

IN3050/IN4050, Lecture 12

Reinforcement Learning

1: Introduction

Ole Christian Lingjærde

Supervised learning

training set $(x_i, y_i), i=1, \dots, n$

Unsupervised learning

training set $x_i, i=1, \dots, n$

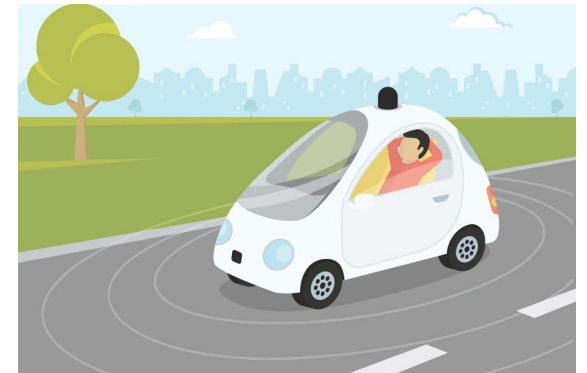
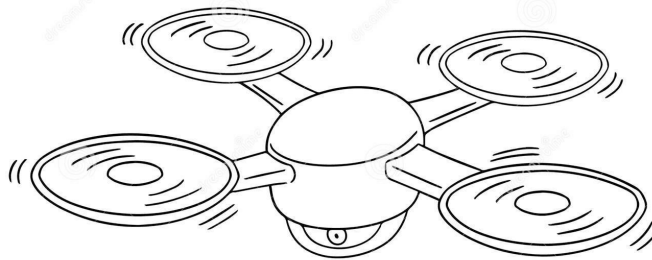
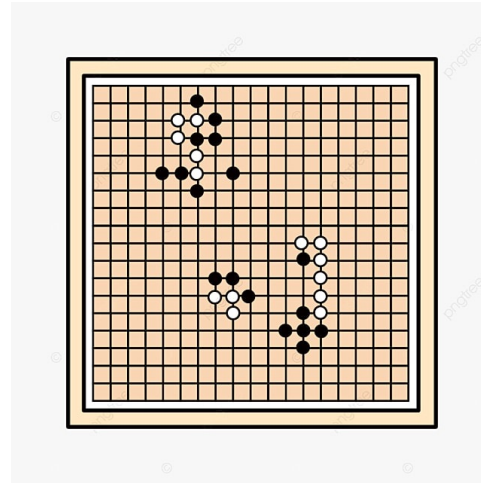
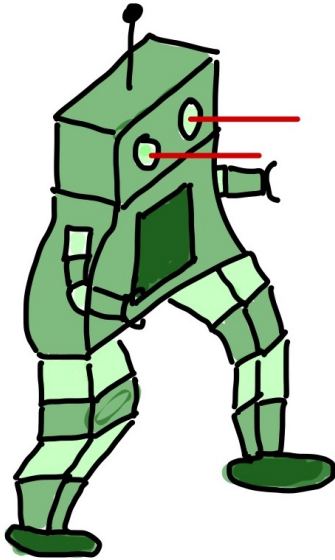
Reinforcement learning

training set generated dynamically

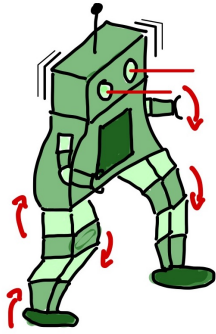
told if actions are good or bad

exploration to find right actions

Reinforcement learning (RL) applications



Robot with sensors and motors:



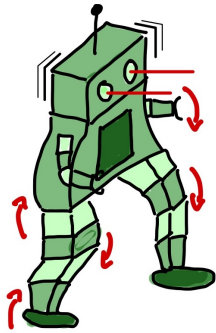
Observations:

- camera sensors
- touch sensors
- accelerometers
- microphones

Actions:

- motor torques

Robot with sensors and motors:



observation

camera image
touch sensor reading
sound
accel



action

motor
torques

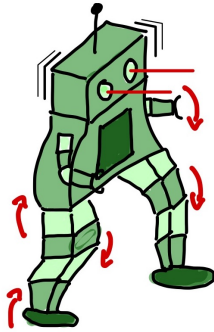
Observations:

- camera sensors
- touch sensors
- accelerometers
- microphones

Actions:

- motor torques

Robot with sensors and motors:



observation

camera image
touch sensor reading
sound
accel



action

motor
torques

Observations:

- camera sensors
- touch sensors
- accelerometers
- microphones

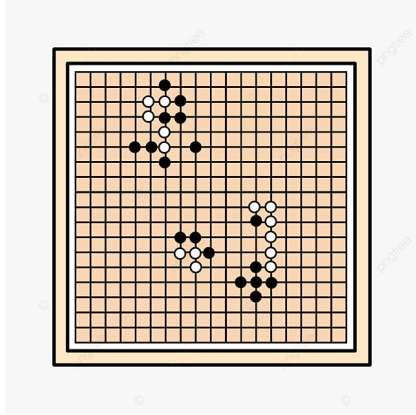
Actions:

- motor torques

Goal:

- keep balance
- walk past obstacles
- reach a destination

Computer playing a game



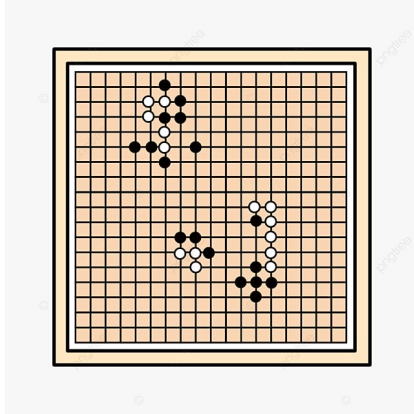
Observation:

- current state

Action:

- next move

Computer playing a game



observation

position of all pieces
on the board



action

next
move

Observation:

- current state

Action:

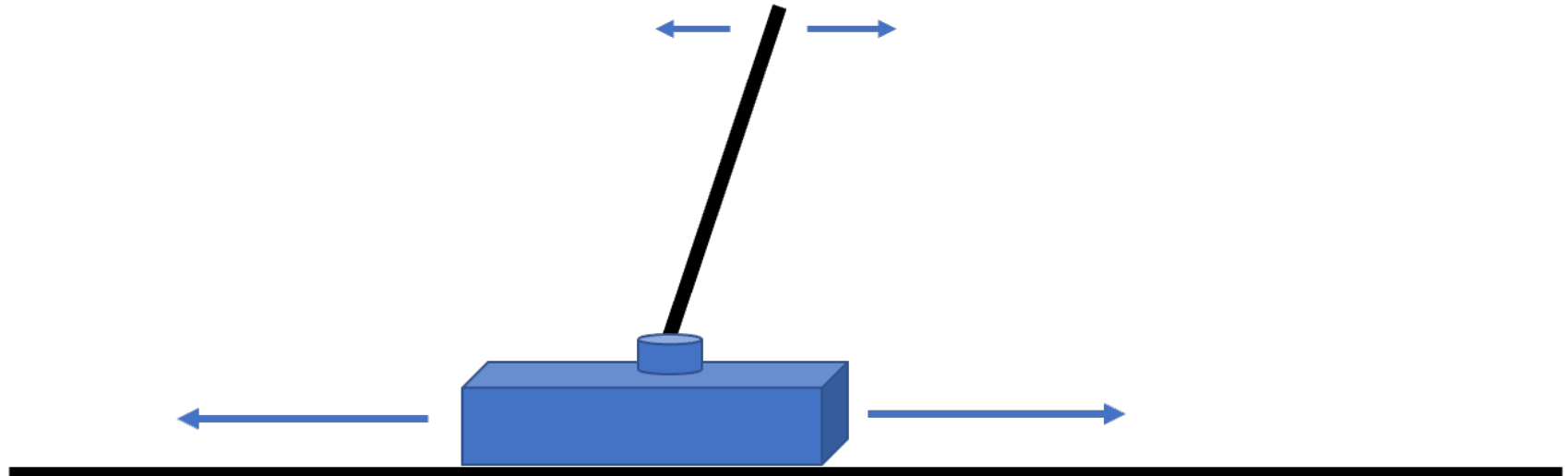
- next move

Goal: to win the game

Balancing a pole

Observations: angle between pole and cart

Actions: moving left or right





youtube.com is now full screen

Exit Full Screen (esc)



Performance after 0 episodes

Play (k)



0:09 / 2:25

Scroll for details

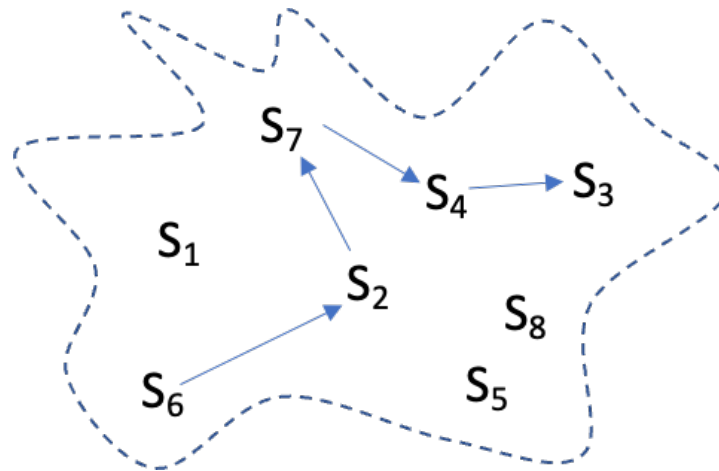


<https://www.youtube.com/watch?v=lfHX2hHRMVQ>

IN3050/IN4050, Lecture 12

Reinforcement Learning

2: Policy and states
Ole Christian Lingjærde



Policy

The way observations are mapped to actions defines a policy:



A policy can be **deterministic** or **stochastic**.

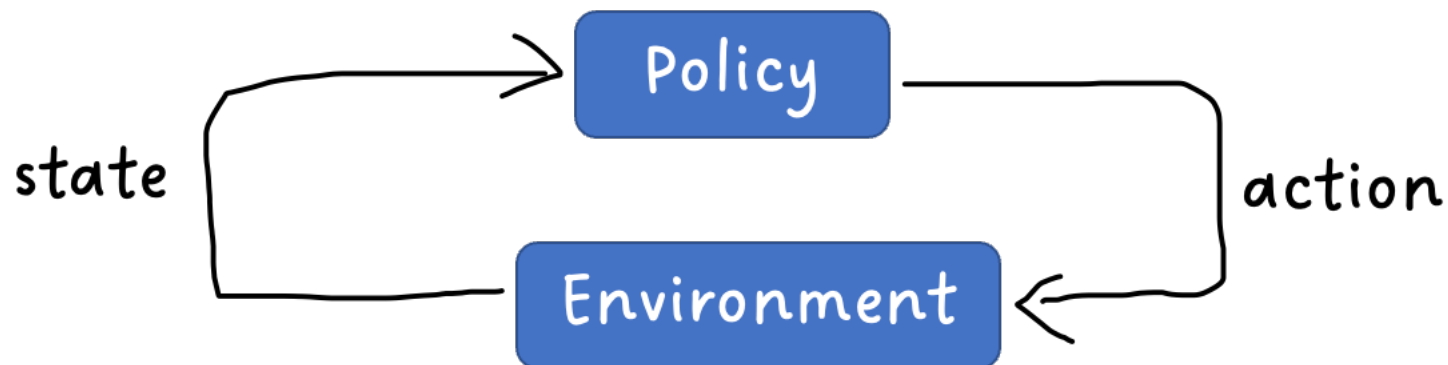
Exploitation vs. exploration

States

All RL systems work within an **environment**.

