q-learning

Learn  $Q$  values as the robot/agent goes (temporal difference). If some  $Q$  quantity not known, initialize.

- ▶ time  $t$ , at  $S_t$
- ▶ take  $A_t \in \mathcal{A}(S_t)$ , observe  $R_{t+1}, S_{t+1}$
- ▶ compute trial/sample estimate at time  $t$   
trial =  $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶  $\alpha$  temporal difference update  
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{trial} - Q(S_t, A_t))$
- ▶  $S_t \leftarrow S_{t+1}$  and repeat (unless  $S_t$  is terminal)



*max steps*

In each step  $Q$  approximates the optimal  $q^*$  function.



# Q-learning

Learn  $Q$  values as the robot/agent goes (temporal difference). If some  $Q$  quantity not known, initialize.

- ▶ time  $t$ , at  $S_t$
- ▶ take  $A_t \in \mathcal{A}(S_t)$ , observe  $R_{t+1}, S_{t+1}$
- ▶ compute trial/sample estimate at time  $t$   
trial =  $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶  $\alpha$  temporal difference update  
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{trial} - Q(S_t, A_t))$
- ▶  $S_t \leftarrow S_{t+1}$  and repeat (unless  $S_t$  is terminal)

In each step  $Q$  approximates the optimal  $q^*$  function.

# Q-learning: algorithm

step size  $0 < \alpha \leq 1$

initialize  $Q(s, a)$  for all  $s \in \mathcal{S}, a \in \mathcal{S}(s)$

**repeat** episodes:

    initialize  $S$

**for** for each step of episode: **do**

        choose  $A$  from  $S$

        take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$

**end for** until  $S$  is terminal

**until** Time is up, ...

## How to select $A_t$ in $S_t$ ?

- ▶ time  $t$ , at  $S_t$
- ▶ take  $A_t \in \mathcal{A}(S_t)$ , observe  $R_{t+1}, S_{t+1}$
- ▶ compute trial/sample estimate at time  $t$   
trial =  $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶  $\alpha$  temporal difference update  
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{trial} - Q(S_t, A_t))$
- ▶  $S_t \leftarrow S_{t+1}$  and repeat (unless  $S_t$  is terminal)

## How to select $A_t$ in $S_t$ ?

- ▶ time  $t$ , at  $S_t$
- ▶ take  $A_t \in \mathcal{A}(S_t)$ , observe  $R_{t+1}, S_{t+1}$
- ▶ compute trial/sample estimate at time  $t$   
trial =  $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶  $\alpha$  temporal difference update  
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{trial} - Q(S_t, A_t))$
- ▶  $S_t \leftarrow S_{t+1}$  and repeat (unless  $S_t$  is terminal)

## How to select $A_t$ in $S_t$ ?

- ▶ time  $t$ , at  $S_t$
- ▶ take  $A_t$  derived from  $Q$ , observe  $R_{t+1}, S_{t+1}$
- ▶ compute trial/sample estimate at time  $t$   
trial =  $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶  $\alpha$  temporal difference update  
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{trial} - Q(S_t, A_t))$
- ▶  $S_t \leftarrow S_{t+1}$  and repeat (unless  $S_t$  is terminal)

## How to select $A_t$ in $S_t$ ?

- ▶ time  $t$ , at  $S_t$
- ▶ take  $A_t \in \mathcal{A}(S_t)$ , observe  $R_{t+1}, S_{t+1}$
- ▶ compute trial/sample estimate at time  $t$   
trial =  $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶  $\alpha$  temporal difference update  
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{trial} - Q(S_t, A_t))$
- ▶  $S_t \leftarrow S_{t+1}$  and repeat (unless  $S_t$  is terminal)

# How to select $A_t$ in $S_t$ ?

- ▶ time  $t$ , at  $S_t$
- ▶ take  $A_t \in \mathcal{A}(S_t)$ , observe  $R_{t+1}, S_{t+1}$
- ▶ compute trial/sample estimate at time  $t$   
trial =  $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶  $\alpha$  temporal difference update  
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{trial} - Q(S_t, A_t))$
- ▶  $S_t \leftarrow S_{t+1}$  and repeat (unless  $S_t$  is terminal)

random	█	22%
$A_t = \operatorname{argmax}_a Q(S_t, a)$	█	37%
nějaká kombinace	█	40%

## How to select $A_t$ in $S_t$ ?

- ▶ time  $t$ , at  $S_t$
- ▶ take  $A_t$  derived from  $Q$ , observe  $R_{t+1}, S_{t+1}$
- ▶ compute trial/sample estimate at time  $t$   
trial =  $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶  $\alpha$  temporal difference update  
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{trial} - Q(S_t, A_t))$
- ▶  $S_t \leftarrow S_{t+1}$  and repeat (unless  $S_t$  is terminal)



...  $A_t$  derived from  $Q$

What about keeping optimality, taking max?

$$A_t = \arg \max_a Q(S_t, a)$$

see the demo run of `rl_agents.py`.

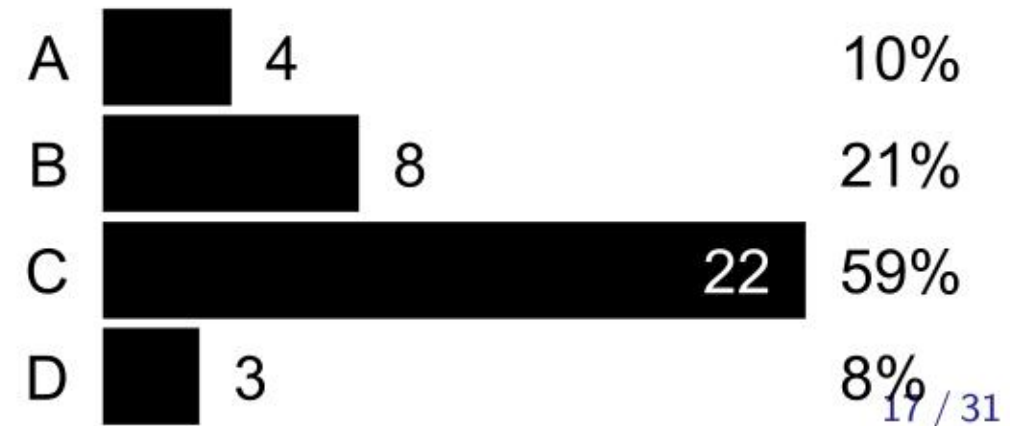
...  $A_t$  derived from  $Q$

What about keeping optimality, taking max?

*exploitation*

$$A_t = \arg \max_a Q(S_t, a)$$

see the demo run of `rl_agents.py`.



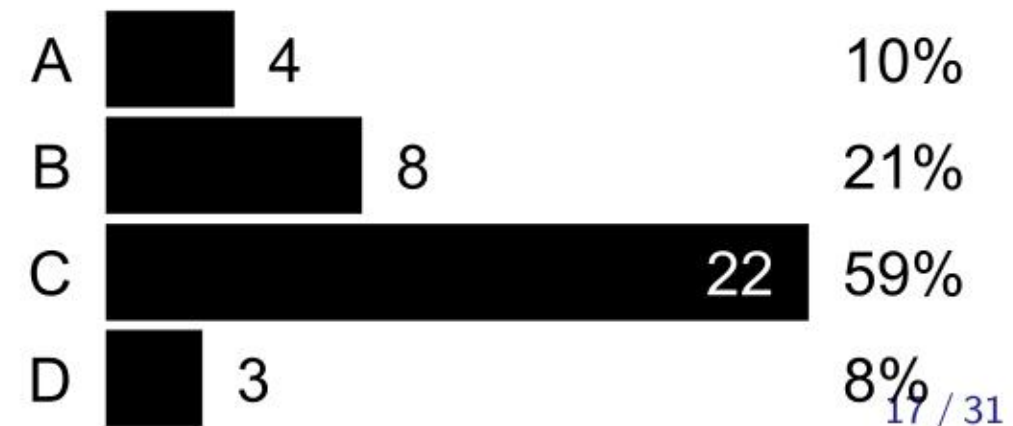
...  $A_t$  derived from  $Q$

What about keeping optimality, taking max?

