

Reinforcement Learning

Peerapon S.

Machine Learning (1/67)

Topics

Examples and problem formulation

Policy and value function

Off-line learning

Q-learning

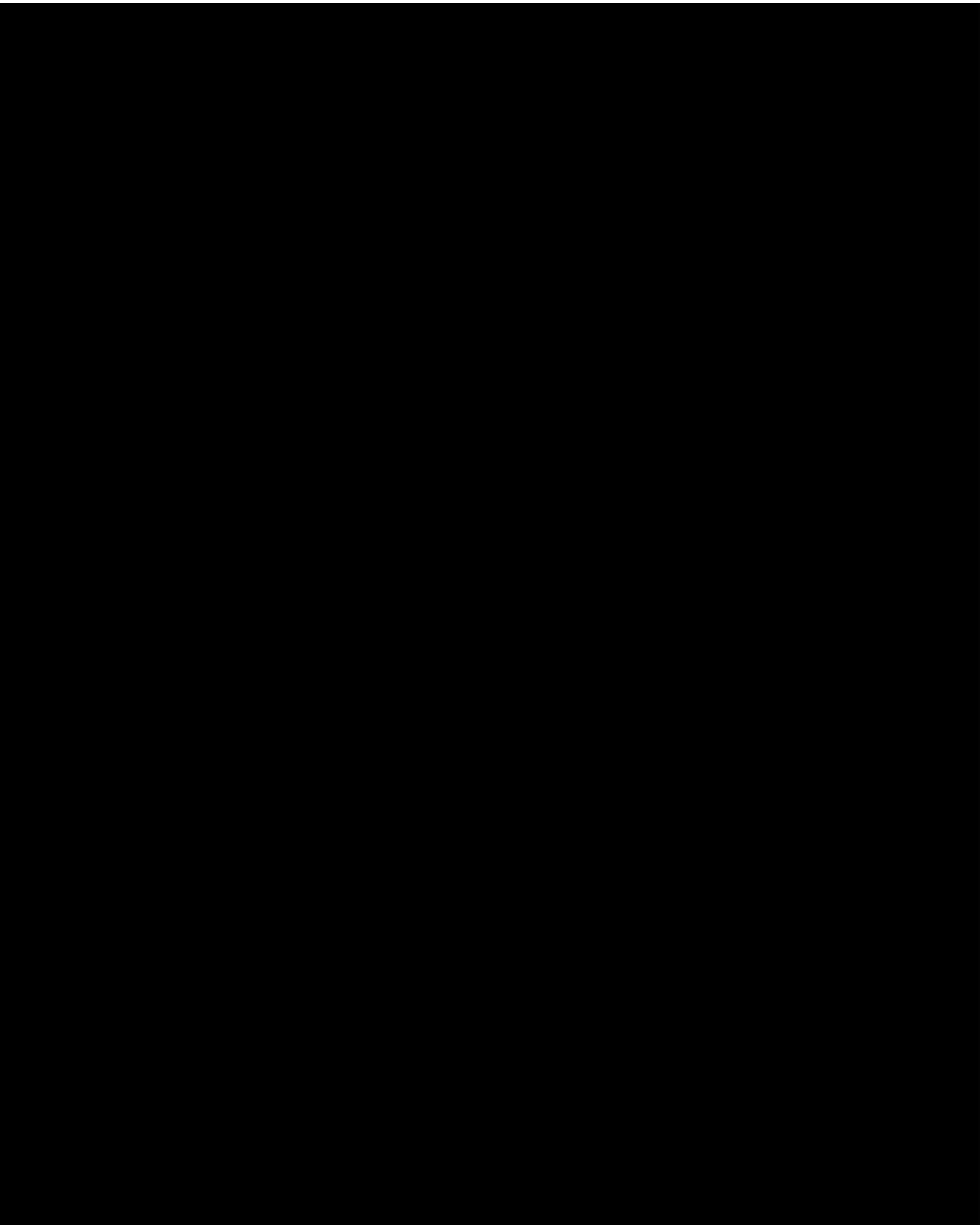