At any time the environment has a particular **state** and this is what we observe.

Actions can **change** the state:



Algorithm:

$s[0] = \text{<initial state>}$

for i in range(N):

$a = \text{action}(s[i])$

$s[i+1] = \text{newstate}(s[i], a)$

Algorithm:

$s[0] = \text{<initial state>}$

for i in range(N):

$a = \text{action}(s[i])$  implements the policy


$s[i+1] = \text{newstate}(s[i], a)$

Algorithm:

$s[0] = \text{<initial state>}$

for i in range(N):

$a = \text{action}(s[i])$  implements the policy


$s[i+1] = \text{newstate}(s[i], a)$  implements the environment -
internal logic not known

Algorithm:

$s[0] = \text{<initial state>}$

for i in range(N):

$a = \text{action}(s[i])$  implements the policy

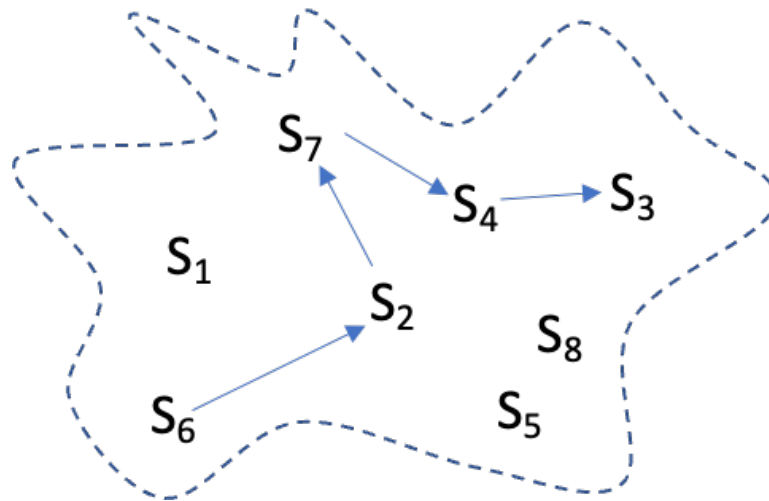
$s[i+1] = \text{newstate}(s[i], a)$  implements the
environment -
internal logic not
known

Result: $s[0]$ $s[1]$ $s[2]$ $s[N-1]$

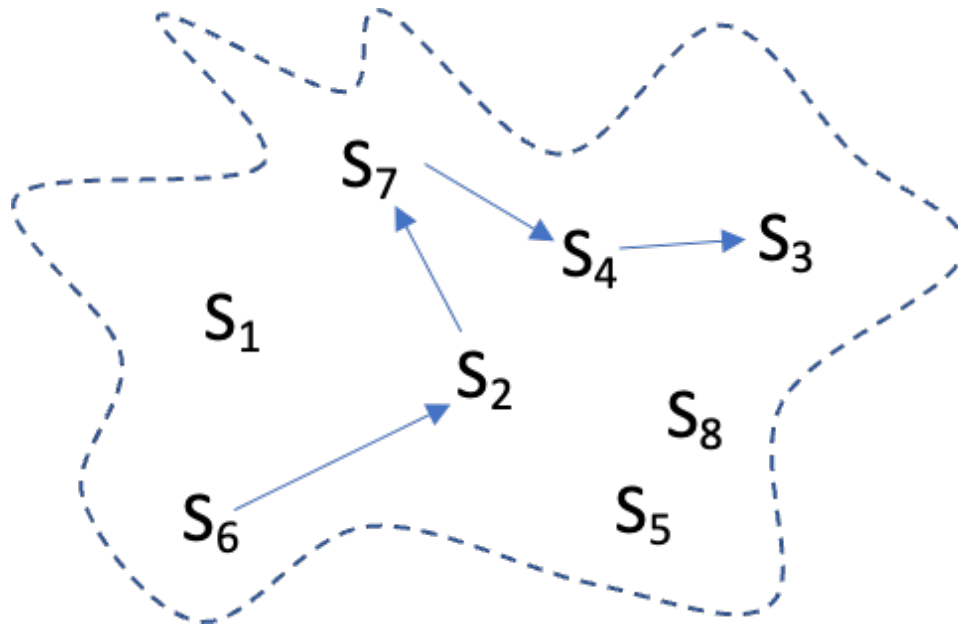
An RL session involves a series of state changes each being the result of an action:

$$s^1 \xrightarrow{a^1} s^2 \xrightarrow{a^2} s^3 \xrightarrow{a^3} s^4 \xrightarrow{a^4} s^5 \xrightarrow{a^5} \dots$$

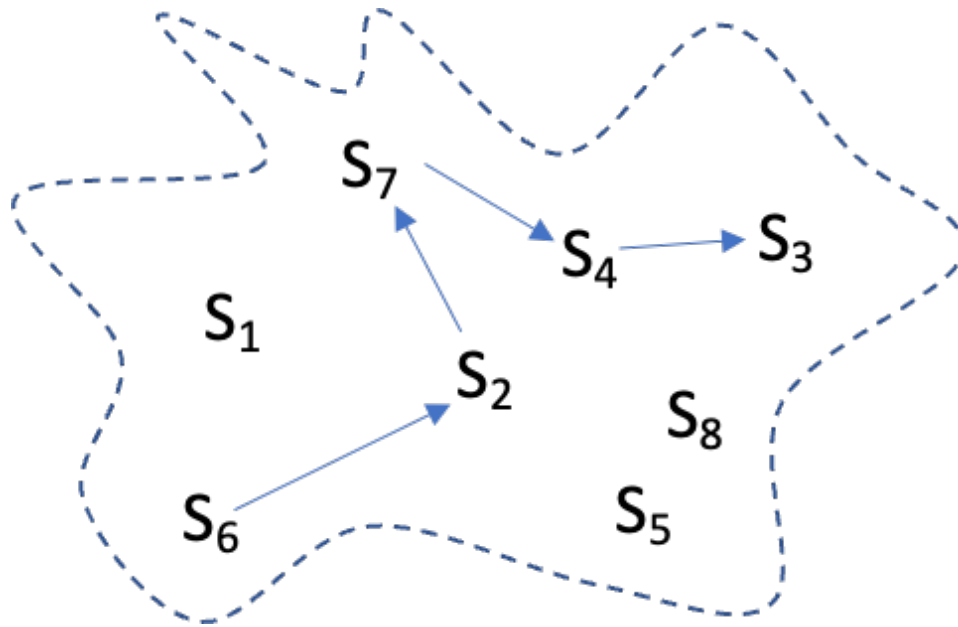
This trails a path in state space:



Typical problem in RL:



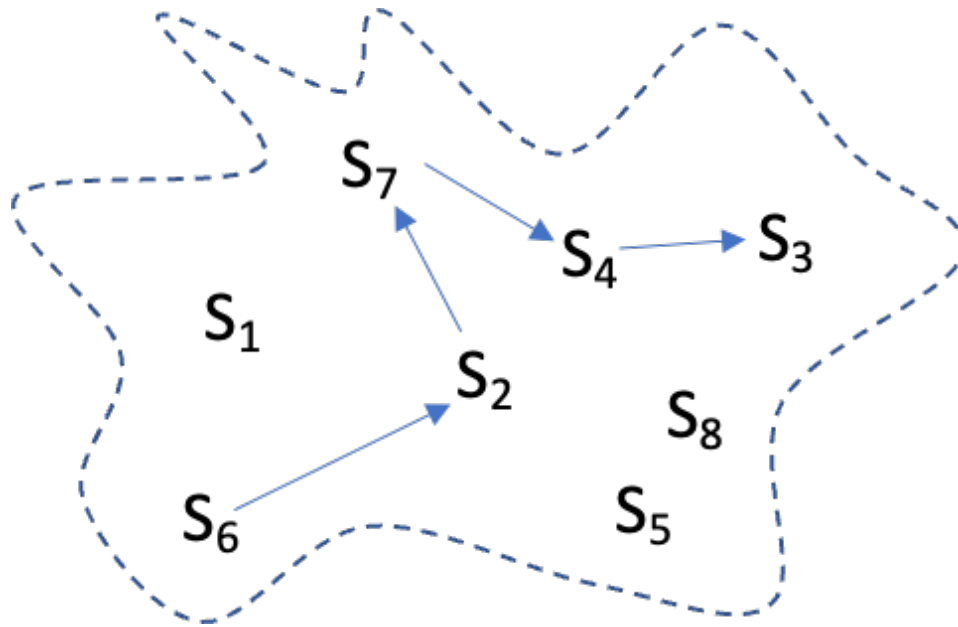
Typical problem in RL:



Goal:

to get from a start state
to a destination state in
the least number of moves

Typical problem in RL:



Goal:

to get from a start state
to a destination state in
the least number of moves

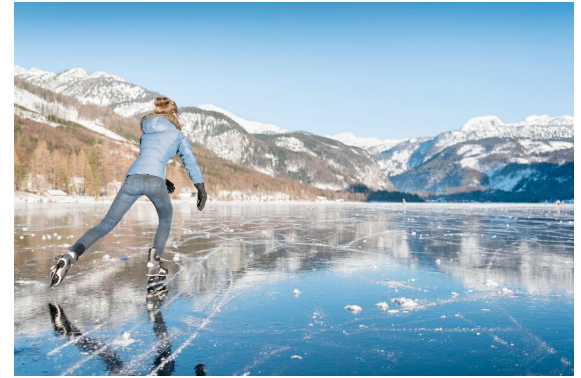
Constraints:

only some state
transitions are allowed

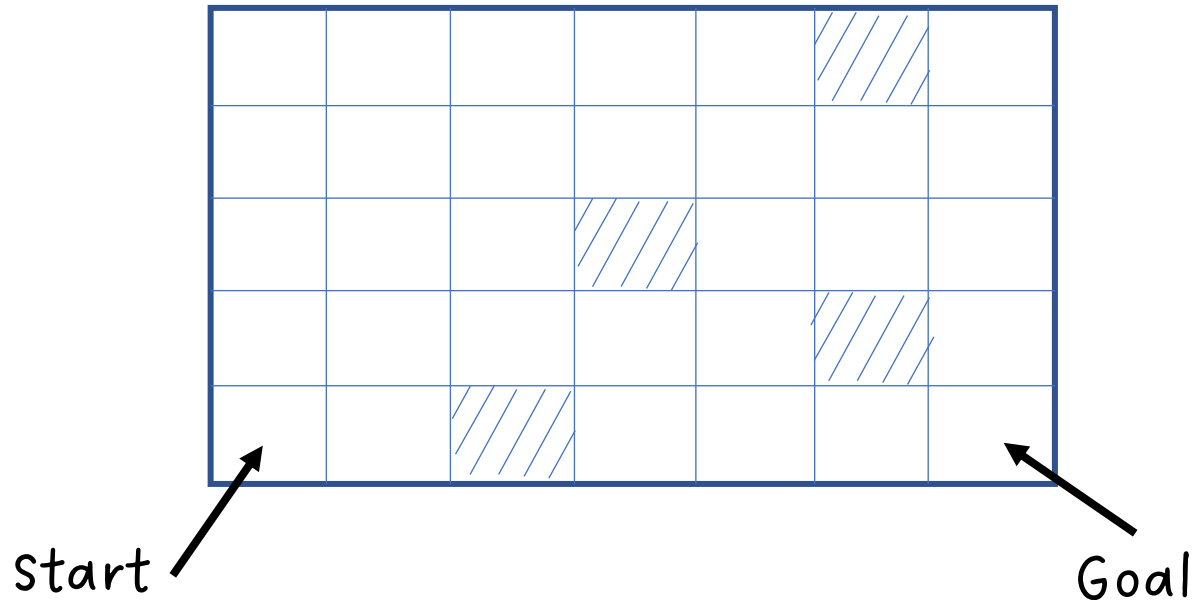
some moves are bad and
some are good

The iced lake

Environment: iced lake with holes



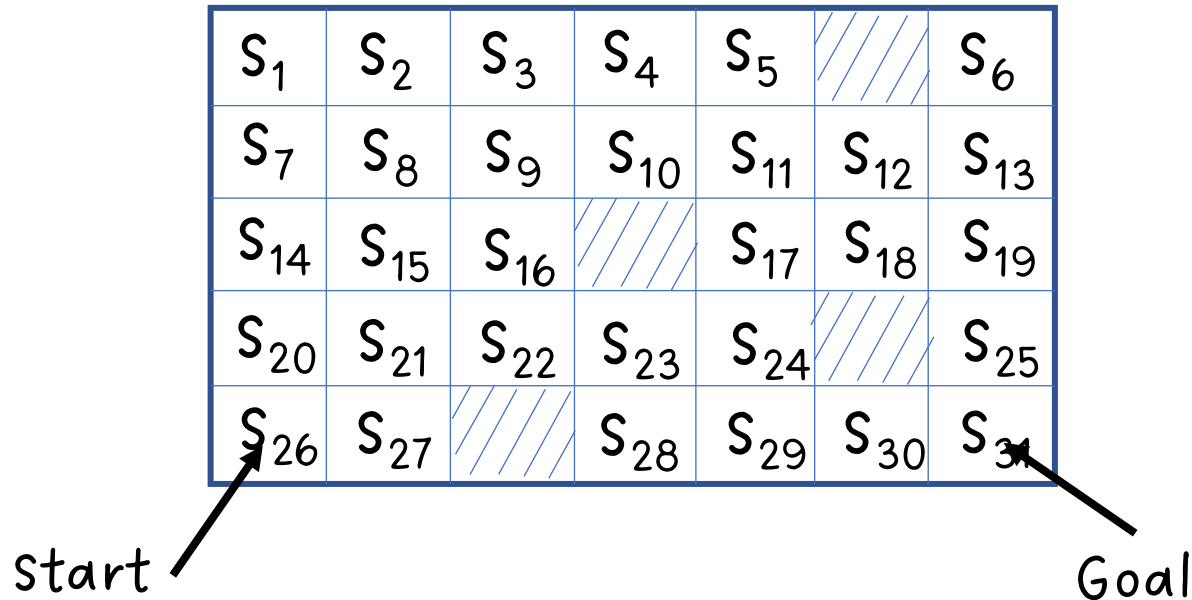
Aim: learn the shortest route from start to goal



The iced lake

The environment is an iced lake with holes

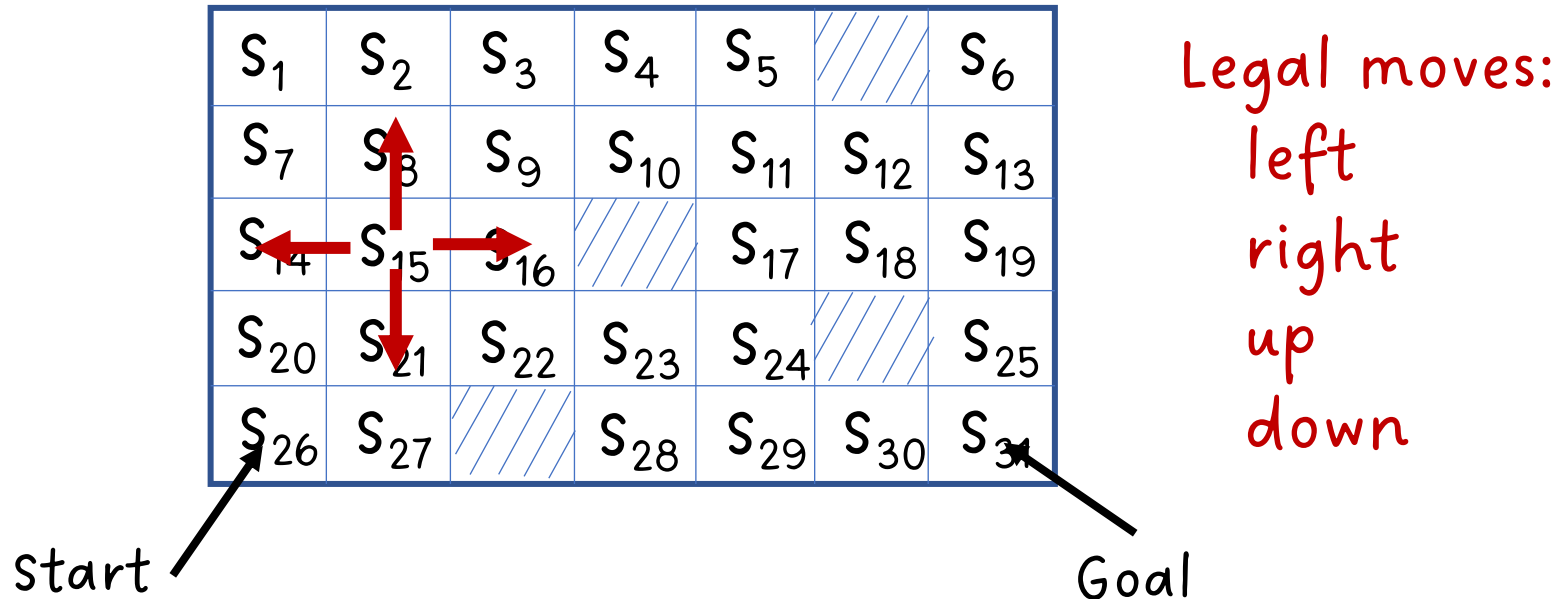
Aim is to learn the shortest route from start to goal



The iced lake

The environment is an iced lake with holes

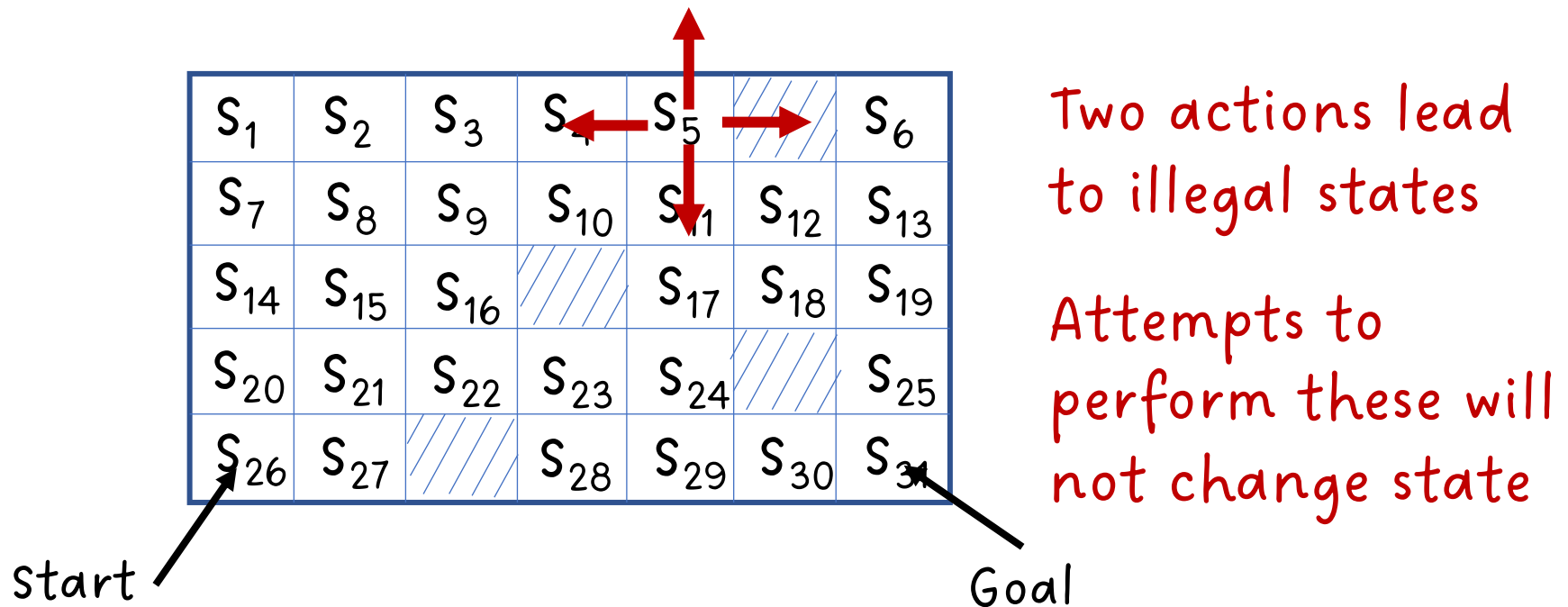
Aim is to learn the shortest route from start to goal