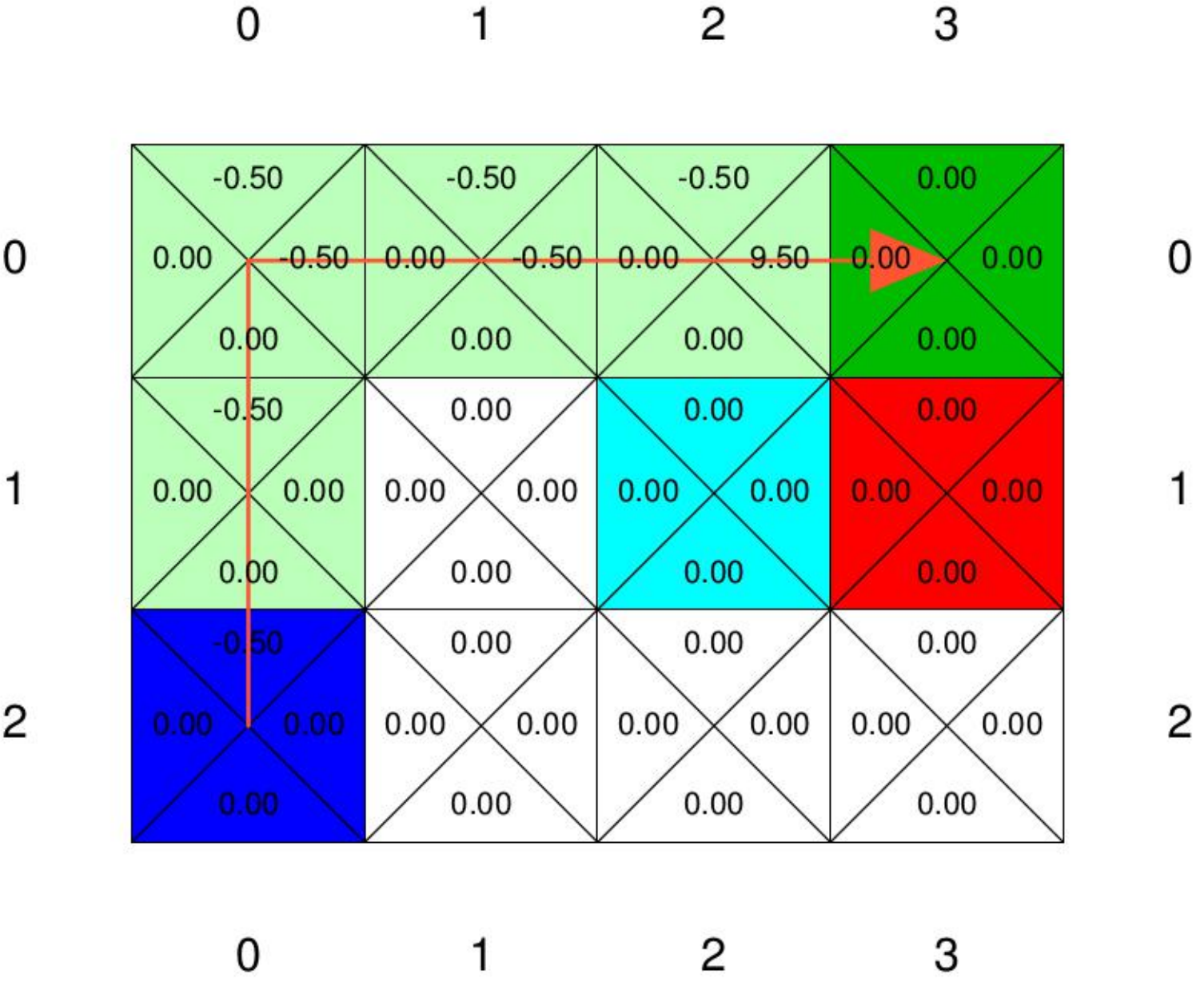
*exploitation*

$$A_t = \arg \max_a Q(S_t, a)$$

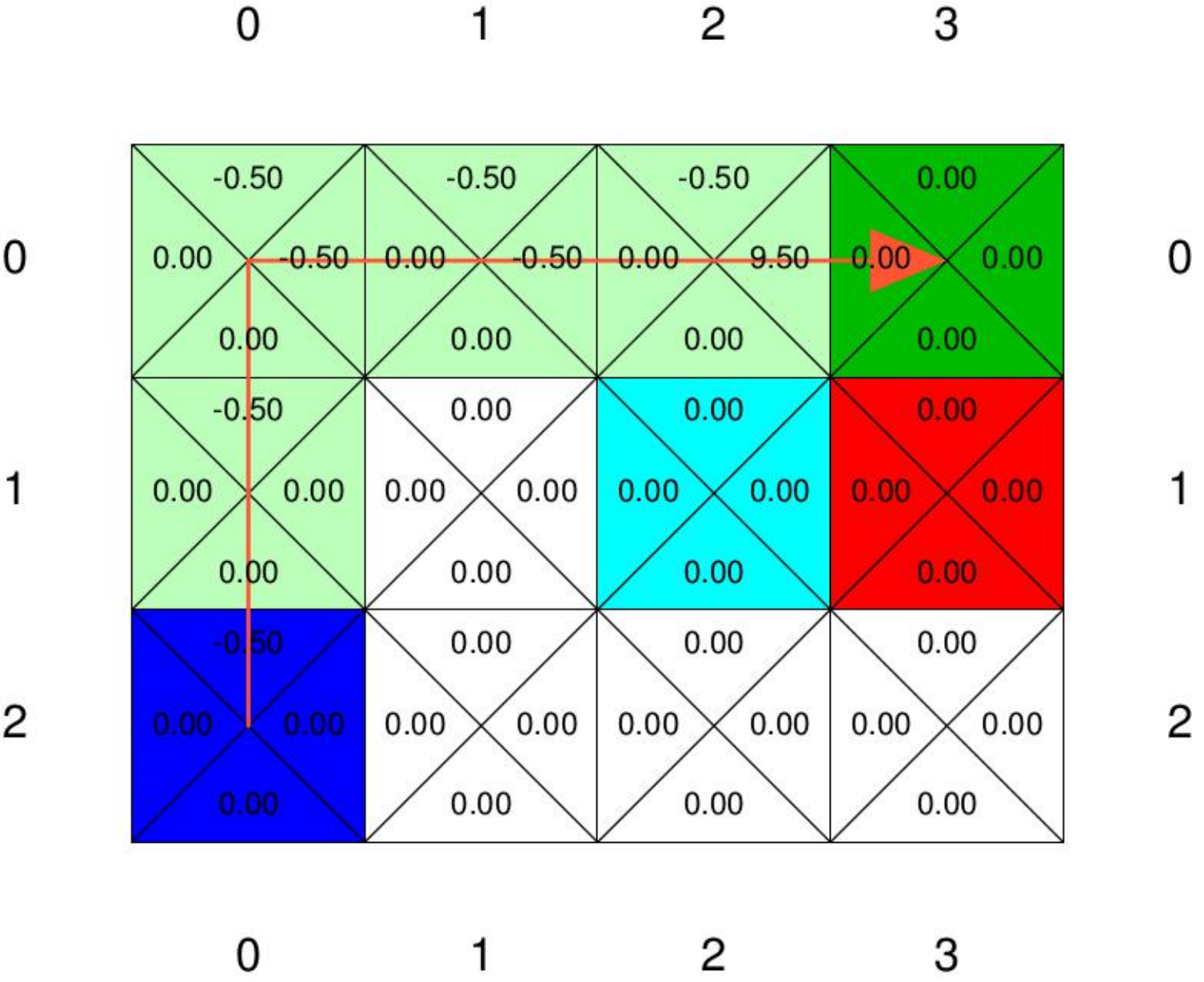
see the demo run of `rl_agents.py`.



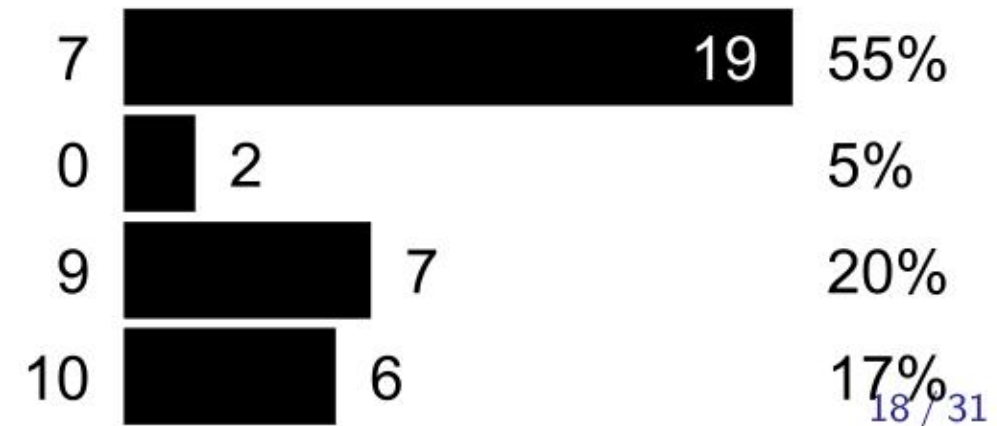
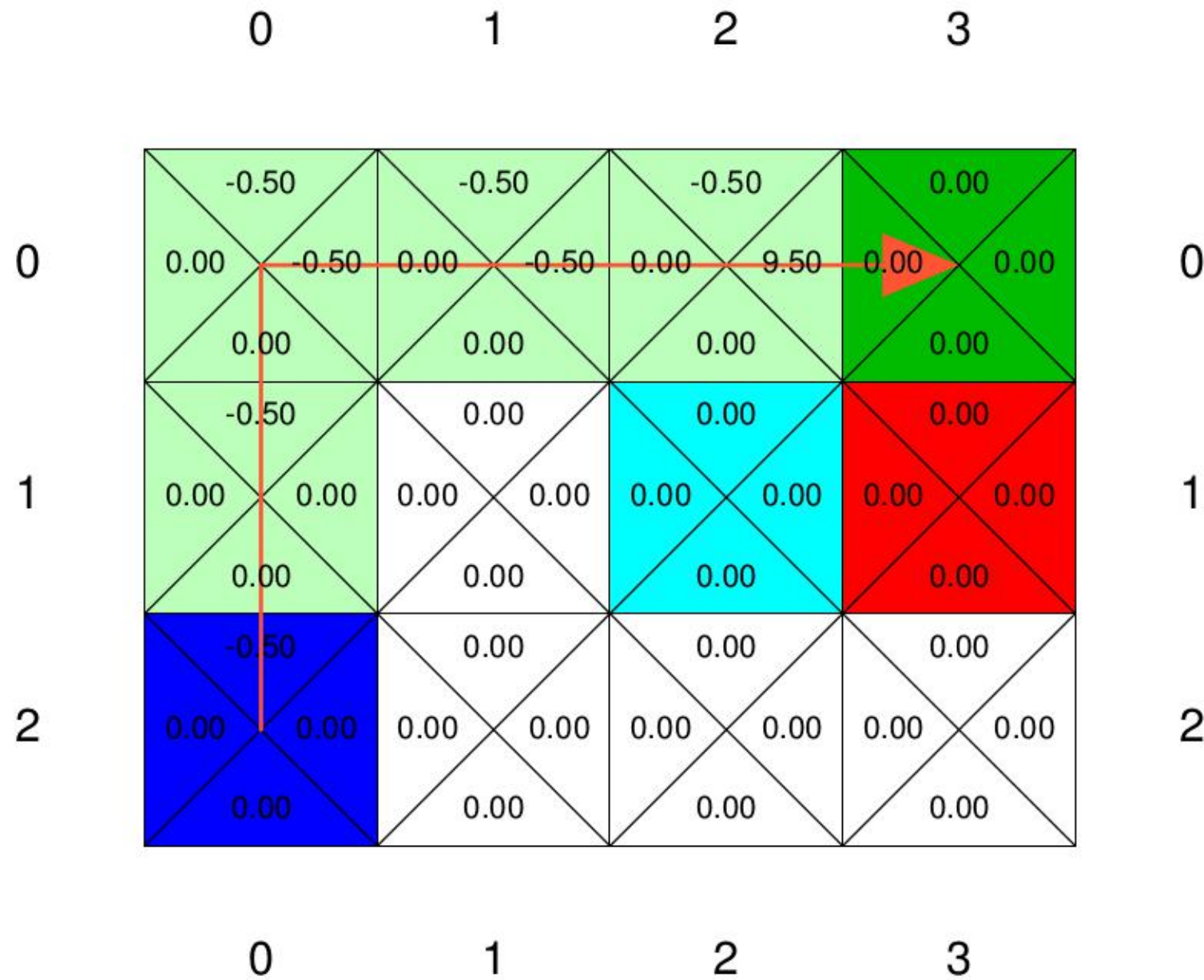
# Two good goal states