The iced lake

The environment is an iced lake with holes

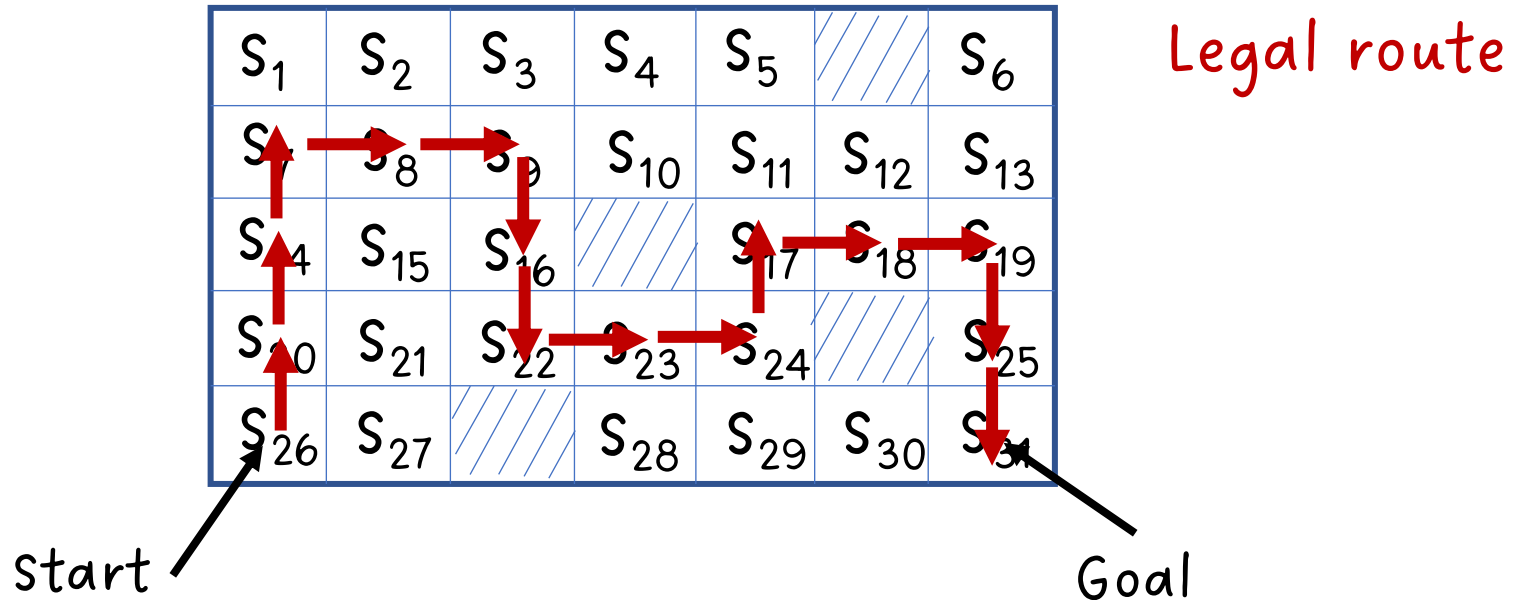
Aim is to learn the shortest route from start to goal



The iced lake

The environment is an iced lake with holes

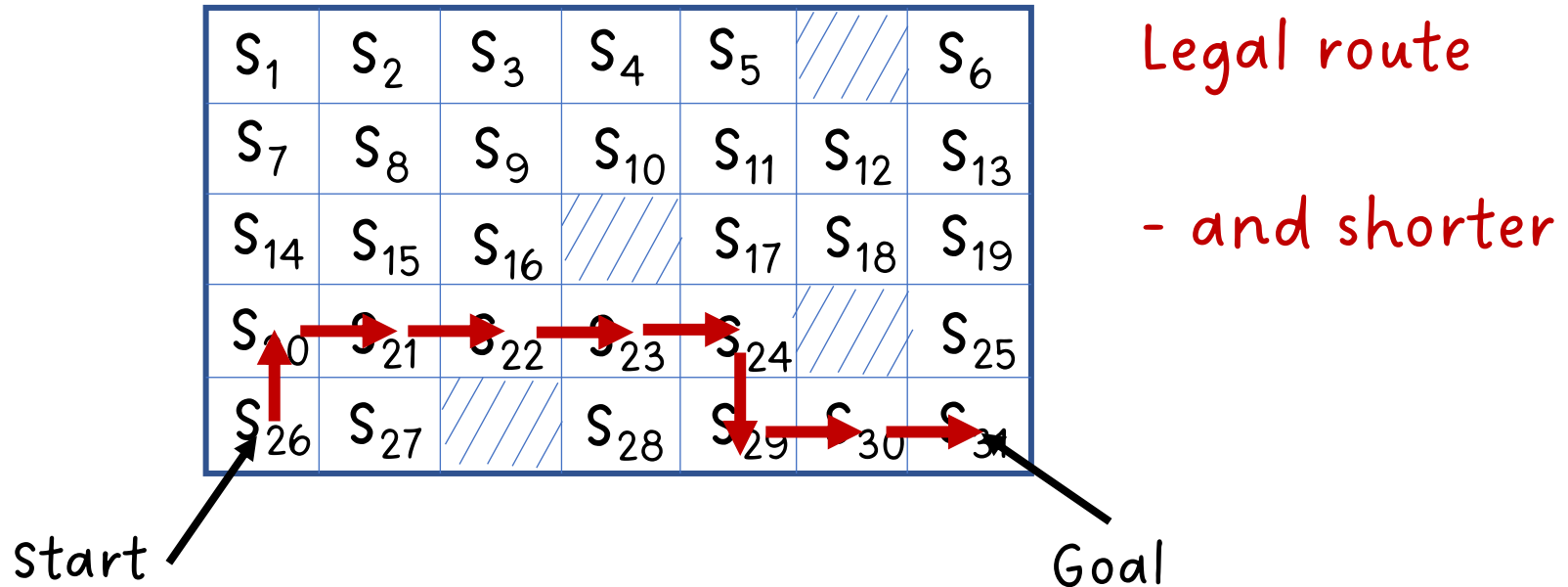
Aim is to learn the shortest route from start to goal



The iced lake

The environment is an iced lake with holes

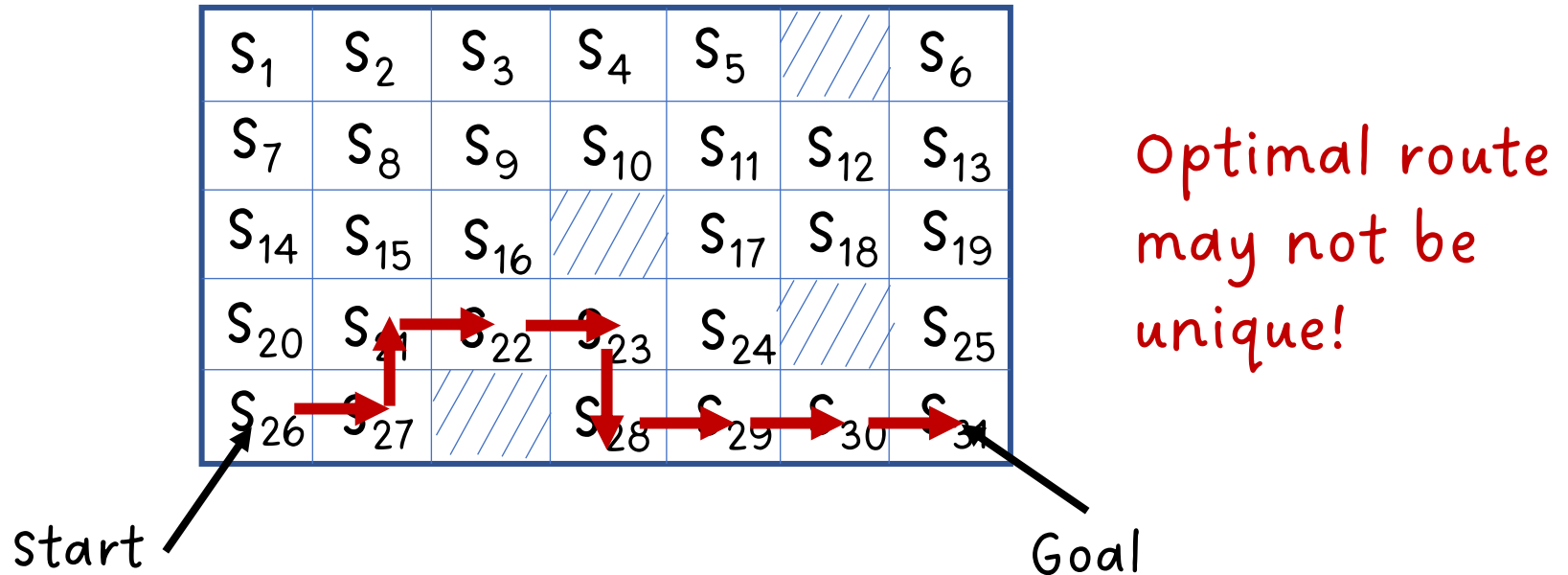
Aim is to learn the shortest route from start to goal



The iced lake

The environment is an iced lake with holes

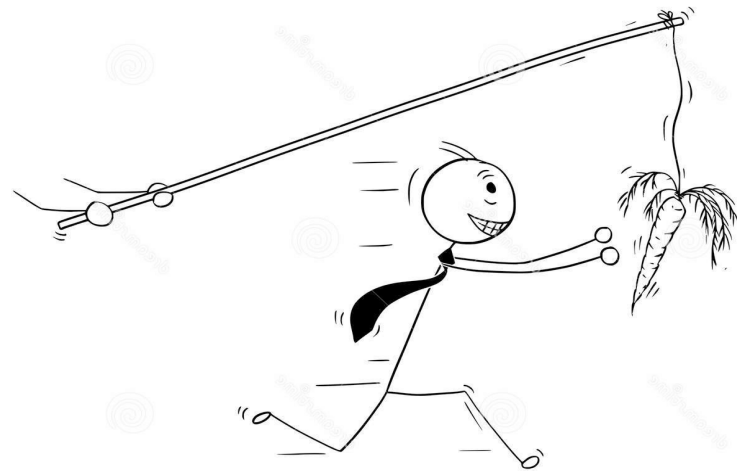
Aim is to learn the shortest route from start to goal