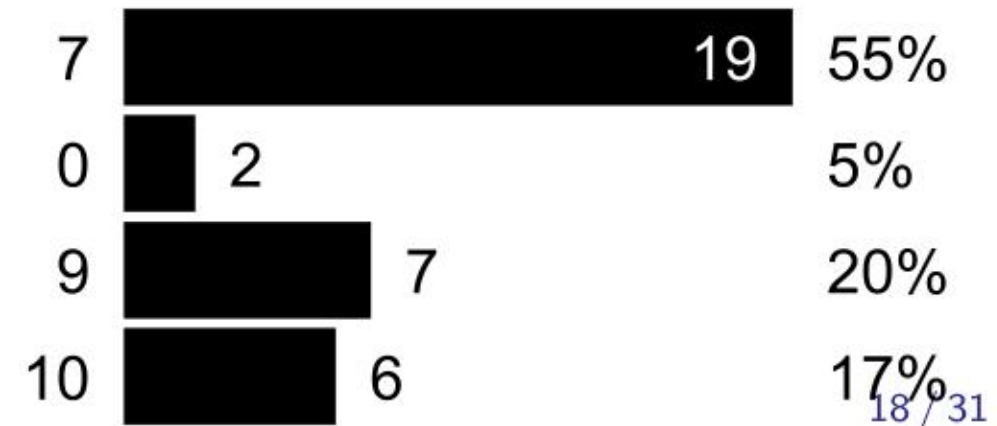
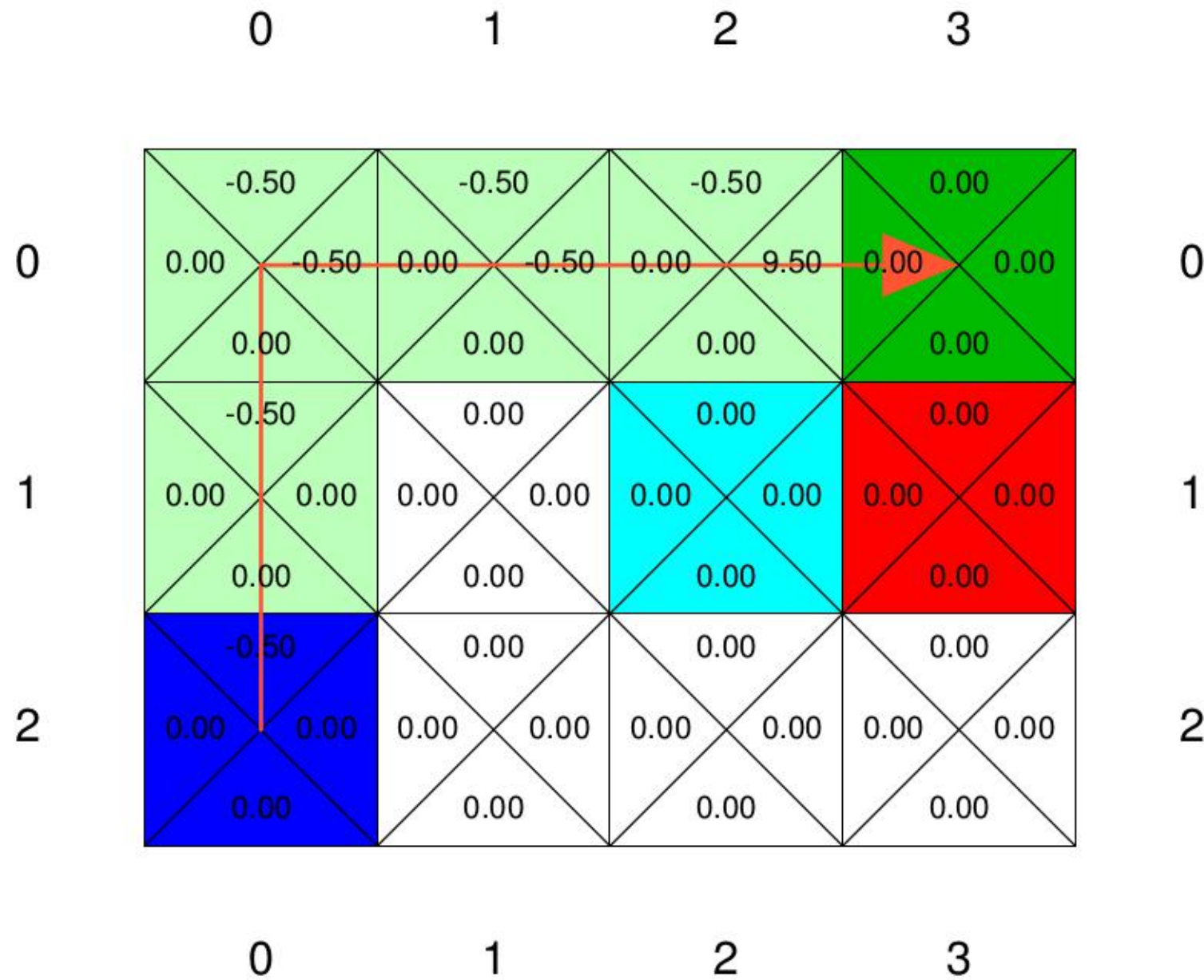
# Two good goal states



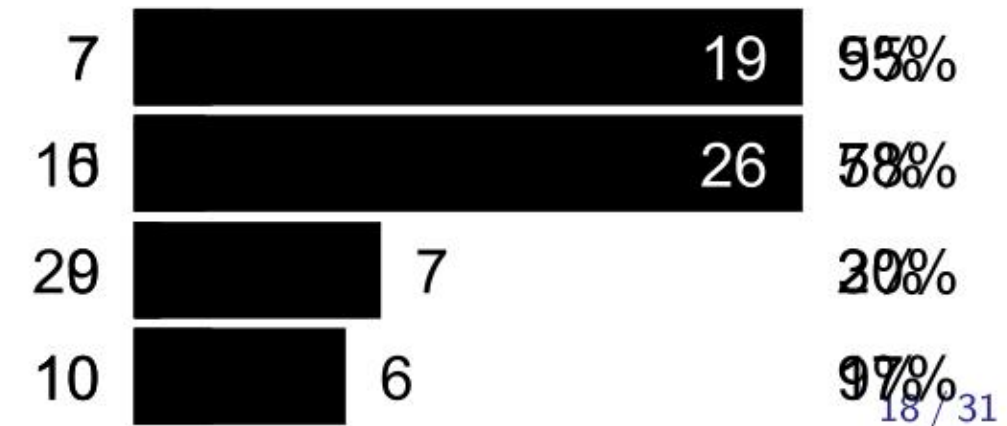
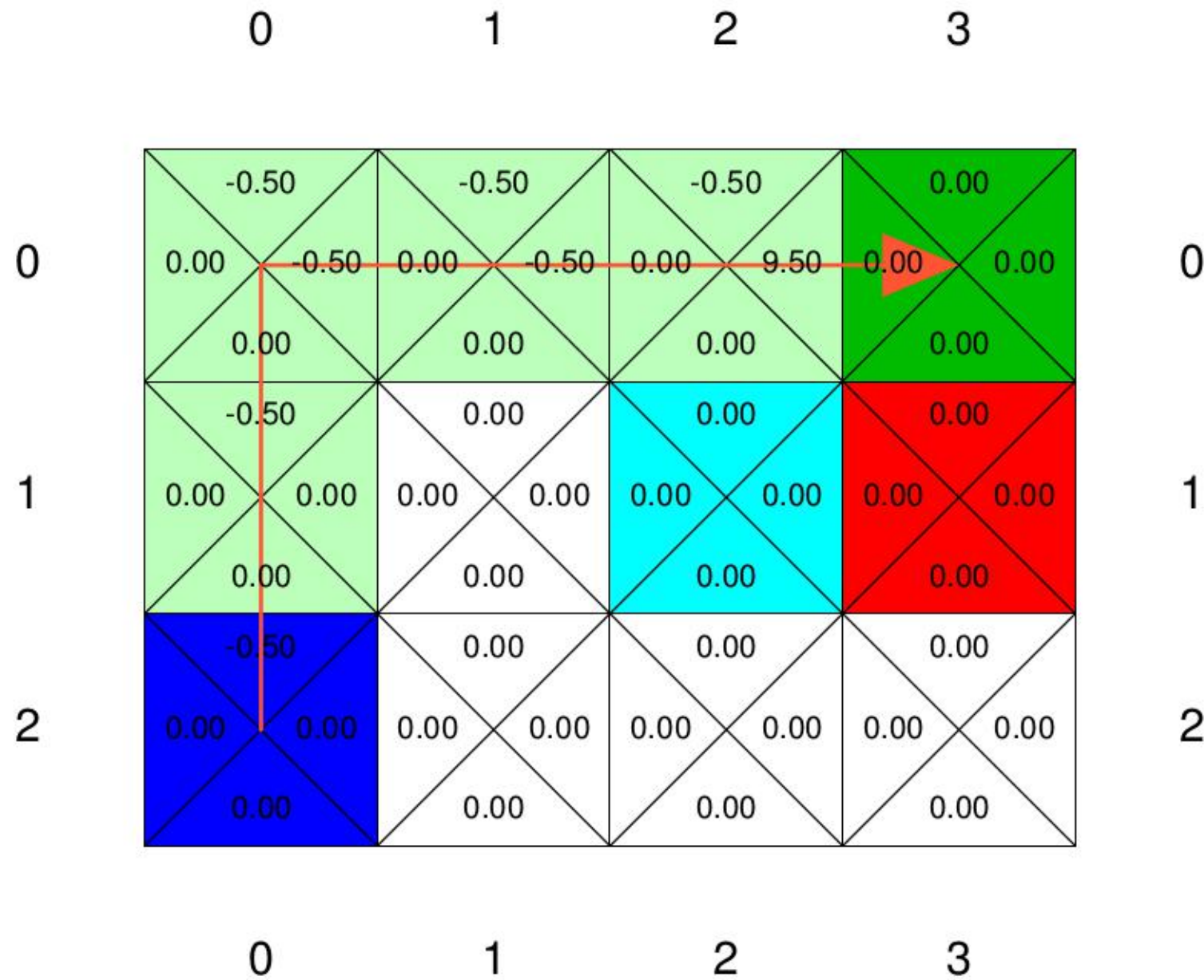
# Two good goal states