IN3050/IN4050, Lecture 12

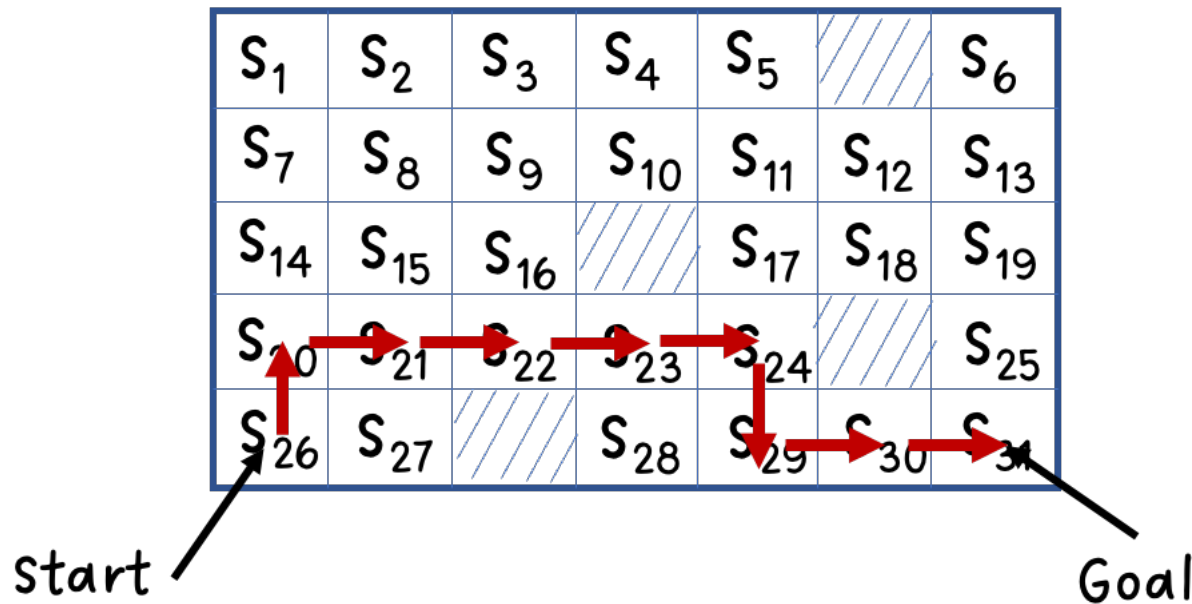
Reinforcement Learning

3: Learning from rewards
Ole Christian Lingjærde



Finding the best action

It may seem easy for a human to find the best action in the iced lake example:



But can we train a computer to do it, too?

Finding the best action

If we could make a computer tackle the iced lake example, we may use same principle to solve much harder problems:

robotics

self-driving cars

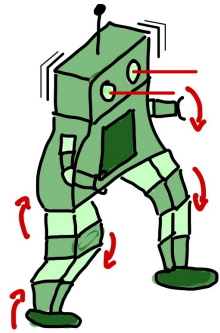
health care

finance

...

This is where reinforcement learning (RL) comes in.

Supervised learning is not practical

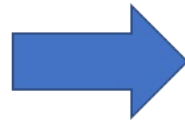
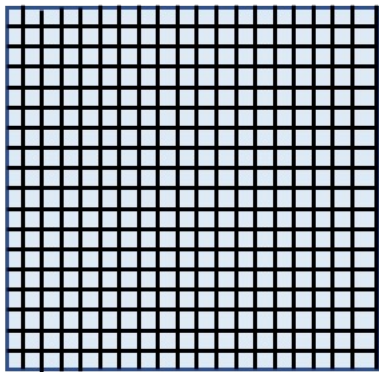


Training set: (x_i, y_i) , $i=1, \dots, n$

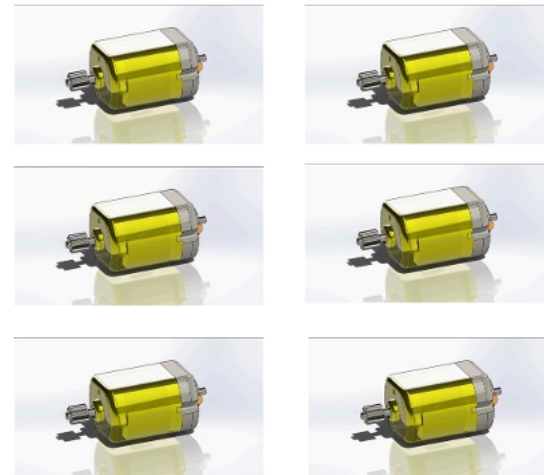
x_i : current sensor readings

y_i : correct/best action

Pixel values, sensor readings, ...




Correct torque for each motor



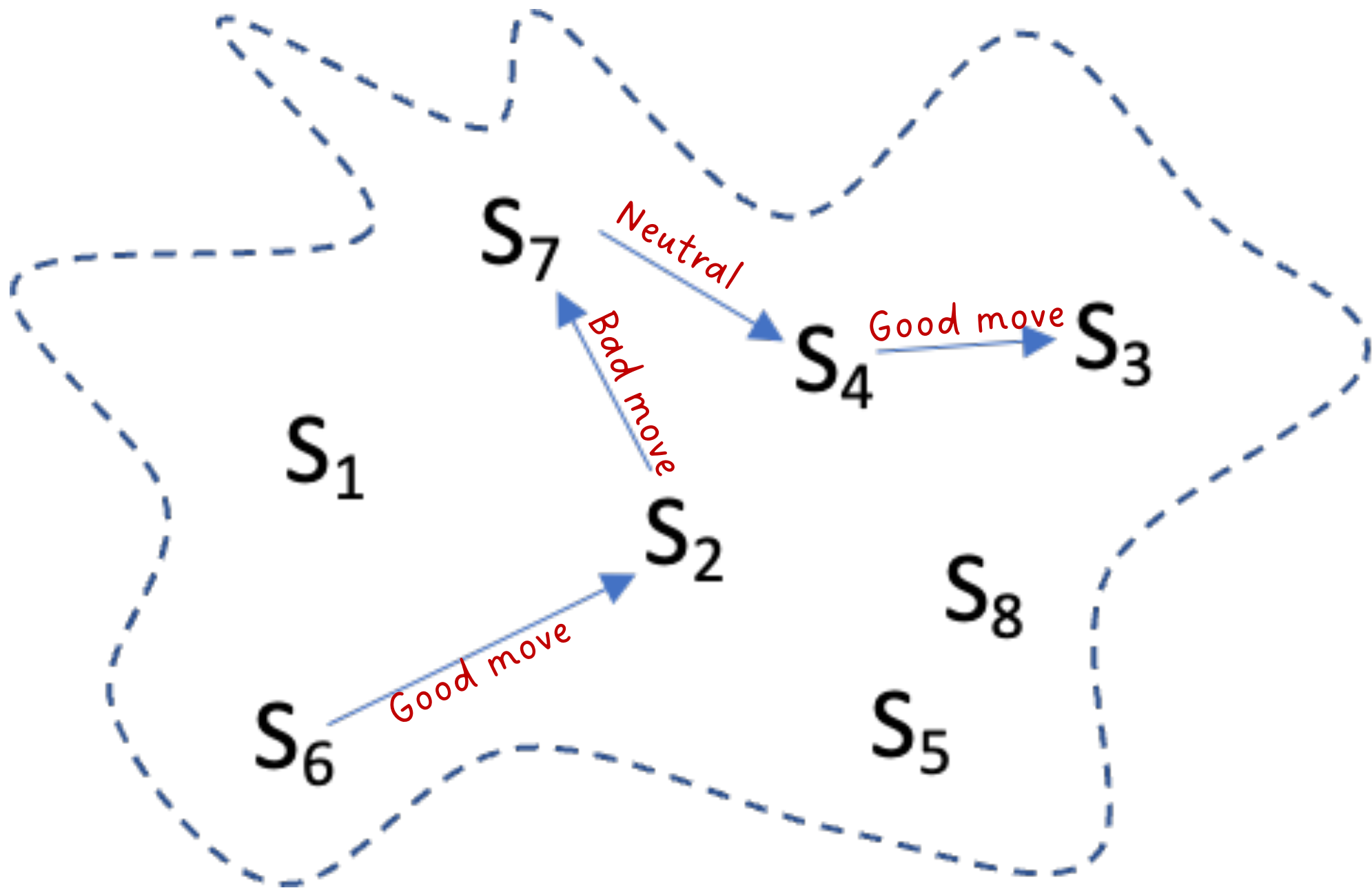
Difficult and extremely time consuming to create such data

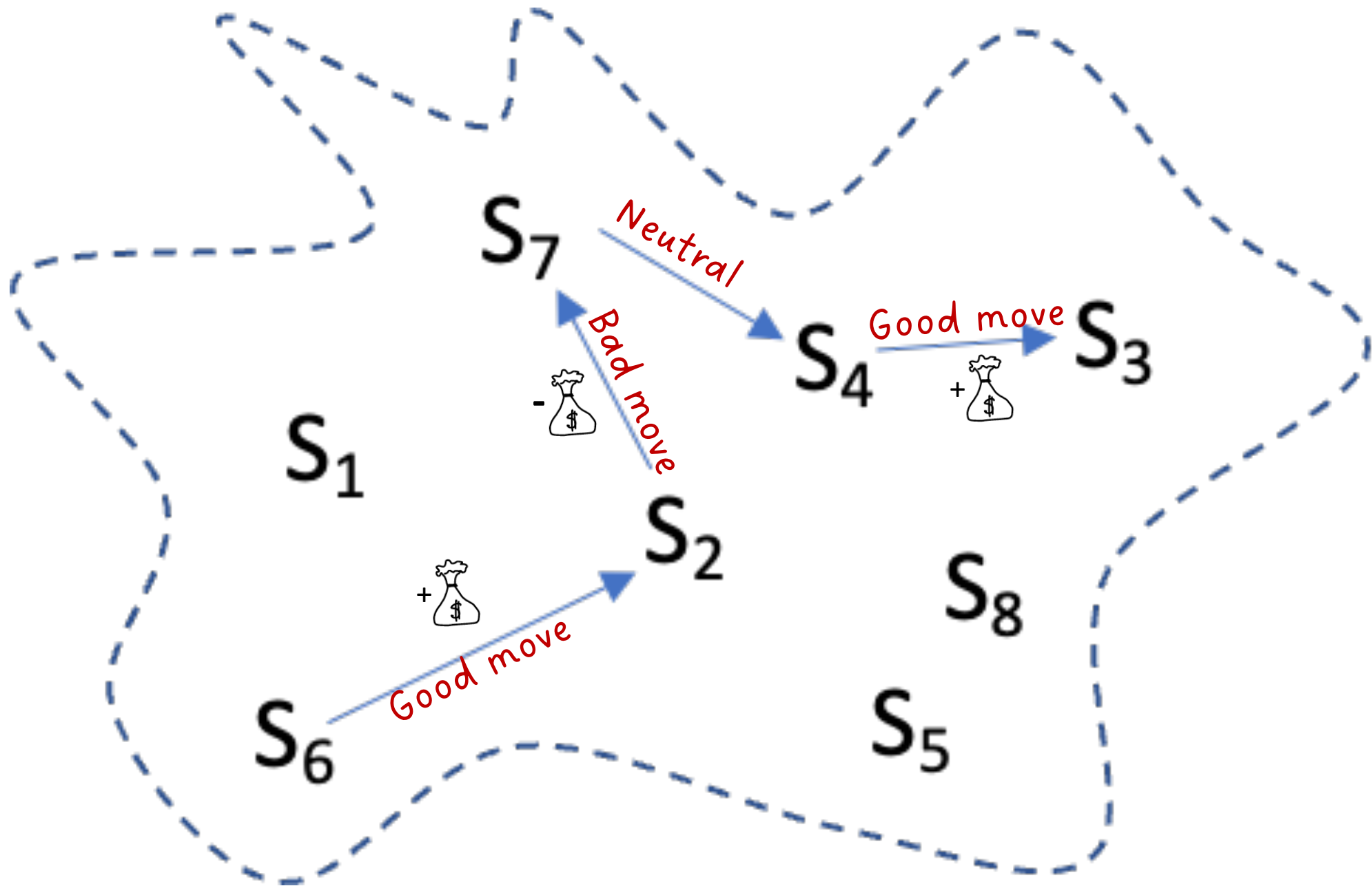
Rewards

Instead of supplying the learner with the correct outputs (which we do not possess) we will tell if the result of an action was **good** or **bad**.

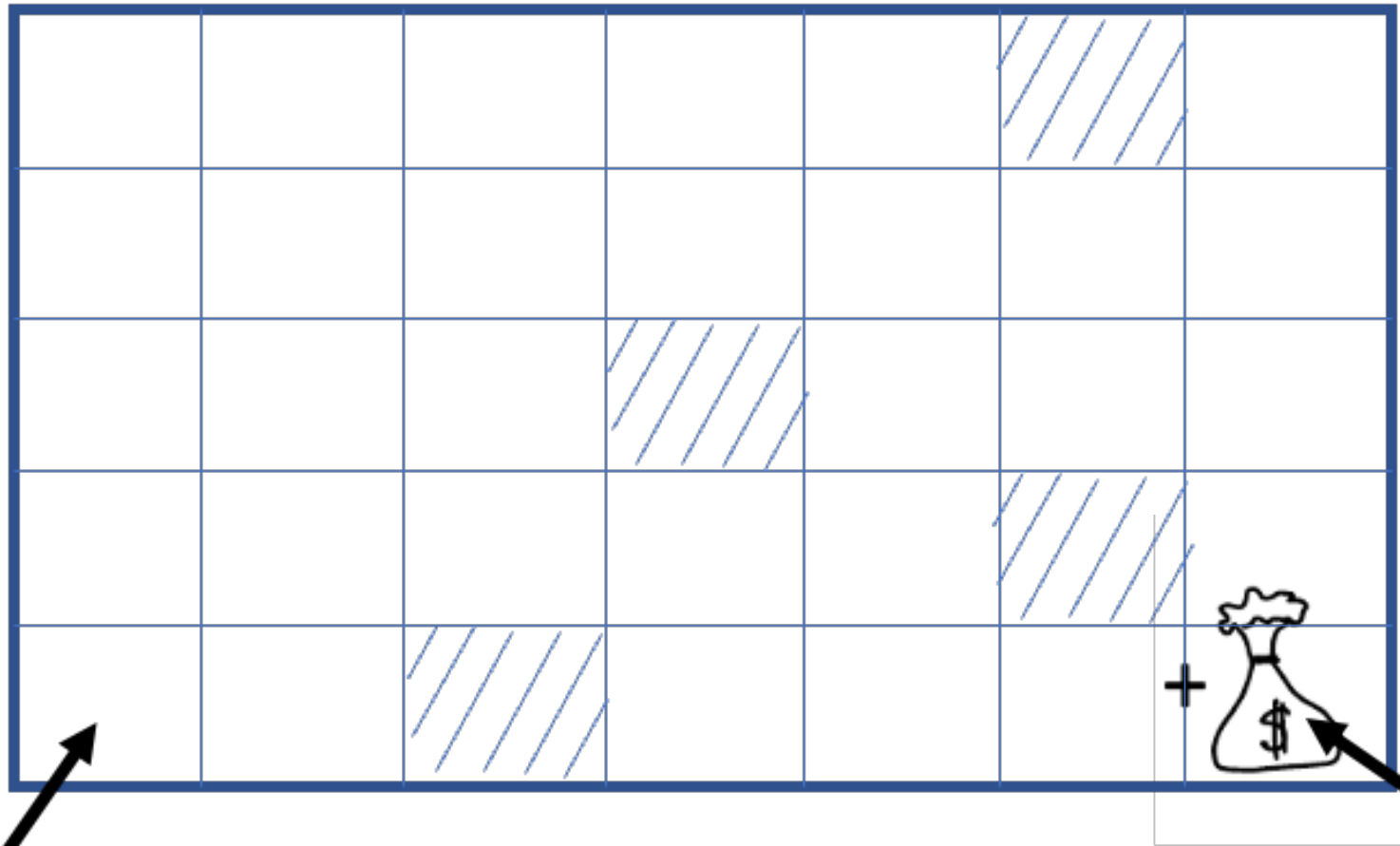
+  = reward

-  = negative reward (penalty)