# Two good goal states

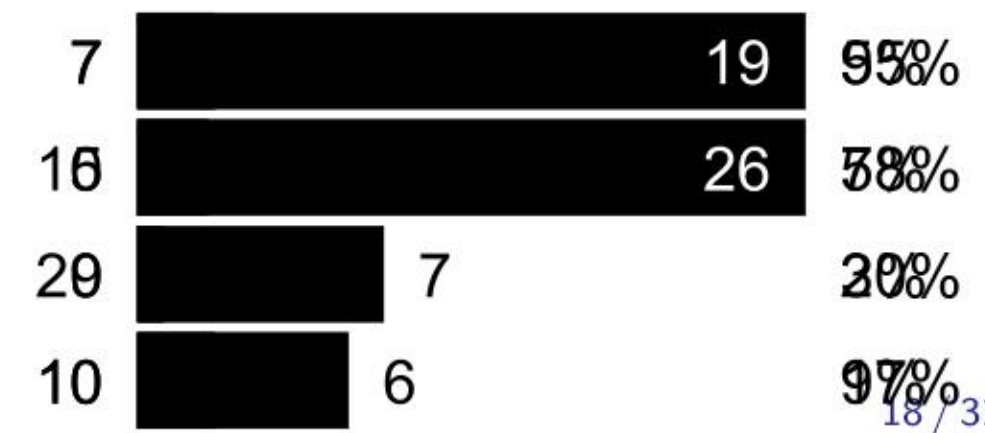
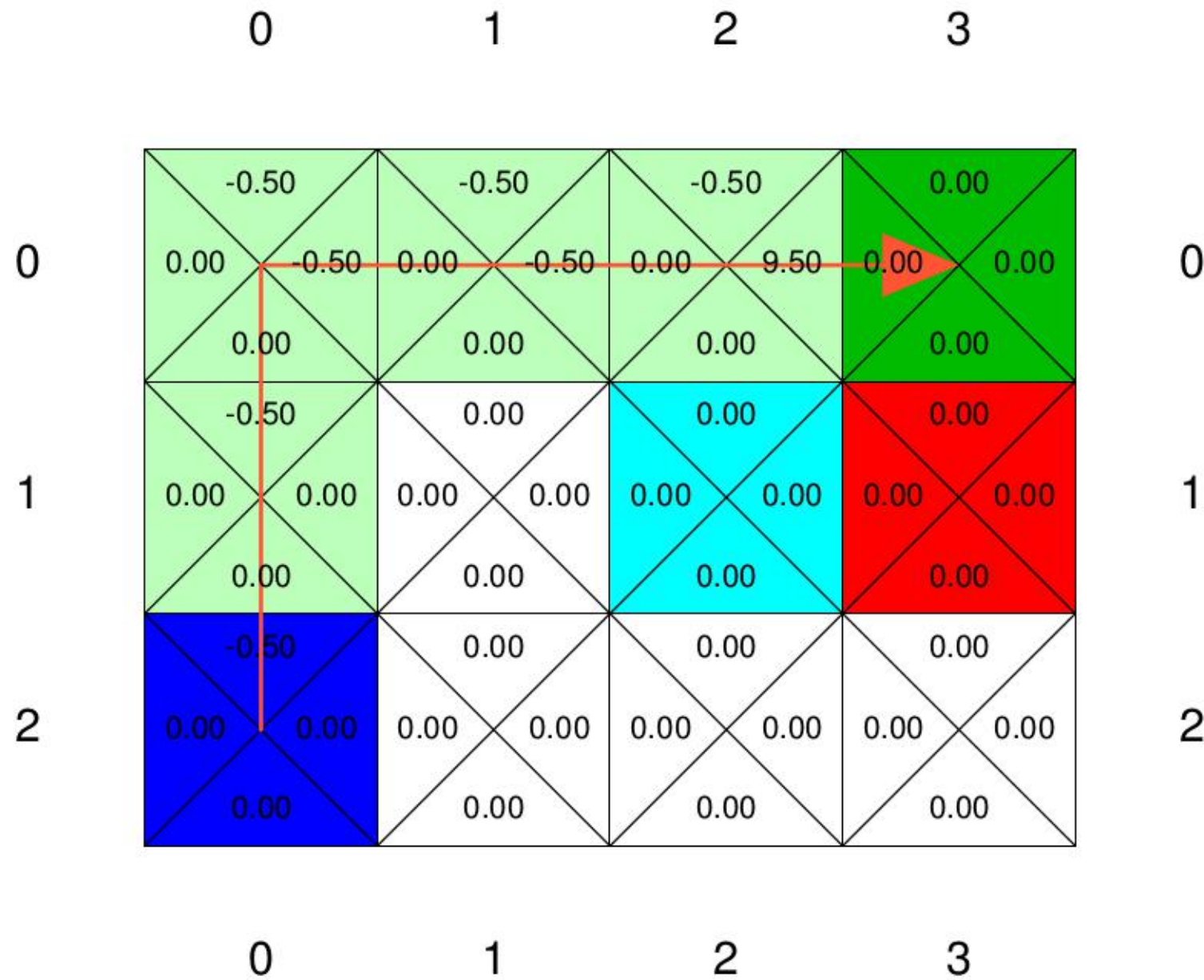


# Two good goal states

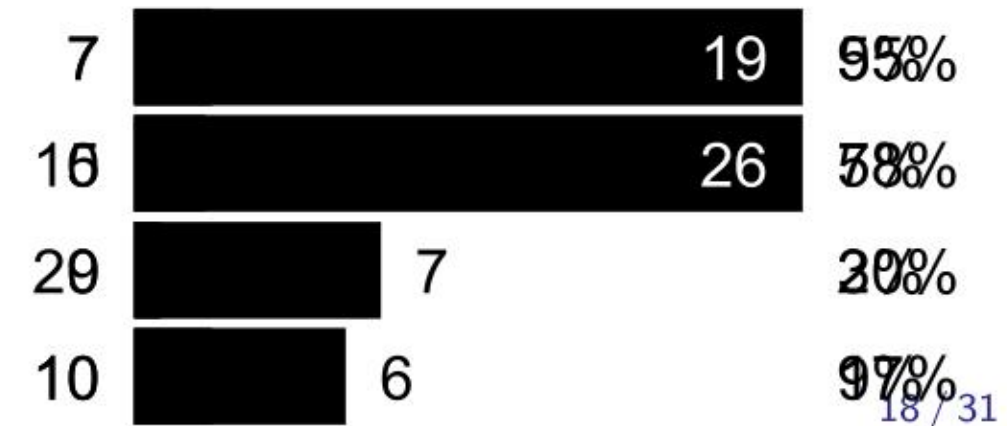
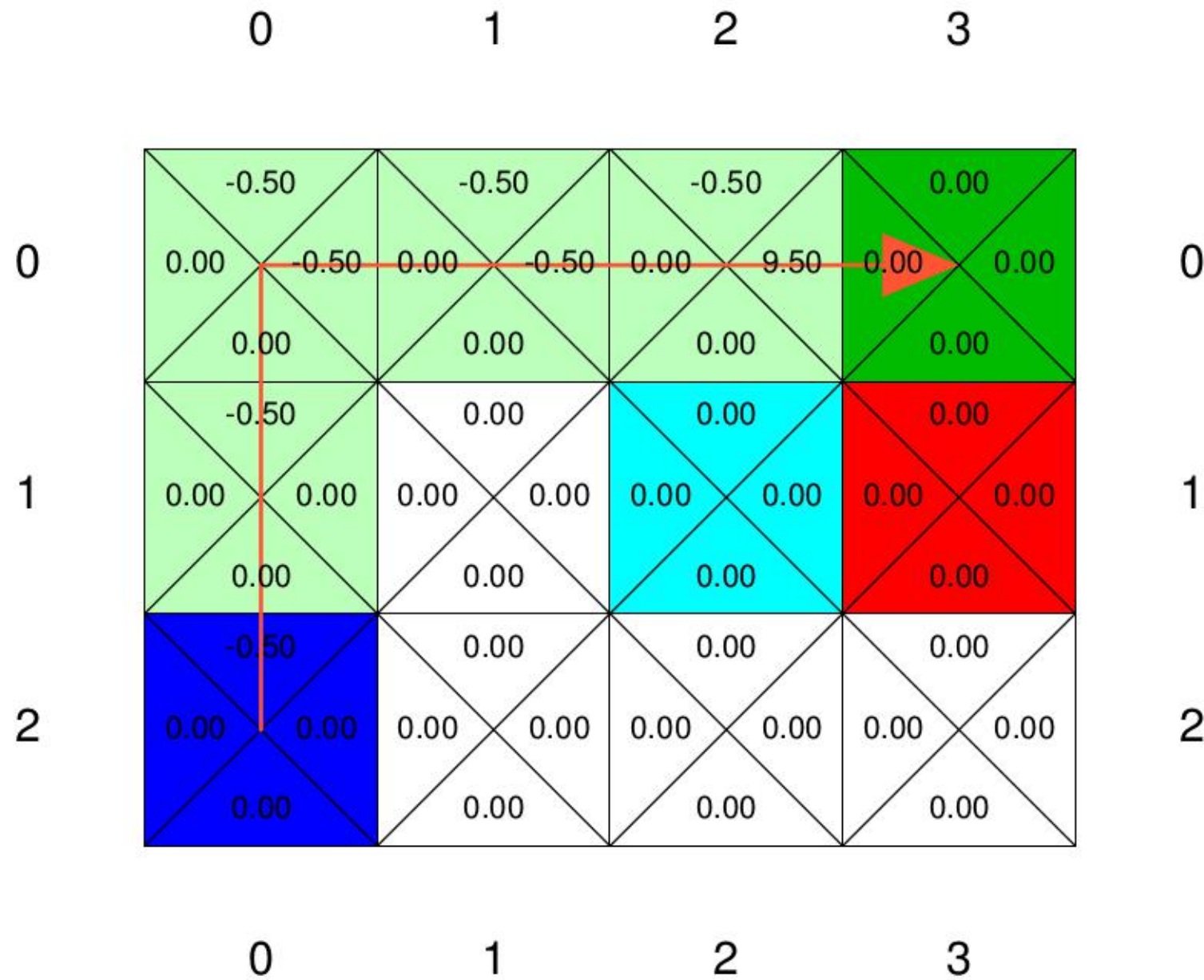




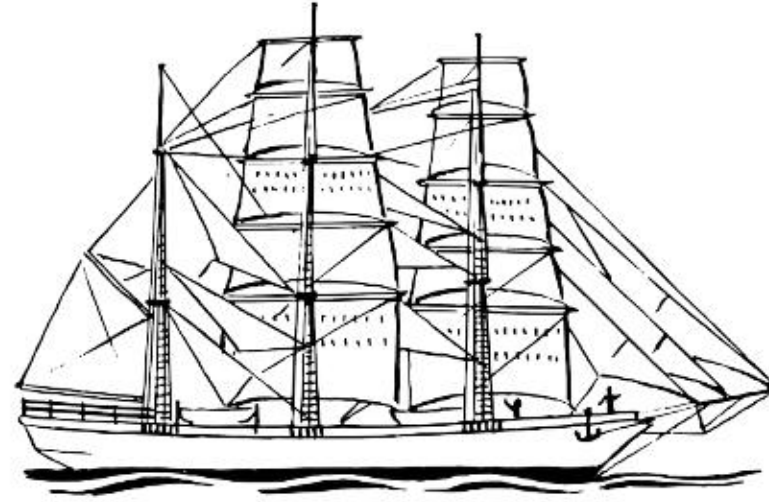
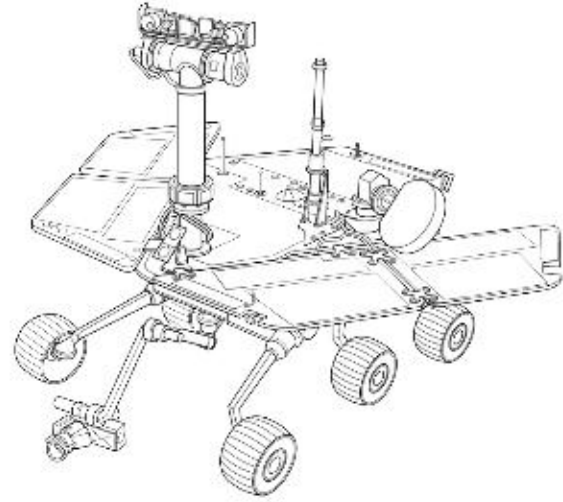
# Two good goal states