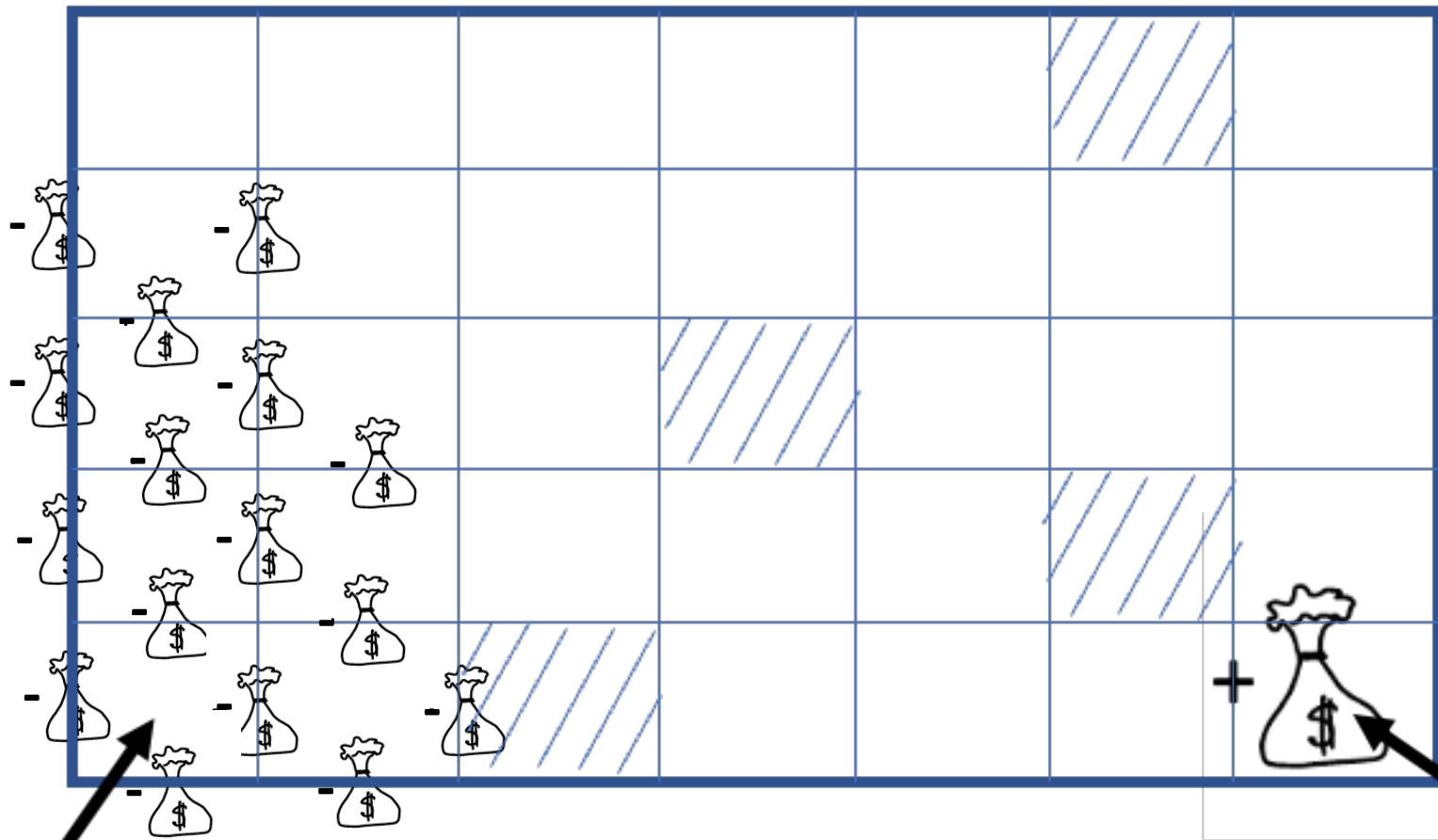
The iced lake



Start

Goal

The iced lake

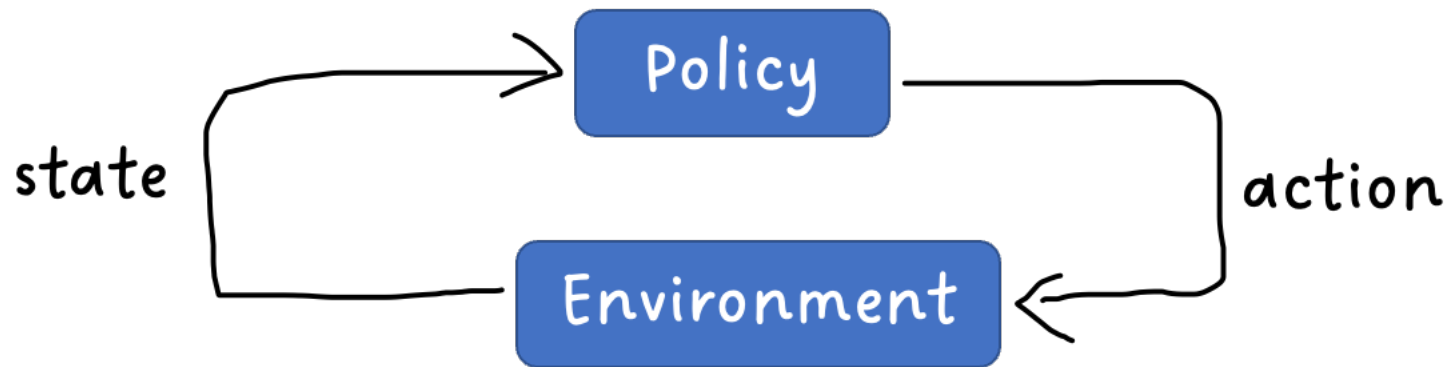


Start

Goal

A static policy means no improvement over time

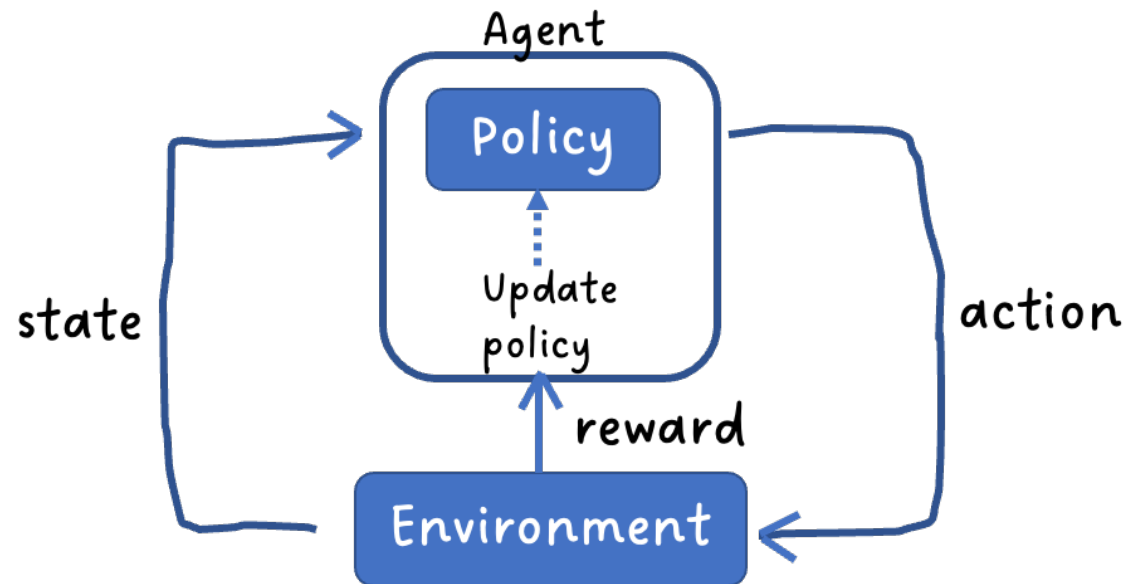
In the static system described earlier in this lecture, the policy was fixed and the system could not improve over time:



Learning means changing the policy

By allowing the policy to change over time, the system can adapt to the environment.

An **agent** repeatedly adjusts the policy based on the previous state, the current state and the reward received:



Reinforcement learning cycle:

```
s[0] = <initial state>  
w = <initial weights>  
for i in range(N):  
    a = action(s[i], w)  
    s[i+1] = newstate(s[i], a)  
    r = reward(s[i], s[i+1], a)  
    <update w>
```


Reinforcement learning cycle:

```
s[0] = <initial state>
w = <initial weights>
for i in range(N):
    a = action(s[i], w)           # POLICY
    s[i+1] = newstate(s[i], a)
    r = reward(s[i], s[i+1], a)
    <update w>
```

Reinforcement learning cycle:

```
s[0] = <initial state>
w = <initial weights>
for i in range(N):
    a = action(s[i], w)          # POLICY
    s[i+1] = newstate(s[i], a)  # ENVIRONMENT
    r = reward(s[i], s[i+1], a)
    <update w>
```

Reinforcement learning cycle:

```
s[0] = <initial state>
w = <initial weights>
for i in range(N):
    a = action(S[i], w)          # POLICY
    S[i+1] = newstate(S[i], a)  # ENVIRONMENT
    r = reward(S[i], S[i+1], a) # ENVIRONMENT
    <update w>
```

Reinforcement learning cycle:

```
s[0] = <initial state>
w = <initial weights>
for i in range(N):
    a = action(S[i], w)           # POLICY
    S[i+1] = newstate(S[i], a)   # ENVIRONMENT
    r = reward(S[i], S[i+1], a)  # ENVIRONMENT
    <update w>                   # POLICY LEARNING
```

Reinforcement learning cycle:

$s[0] = \text{<initial state>}$

$w = \text{<initial weights>}$

for i in range(N):

$a = \text{action}(s[i], w)$

MUST DEFINE

$s[i+1] = \text{newstate}(s[i], a)$

$r = \text{reward}(s[i], s[i+1], a)$

$\text{<update } w\text{>}$

MUST DEFINE

Reinforcement learning is stochastic

Just as in evolutionary learning, there are several ways to introduce randomness into RL:

$S[0] = \text{<initial state>}$

Can be stochastic

$w = \text{<initial weights>}$

Can be stochastic

for i in range(N):

$a = \text{action}(S[i], w)$

Can be stochastic

$S[i+1] = \text{newstate}(S[i], a)$

Can be stochastic

$r = \text{reward}(S[i], S[i+1], a)$

Can be stochastic

 <update w >

Can be stochastic

Reinforcement learning involves exploration

Reinforcement learning receives less guidance from trainer than supervised learning.

Exploration is essential:

- explore different actions
- see where they lead
- accumulate knowledge



As training progresses, exploration can be replaced by exploitation.

IN3050/IN4050, Lecture 12

Reinforcement Learning

4: Defining the policy
Ole Christian Lingjærde



Recall the reinforcement learning cycle:

$s[0] = \text{<initial state>}$

$w = \text{<initial weights>}$

for i in range(N):

$a = \text{action}(s[i], w)$

$s[i+1] = \text{newstate}(s[i], a)$

$r = \text{reward}(s[i], s[i+1], a)$

 <update w >

Recall the reinforcement learning cycle:

$s[0] = \text{<initial state>}$

$w = \text{<initial weights>}$

Defines the current policy

for i in range(N):

Adjustable through the parameter vector w

$a = \text{action}(s[i], w)$

$s[i+1] = \text{newstate}(s[i], a)$

$r = \text{reward}(s[i], s[i+1], a)$

$\text{<update } w\text{>}$

How do we design this function?

The policy

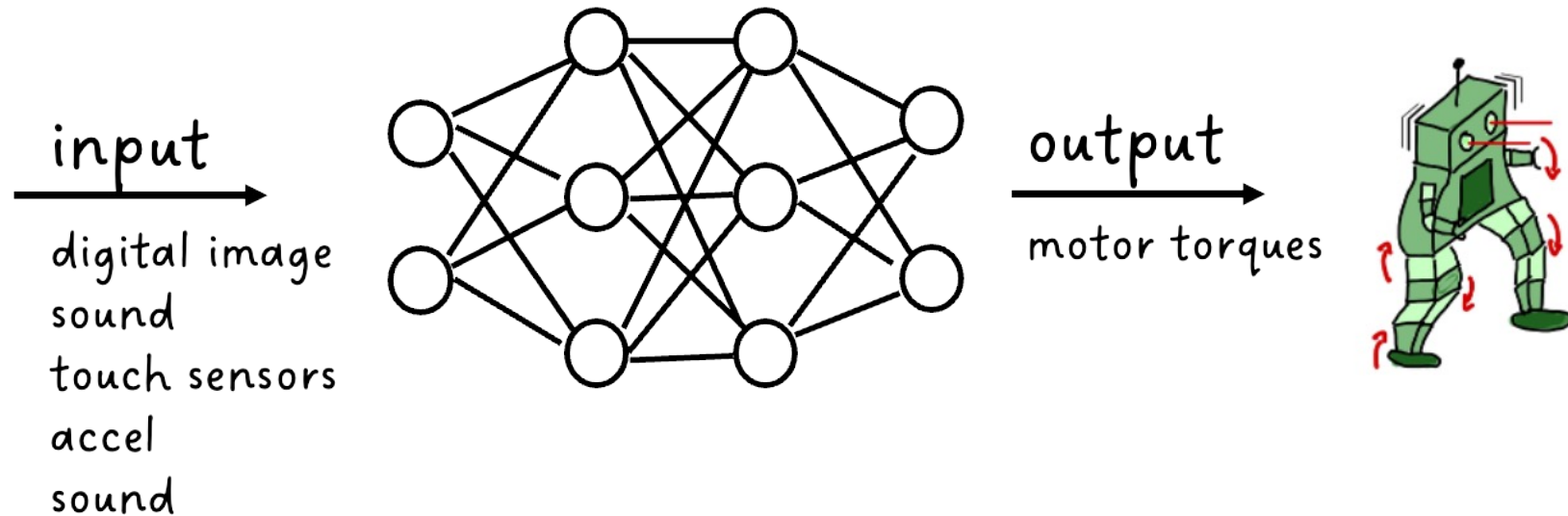
Decides what action(S, w) to take in a given state S , using prior knowledge collected in a parameter vector w .

Two common approaches:

- Represent policy as a neural network
- Represent policy as a table



Representing the policy as a neural network



The weights in the network are given by the parameter vector w in $\text{action}(S, w)$

Representing the policy as a table

An alternative to neural networks and the idea is to store all previous learning experience in a **table** with one entry for each (state, action) combination:

	a_1	a_2	a_3	a_m
s_1					
s_2					
.					
.					
s_m					

Q-table

average reward obtained after taking this action when in this state

Illustration:

	a_1	a_2	a_3
s_1	5	-1	3
s_2	0	2	11
s_3	7	9	-5
s_m	8	2	0

→ In state s_3 the highest average reward is obtained when taking action a_2

$S = 0, 1, 2, \dots$ is a state and Q is a numpy table.

```
def action(S, Q):
```

```
    return np.where(Q[S,:] == max(Q[S,:]))[0][0]
```

Three different ways of using the Q-table

greedy: $\text{action} = \underset{a}{\operatorname{argmax}} Q[s,a]$

ϵ -greedy: $\text{action} = \begin{cases} \underset{a}{\operatorname{argmax}} Q[s,a], & \text{with prob } 1 - \epsilon \\ \text{random action}, & \text{with prob } \epsilon \end{cases}$

softmax: $\text{action} = \begin{cases} a_1 & \text{with prob } \frac{e^{Q[s,a_1]/\tau}}{\sum e^{Q[s,a_k]/\tau}} \\ \dots & \dots \\ a_m & \text{with prob } \frac{e^{Q[s,a_m]/\tau}}{\sum e^{Q[s,a_k]/\tau}} \end{cases}$

An example

Suppose we are in state s_2 .
How do we choose the next action?



	a_1	a_2	a_3
s_1	5	-1	3
s_2	0	2	11
s_3	7	9	-5
s_m	8	2	0

greedy: choose a_3

ϵ -greedy: Draw a random number $X \sim U[0,1]$
if $X \geq \epsilon$: choose a_3
else: select action at random

softmax: Draw $(X_1, X_2, X_3) \sim \text{multinomial}(1; p_1, p_2, p_3)$
if $X_1 = 1$: choose a_1
else if $X_2 = 1$: choose a_2
else if $X_3 = 1$: choose a_3

$$p_1 = \frac{\exp(0)}{\exp(0) + \exp(2) + \exp(11)}$$

$$p_2 = \frac{\exp(2)}{\exp(0) + \exp(2) + \exp(11)}$$

$$p_3 = \frac{\exp(11)}{\exp(0) + \exp(2) + \exp(11)}$$

IN3050/IN4050, Lecture 12

Reinforcement Learning

5: Learning the policy
Ole Christian Lingjærde



Main reinforcement learning cycle:

$s[0] = \text{<initial state>}$

$w = \text{<initial weights>}$

for i in range(N):

$a = \text{action}(s[i], w)$

$s[i+1] = \text{newstate}(s[i], a)$

$r = \text{reward}(s[i], s[i+1], a)$

 <update w >

Main reinforcement learning cycle:

$s[0] = \text{<initial state>}$

$w = \text{<weights from a previous learning cycle (episode)>}$

for i in range(N):

$a = \text{action}(s[i], w)$

$s[i+1] = \text{newstate}(s[i], a)$

$r = \text{reward}(s[i], s[i+1], a)$

 <update w >

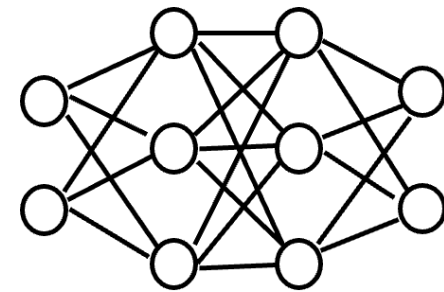
Learning (= updating w)

Neural network implementation of policy:

w = network weights

policy-gradient methods

actor-critic methods



Q-table implementation of policy:

w = the table Q

SARSA

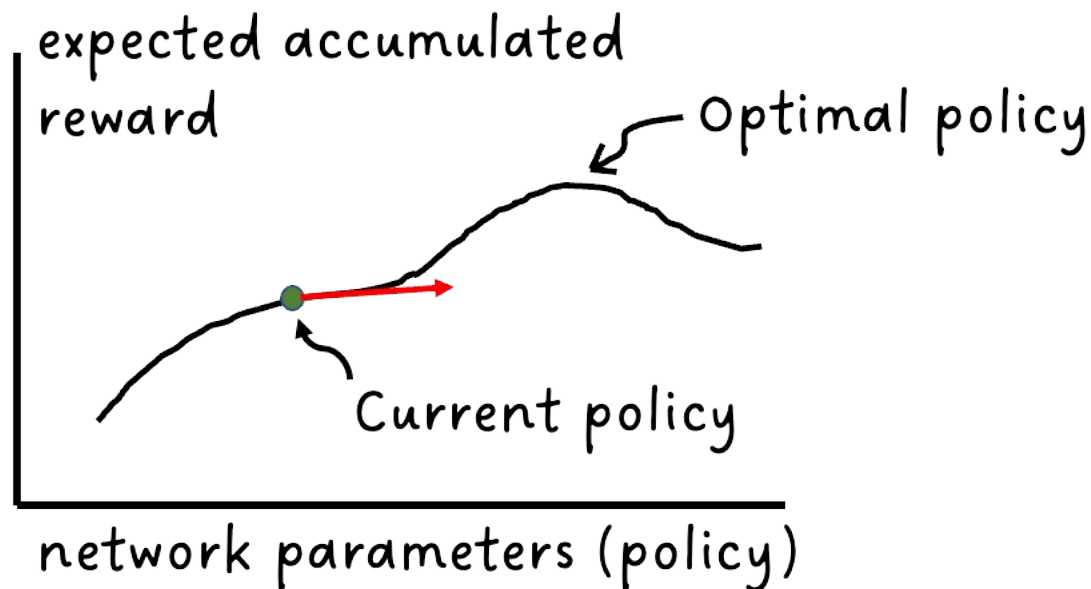
Q-learning

	a_1	a_2	a_3	a_m
s_1					
s_2					
.					
.					
s_m					

An arrow points to the cell at the intersection of row s_2 and column a_m .

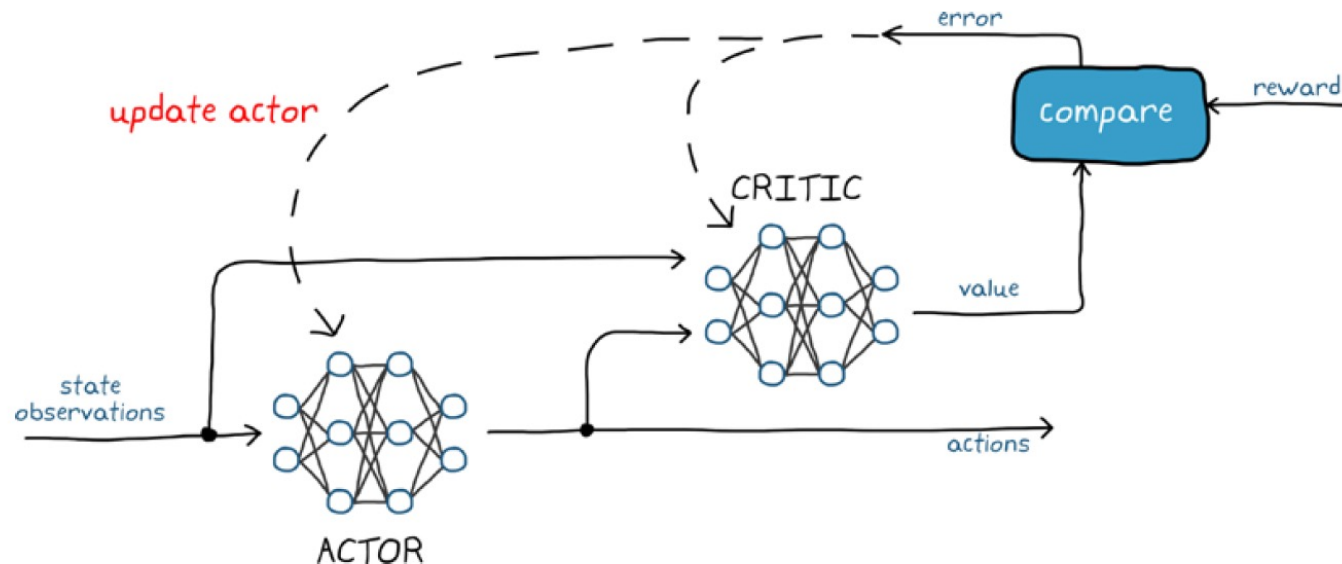
Policy-gradient methods

- Execute current policy
- Collect rewards
- After several iterations adjust the weights in the direction of increased expected reward



Actor-critic methods

- Two neural networks interact and are trained together
- **The critic network** learns the values of different (state, action) pairs from the received rewards
- **The actor network** learns correct actions using feedback from the critic network



Learning algorithms for Q-tables

Store all previous learning experience as a table.

	a_1	a_2	a_3	a_m
s_1					
s_2					
.					
.					
s_m					

average reward
obtained after
taking this action
when in this state

Suppose we have experienced (s_i, a_j) a total of k times and have received rewards

$$r_1, r_2, \dots, r_k$$

Seems reasonable to let

$$Q[s_i, a_j] = (r_1 + r_2 + \dots + r_k) / k$$

NOTE: If the environment is completely deterministic then $\text{newstate}(\dots)$ and $\text{reward}(\dots)$ are deterministic, too.

Then $r_1 = r_2 = \dots = r_k$.

Problem 1

The environment may change over time and in that case we would trust recent rewards more than older ones

Solution:

Suppose we are in state s , current policy says take action a , and the reward is r . Then we let

$$Q[s,a] = (1-\mu)Q[s,a] + \mu r$$

Updates (first to last):

$$Q[s,a] = \mu r_1$$

$$Q[s,a] = (1-\mu)\mu r_1 + \mu r_2$$

$$Q[s,a] = (1-\mu)^2\mu r_1 + (1-\mu)\mu r_2 + \mu r_3$$

Problem 2

The update rule only considers the reward for the next move. We want $Q[s,a]$ to reflect the total accumulated reward from s and all the way to the final destination.


Solution:

Suppose we are in state s , policy says take action a , the reward is r , the new state is s' , and policy says we should then take action a' .

$$Q[s,a] = (1-\mu)Q[s,a] + \mu (r + \gamma Q[s',a'])$$

This learning rule is called SARSA.

Discount rate - a number $\gamma \in (0,1)$ close to 1



SARSA actually takes into account all future steps!

Suppose the (state,action) sequence is $(s,a), (s',a'), (s'',a''), \dots$

Updates in (s,a) :

$$Q[s,a] = (1-\mu)Q[s,a] + \mu [r + \gamma Q[s',a']]$$

Updates in (s',a') :

$$Q[s',a'] = (1-\mu)Q[s',a'] + \mu [r' + \gamma Q[s'',a'']]$$

Thus (assuming $\mu=1$ for simplicity):

$$Q[s,a] = r + \gamma r' + \gamma^2 r'' + \gamma^3 r''' + \dots$$

(called the discounted future reward when $\gamma < 1$)

Discounted future reward

The expression

$$Q[s,a] = r + \gamma r' + \gamma^2 r'' + \gamma^3 r''' + \dots$$

tells us that rewards received later on the path to the destination are reduced (discounted). Why?



Reason: there is some uncertainty in each action we take and that uncertainty accumulates along the path.

Q-learning

This is an alternative to SARSA.

In Q-learning we replace the term $Q[s', a']$ in SARSA with the maximal value of $Q[s', :]$

$$Q[s, a] = (1 - \mu)Q[s, a] + \mu (r + \gamma \max Q[s', :])$$

If the policy used is the greedy algorithm then Q-learning and SARSA are identical.

The Q-Learning Algorithm

- Initialisation
 - set $Q(s, a)$ to small random values for all s and a
 - Repeat:
 - initialise s
 - repeat:
 - * select action a using ϵ -greedy or another policy
 - * take action a and receive reward r
 - * sample new state s'
 - * update $Q(s, a) \leftarrow Q(s, a) + \mu(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$
 - * set $s \leftarrow s'$
 - For each step of the current episode
 - Until there are no more episodes
-

The Q-Learning Algorithm

- Initialisation

- set $Q(s, a)$ to small random values for all s and a

- Repeat:

- initialise s

- repeat:

- * select action a using ϵ -greedy or another policy

- * take action a and receive reward r

- * sample new state s'

- * update $Q(s, a) \leftarrow Q(s, a) + \mu(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$

- * set $s \leftarrow s'$

- For each step of the current episode

- Until there are no more episodes

$$Q(s_t, a_t) \leftarrow \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\mu}_{\text{learning rate}} \cdot \left(\underbrace{r_{t+1}}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} - \underbrace{Q(s_t, a_t)}_{\text{old value}} \right)$$