# Two good goal states

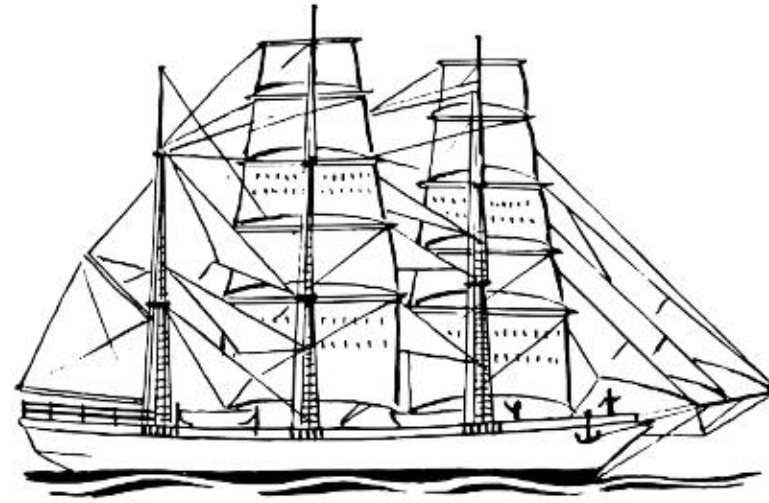
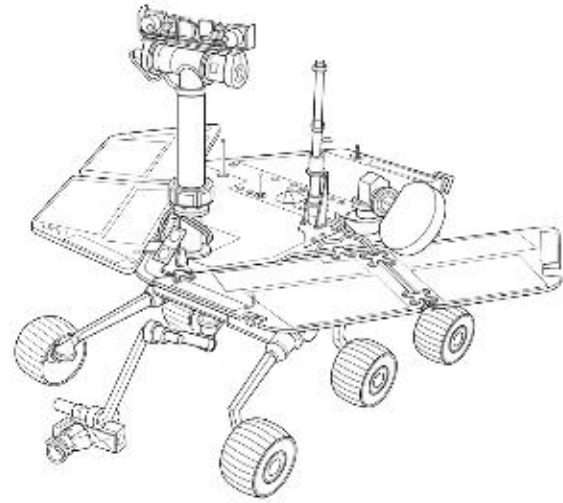


# Exploration vs Exploitation



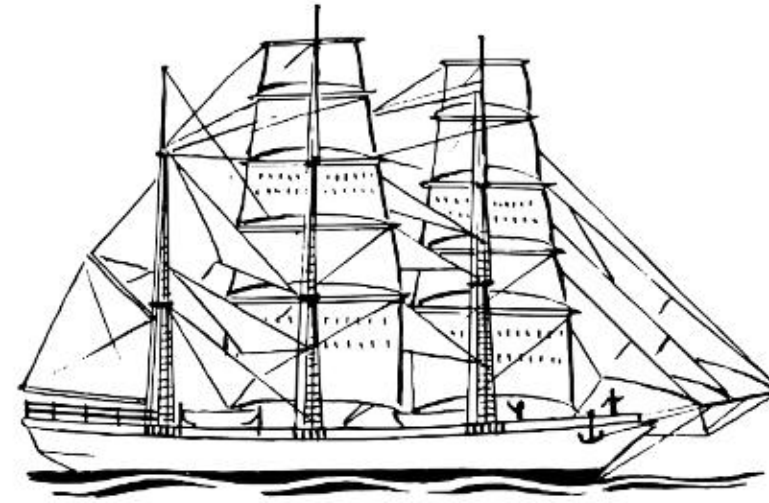
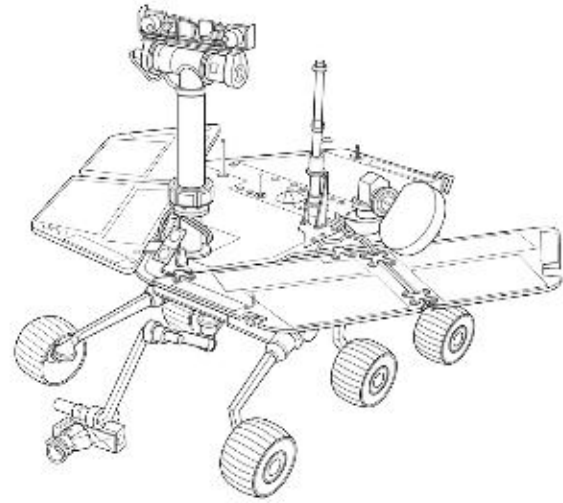
- ▶ Drive the known road or try a new one?
- ▶ Go to the university menza or try a nearby restaurant?
- ▶ Use the SW (operating system) I know or try new one?
- ▶ Go to bussiness or study a demanding program?
- ▶ ...

# Exploration vs Exploitation



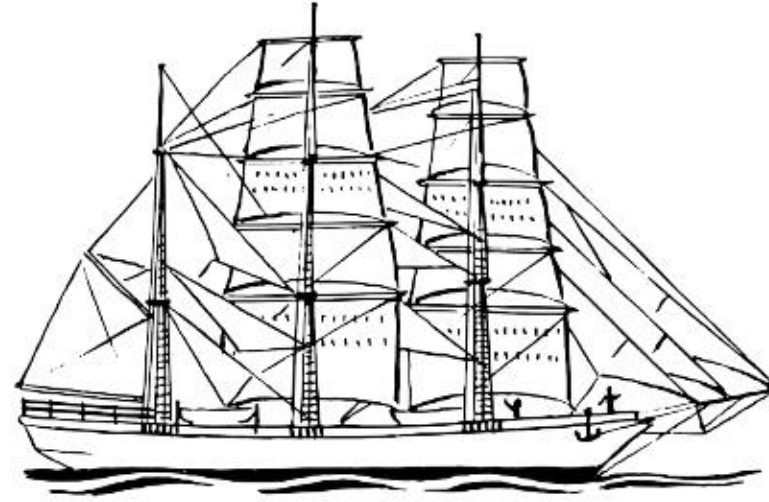
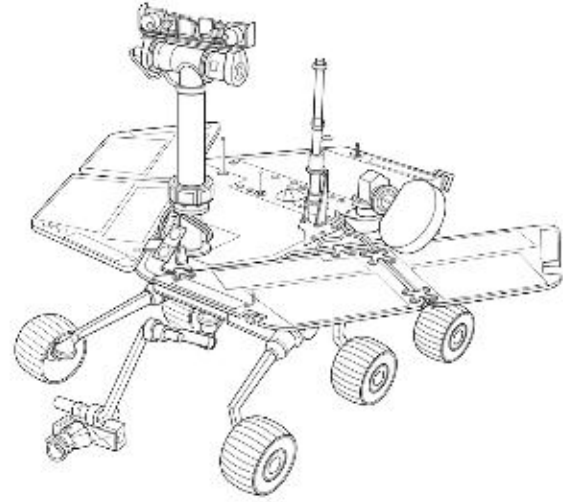
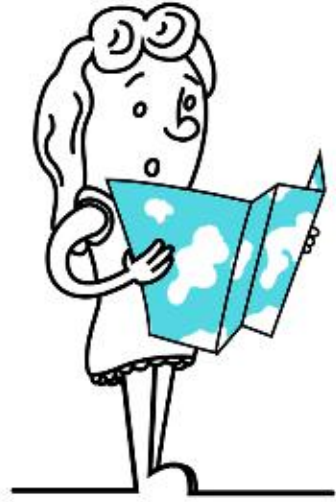
- ▶ Drive the known road or try a new one?
- ▶ Go to the university menza or try a nearby restaurant?
- ▶ Use the SW (operating system) I know or try new one?
- ▶ Go to bussiness or study a demanding program?
- ▶ ...

# Exploration vs Exploitation



- ▶ Drive the known road or try a new one?
- ▶ Go to the university menza or try a nearby restaurant?
- ▶ Use the SW (operating system) I know or try new one?
- ▶ Go to bussiness or study a demanding program?
- ▶ ...

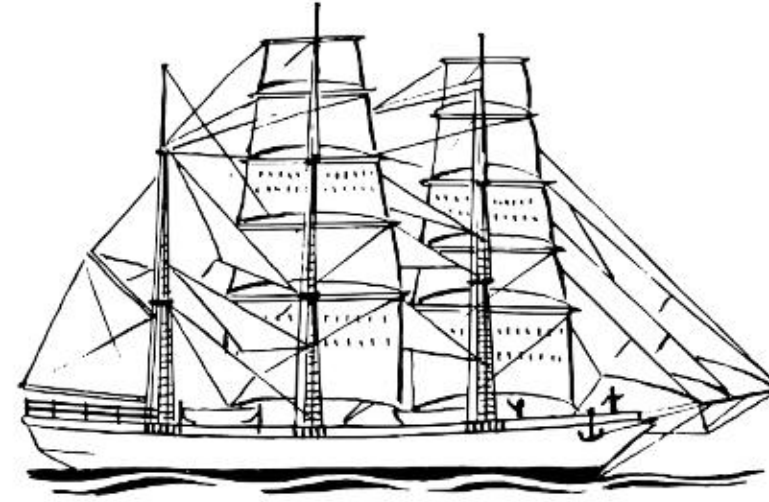
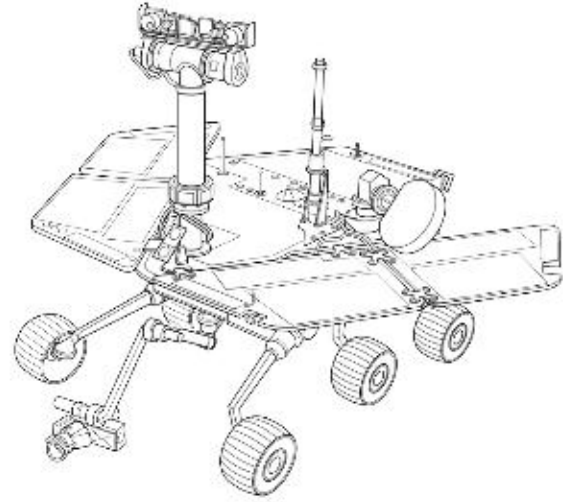
# Exploration vs Exploitation



- ▶ Drive the known road or try a new one?
- ▶ Go to the university menza or try a nearby restaurant?
- ▶ Use the SW (operating system) I know or try new one?
- ▶ Go to bussiness or study a demanding program?



# Exploration vs Exploitation



- ▶ Drive the known road or try a new one?
- ▶ Go to the university menza or try a nearby restaurant?
- ▶ Use the SW (operating system) I know or try new one?
- ▶ Go to bussiness or study a demanding program?
- ▶ ....

# How to explore?

## Random ( $\epsilon$ -greedy):

- ▶ Flip a coin every step.
- ▶ With probability  $\epsilon$ , act randomly.
- ▶ With probability  $1 - \epsilon$ , use the policy.

## Problems with randomness?

- ▶ Keeps exploring forever.
- ▶ Should we keep  $\epsilon$  fixed (over learning)?
- ▶  $\epsilon$  same everywhere?



# How to explore?

## Random ( $\epsilon$ -greedy):

- ▶ Flip a coin every step.
- ▶ With probability  $\epsilon$ , act randomly.
- ▶ With probability  $1 - \epsilon$ , use the policy.

## Problems with randomness?

- ▶ Keeps exploring forever.
- ▶ Should we keep  $\epsilon$  fixed (over learning)?
- ▶  $\epsilon$  same everywhere?

# How to explore?

## Random ( $\epsilon$ -greedy):

- ▶ Flip a coin every step.
- ▶ With probability  $\epsilon$ , act randomly.
- ▶ With probability  $1 - \epsilon$ , use the policy.

## Problems with randomness?

- ▶ Keeps exploring forever.
- ▶ Should we keep  $\epsilon$  fixed (over learning)?
- ▶  $\epsilon$  same everywhere?

# How to explore?

## Random ( $\epsilon$ -greedy):

- ▶ Flip a coin every step.
- ▶ With probability  $\epsilon$ , act randomly.
- ▶ With probability  $1 - \epsilon$ , use the policy.

## Problems with randomness?

- ▶ Keeps exploring forever.
- ▶ Should we keep  $\epsilon$  fixed (over learning)?
- ▶  $\epsilon$  same everywhere?

# How to explore?

## Random ( $\epsilon$ -greedy):

- ▶ Flip a coin every step.
- ▶ With probability  $\epsilon$ , act randomly.
- ▶ With probability  $1 - \epsilon$ , use the policy.

## Problems with randomness?

- ▶ Keeps exploring forever.
- ▶ Should we keep  $\epsilon$  fixed (over learning)?
- ▶  $\epsilon$  same everywhere?

# How to explore?

## Random ( $\epsilon$ -greedy):

- ▶ Flip a coin every step.
- ▶ With probability  $\epsilon$ , act randomly.
- ▶ With probability  $1 - \epsilon$ , use the policy.

## Problems with randomness?

- ▶ Keeps exploring forever.
- ▶ Should we keep  $\epsilon$  fixed (over learning)?
- ▶  $\epsilon$  same everywhere?

# How to explore?

## Random ( $\epsilon$ -greedy):

- ▶ Flip a coin every step.
- ▶ With probability  $\epsilon$ , act randomly.
- ▶ With probability  $1 - \epsilon$ , use the policy.

## Problems with randomness?

- ▶ Keeps exploring forever.
- ▶ Should we keep  $\epsilon$  fixed (over learning)?
- ▶  $\epsilon$  same everywhere?

# How to explore?

## Random ( $\epsilon$ -greedy):

- ▶ Flip a coin every step.
- ▶ With probability  $\epsilon$ , act randomly.
- ▶ With probability  $1 - \epsilon$ , use the policy.

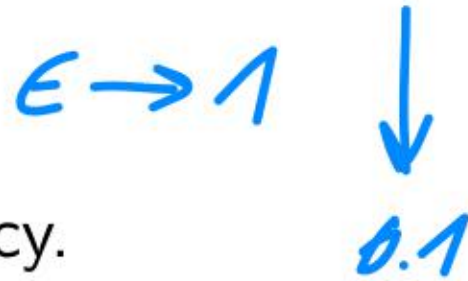
## Problems with randomness?

- ▶ Keeps exploring forever.
- ▶ Should we keep  $\epsilon$  fixed (over learning)?
- ▶  $\epsilon$  same everywhere?

# How to explore?

## Random ( $\epsilon$ -greedy):

- ▶ Flip a coin every step.
- ▶ With probability  $\epsilon$ , act randomly.
- ▶ With probability  $1 - \epsilon$ , use the policy.

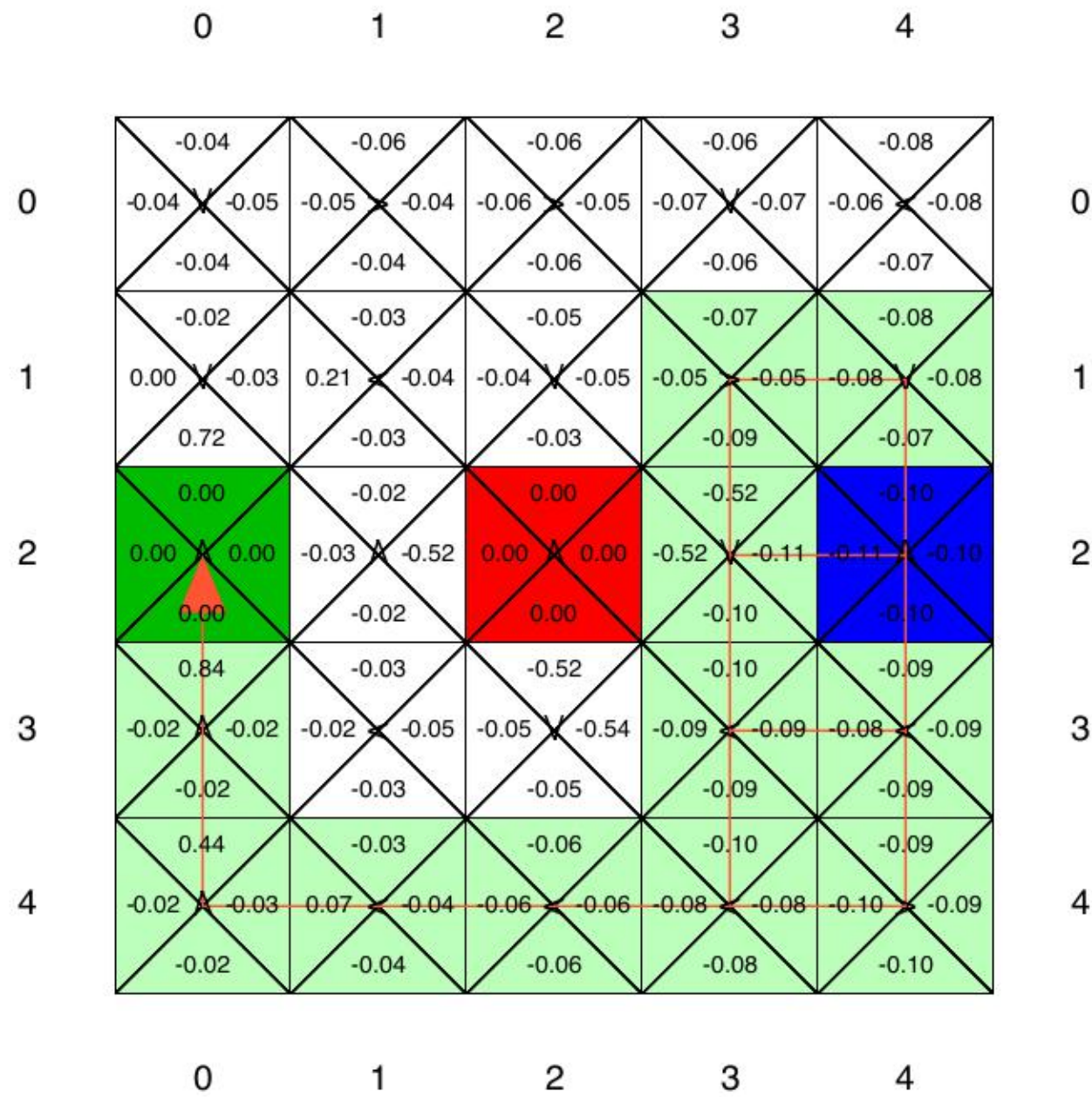


## Problems with randomness?

- ▶ Keeps exploring forever.
- ▶ Should we keep  $\epsilon$  fixed (over learning)?
- ▶  $\epsilon$  same everywhere?



# How to evaluate result, when to stop learning?

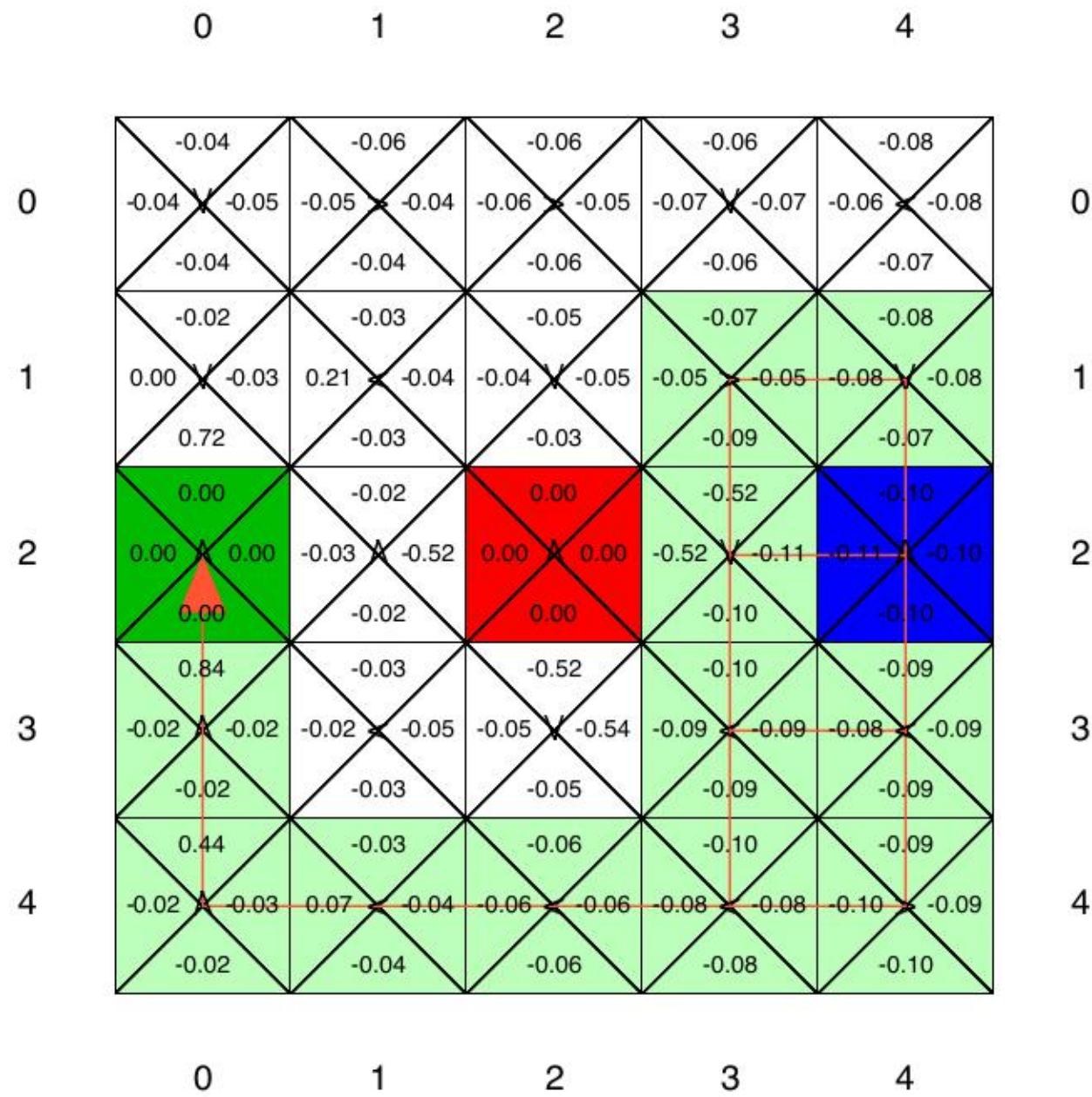


► What is the actual result of the q-learning?

► How to evaluate it?

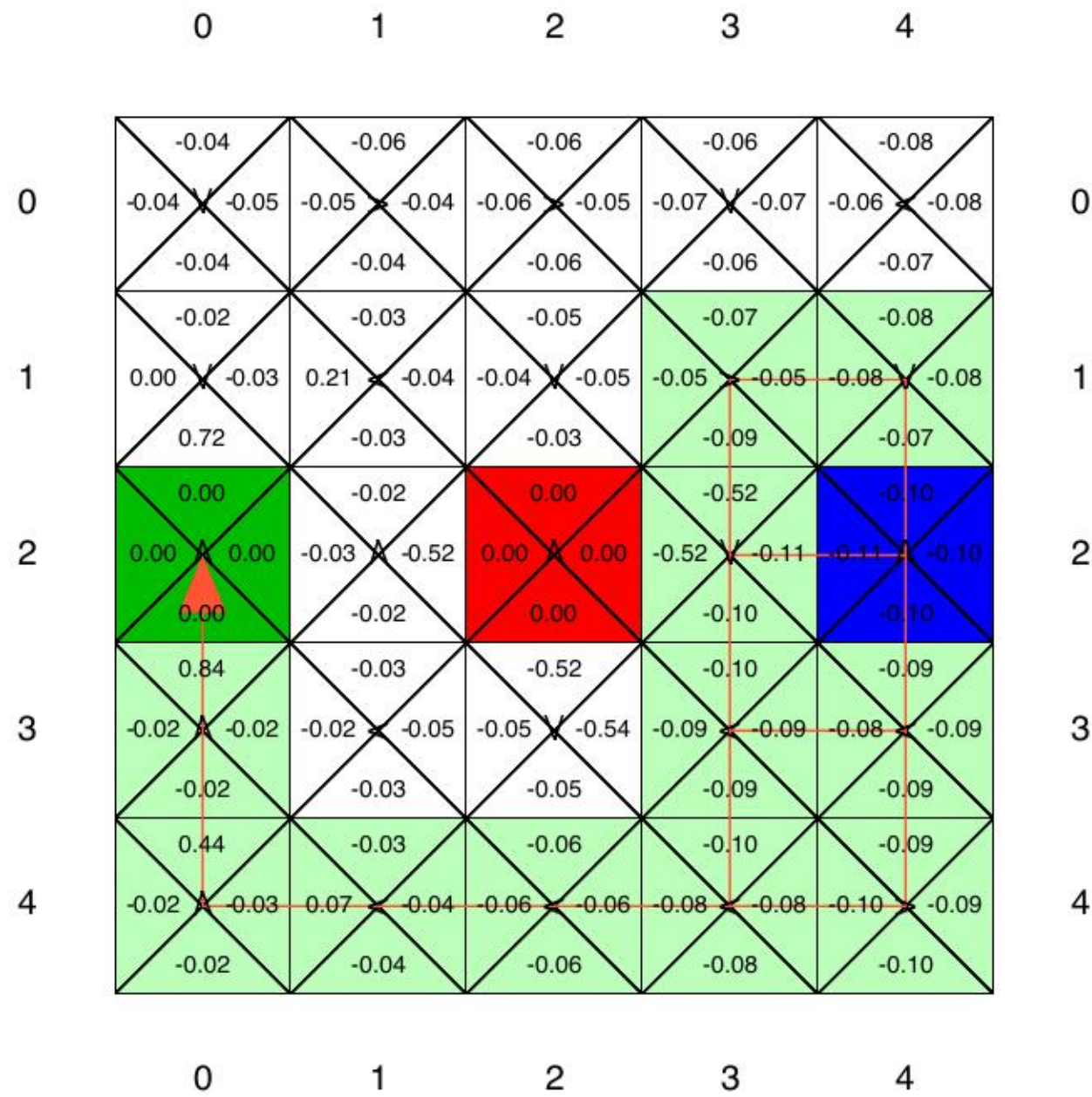
► When to stop learning?

# How to evaluate result, when to stop learning?



- ▶ What is the actual result of the q-learning?
- ▶ How to evaluate it?
- ▶ When to stop learning?

# How to evaluate result, when to stop learning?



- ▶ What is the actual result of the q-learning?
- ▶ How to evaluate it?
- ▶ When to stop learning?

## Exploration function $f(u, n)$

▶ Regular trial/sample estimate:  $\text{trial} = R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$

▶ If  $(S_t, a)$  not yet tried, then perhaps too pessimistic.

▶  $\text{trial} = R_{t+1} + \gamma \max_a f(Q(S_{t+1}, a), N(S_{t+1}, a))$

where  $f(u, n)$

$$\begin{aligned} f(u, n) &= R^+ \text{ if } n < N_e \\ &= u \text{ otherwise} \end{aligned}$$

where

▶  $R^+$  is an optimistic estimate of the best possible reward obtainable in any state

▶  $N_e$  fixed parameter

▶ The function  $f(u, n)$  should be increasing in  $u$  and decreasing in  $n$ .

## Exploration function $f(u, n)$

- ▶ Regular trial/sample estimate:  $\text{trial} = R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶ If  $(S_t, a)$  not yet tried, then perhaps too pessimistic.
- ▶  $\text{trial} = R_{t+1} + \gamma \max_a f(Q(S_{t+1}, a), N(S_{t+1}, a))$

where  $f(u, n)$

$$\begin{aligned} f(u, n) &= R^+ \text{ if } n < N_e \\ &= u \text{ otherwise} \end{aligned}$$

where

- ▶  $R^+$  is an optimistic estimate of the best possible reward obtainable in any state
- ▶  $N_e$  fixed parameter
- ▶ The function  $f(u, n)$  should be increasing in  $u$  and decreasing in  $n$ .

## Exploration function $f(u, n)$

- ▶ Regular trial/sample estimate:  $\text{trial} = R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶ If  $(S_t, a)$  not yet tried, then perhaps too pessimistic.
- ▶  $\text{trial} = R_{t+1} + \gamma \max_a f(Q(S_{t+1}, a), N(S_{t+1}, a))$

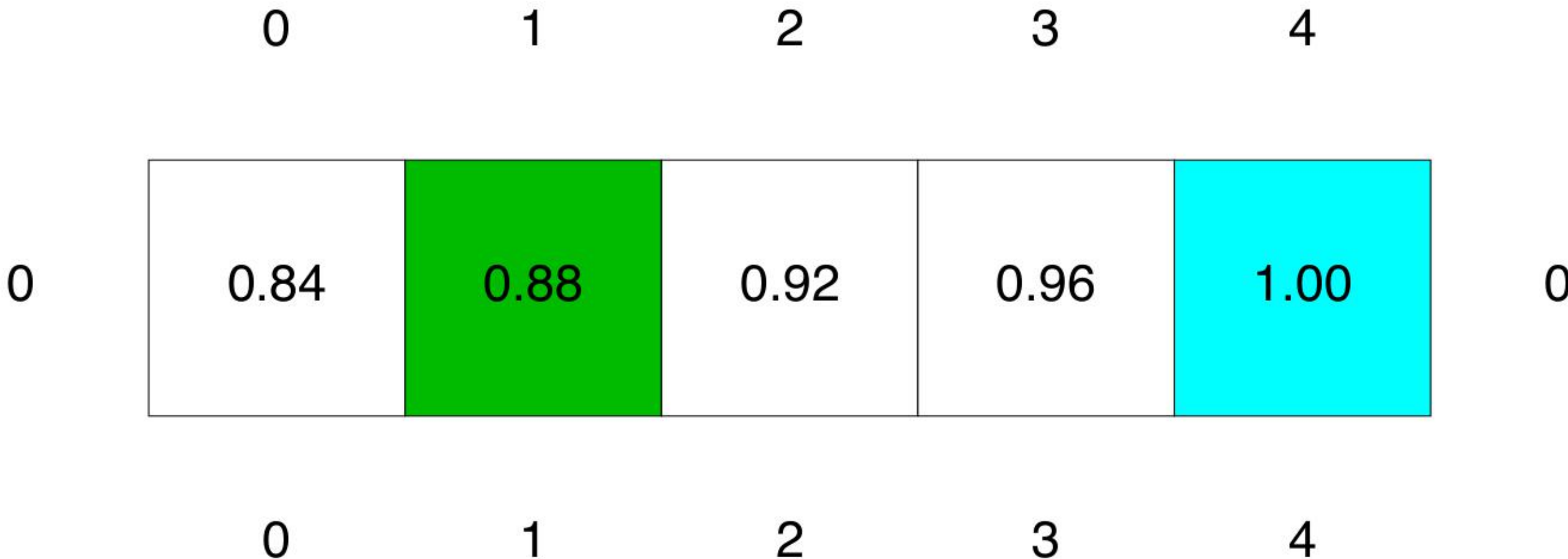
where  $f(u, n)$

$$\begin{aligned} f(u, n) &= R^+ \text{ if } n < N_e \\ &= u \text{ otherwise} \end{aligned}$$

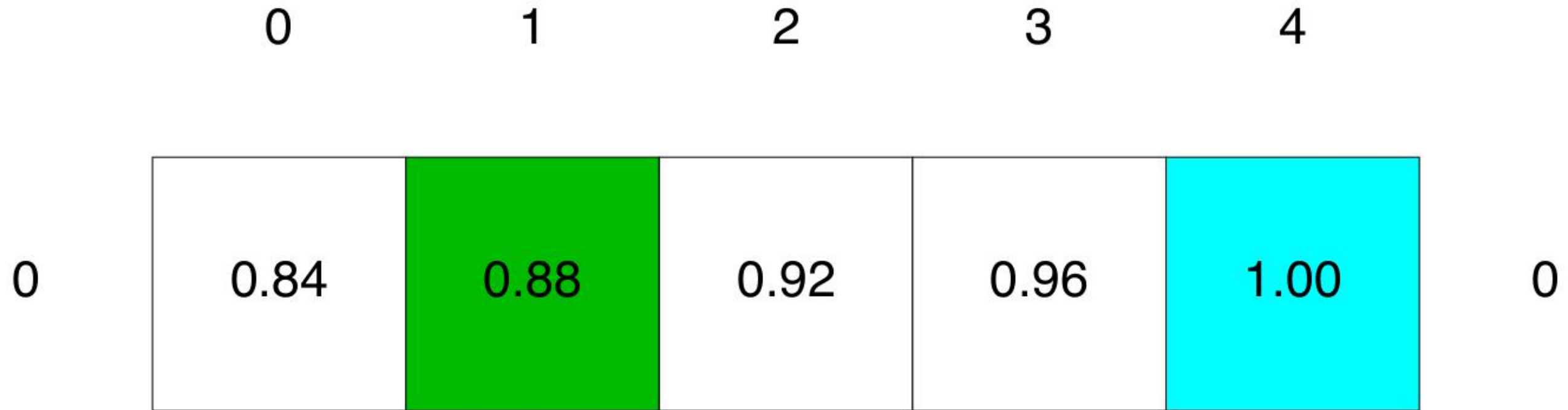
where

- ▶  $R^+$  is an optimistic estimate of the best possible reward obtainable in any state
- ▶  $N_e$  fixed parameter
- ▶ The function  $f(u, n)$  should be increasing in  $u$  and decreasing in  $n$ .

# Going beyond tables - generalizing across states



# Going beyond tables - generalizing across states



0      1      2      3      4

$v(s)$

$V = \begin{bmatrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}$

$Q(s,a) = \begin{bmatrix} \uparrow \\ \swarrow \\ \downarrow \\ \leftarrow \end{bmatrix}$

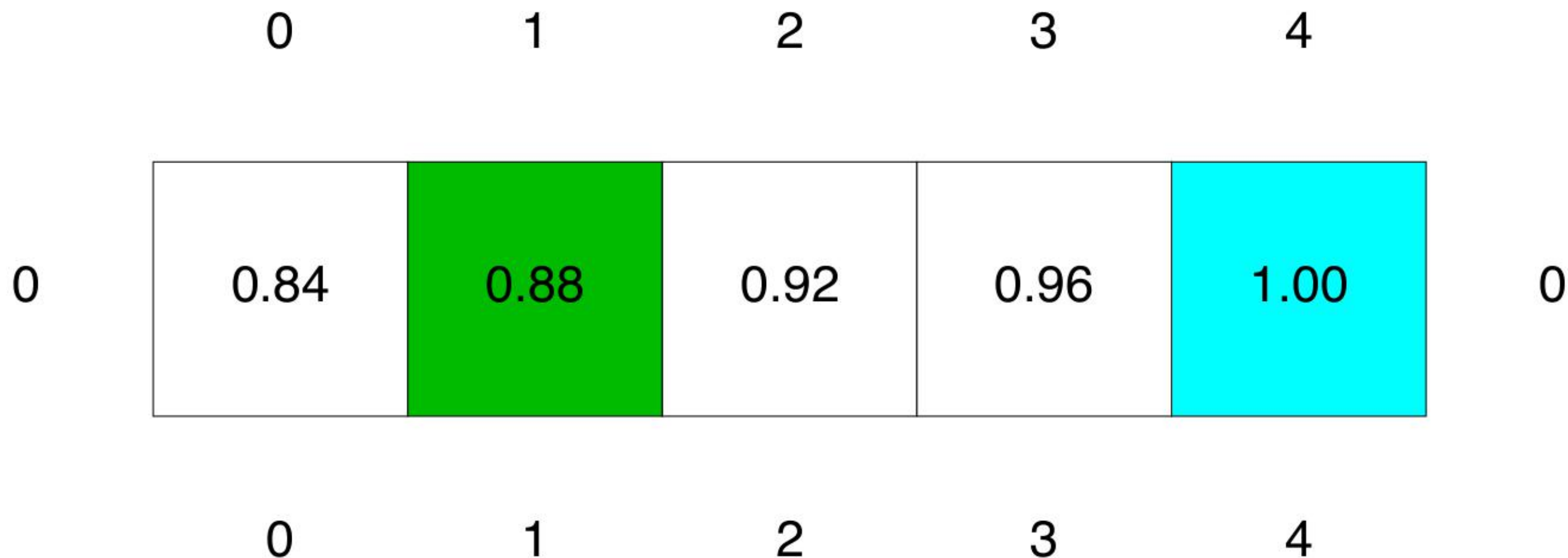
*q-table*



# Going beyond tables - generalizing across states

	0	1	2	3	4	
0	0.84	0.80	0.76	0.72	0.68	0
1	0.88	0.84	0.80	0.76	0.72	1
2	0.92	0.88	0.84	0.80	0.76	2
3	0.96	0.92	0.88	0.84	0.80	3
4	1.00	0.96	0.92	0.88	0.84	4
	0	1	2	3	4	

$v(s)$  not as table but as an approximation function

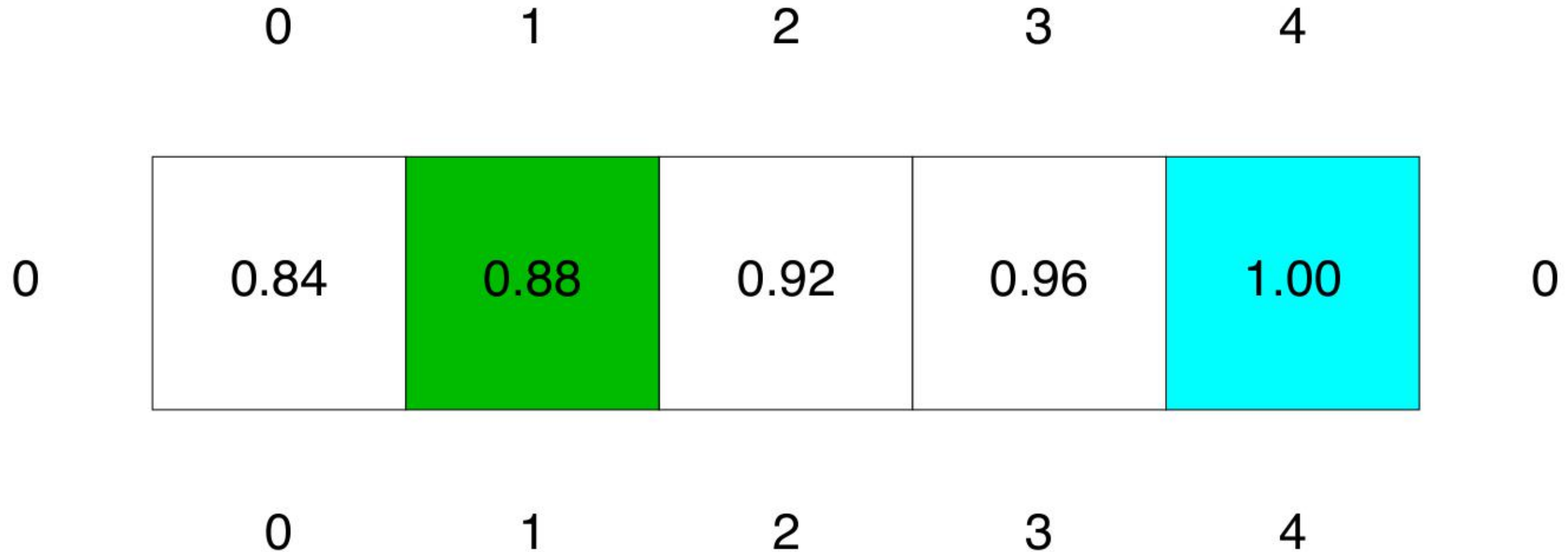


$$v(s, w) = w_0 + w_1 s$$

What are  $w_0, w_1$  equal to?

Instead of the complete table, only 2 parameters to learn  $w = [w_0, w_1]^T$

$v(s)$  not as table but as an approximation function

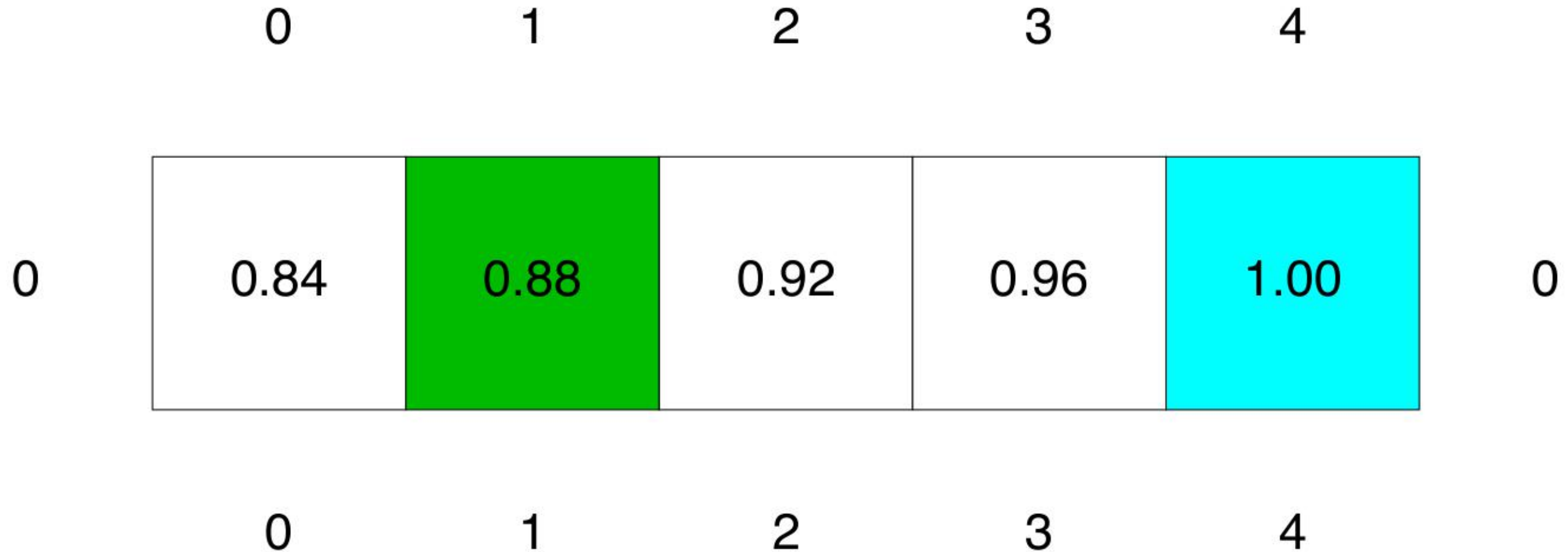


$$\hat{v}(s, \mathbf{w}) = w_0 + w_1 s$$

What are  $w_0, w_1$  equal to?

Instead of the complete table, only 2 parameters to learn  $\mathbf{w} = [w_0, w_1]^T$

$v(s)$  not as table but as an approximation function



$$\hat{v}(s, \mathbf{w}) = w_0 + w_1 s$$

What are  $w_0, w_1$  equal to?

Instead of the complete table, only 2 parameters to learn  $\mathbf{w} = [w_0, w_1]^T$

$v(s)$  not as table but as an approximation function

	0	1	2	3	4	
0	0.84	0.88	0.92	0.96	1.00	0

$$s = [0 \quad 1 \quad 2 \quad 3 \quad 4]$$

$$\hat{v}(s, \mathbf{w}) = w_0 + w_1 s$$

5 → 2

What are  $w_0, w_1$  equal to?

Instead of the complete table, only 2 parameters to learn  $\mathbf{w} = [w_0, w_1]^T$