Temporal difference learning

ANOTHER VIEW OF MACHINE LEARNING

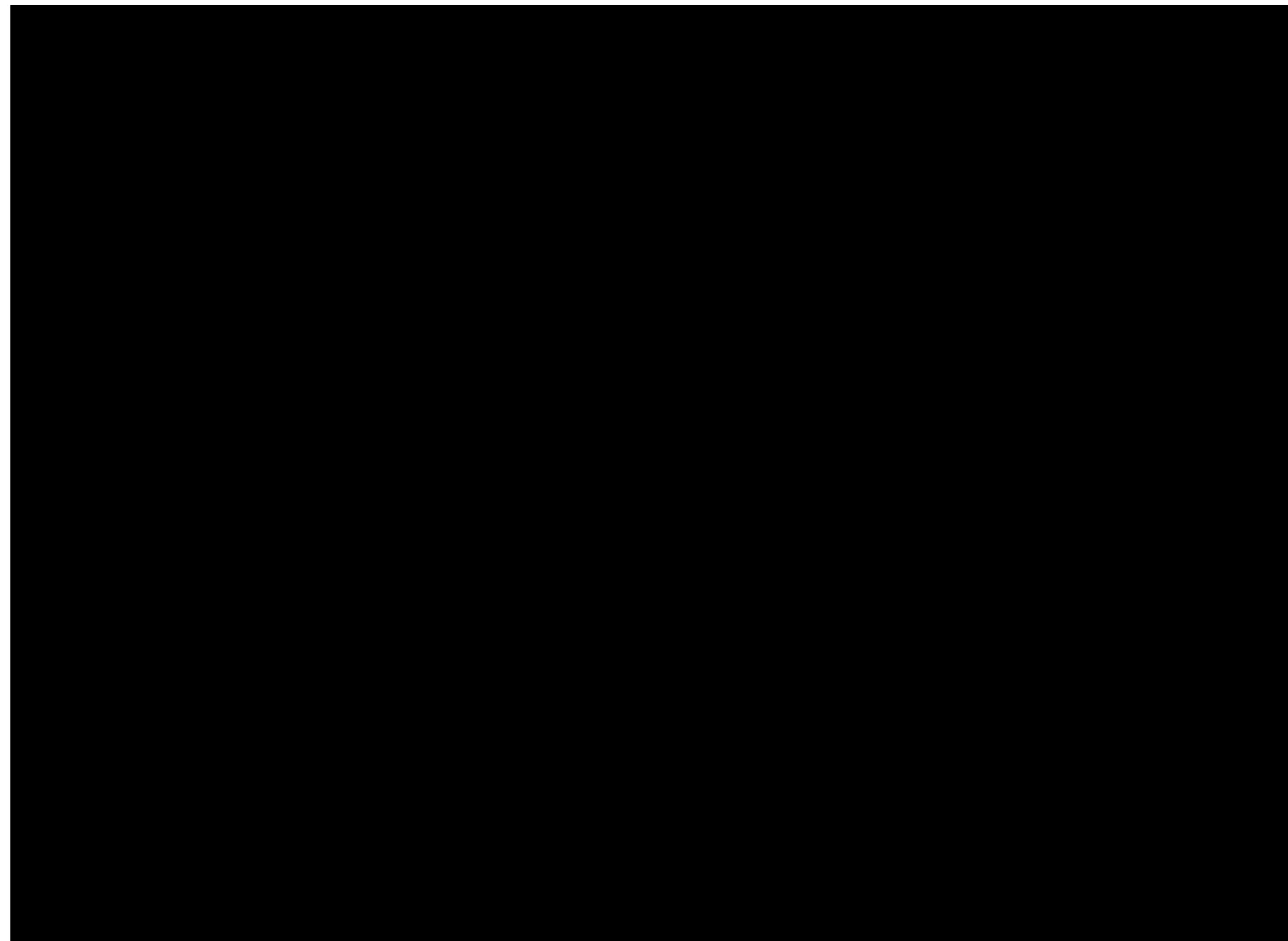


TEACHING THE COMPUTER TO LEARN FROM
EXPERIENCES AND OPTIMIZE A GIVEN
PERFORMANCE INDEX AS THEY PRACTICE.

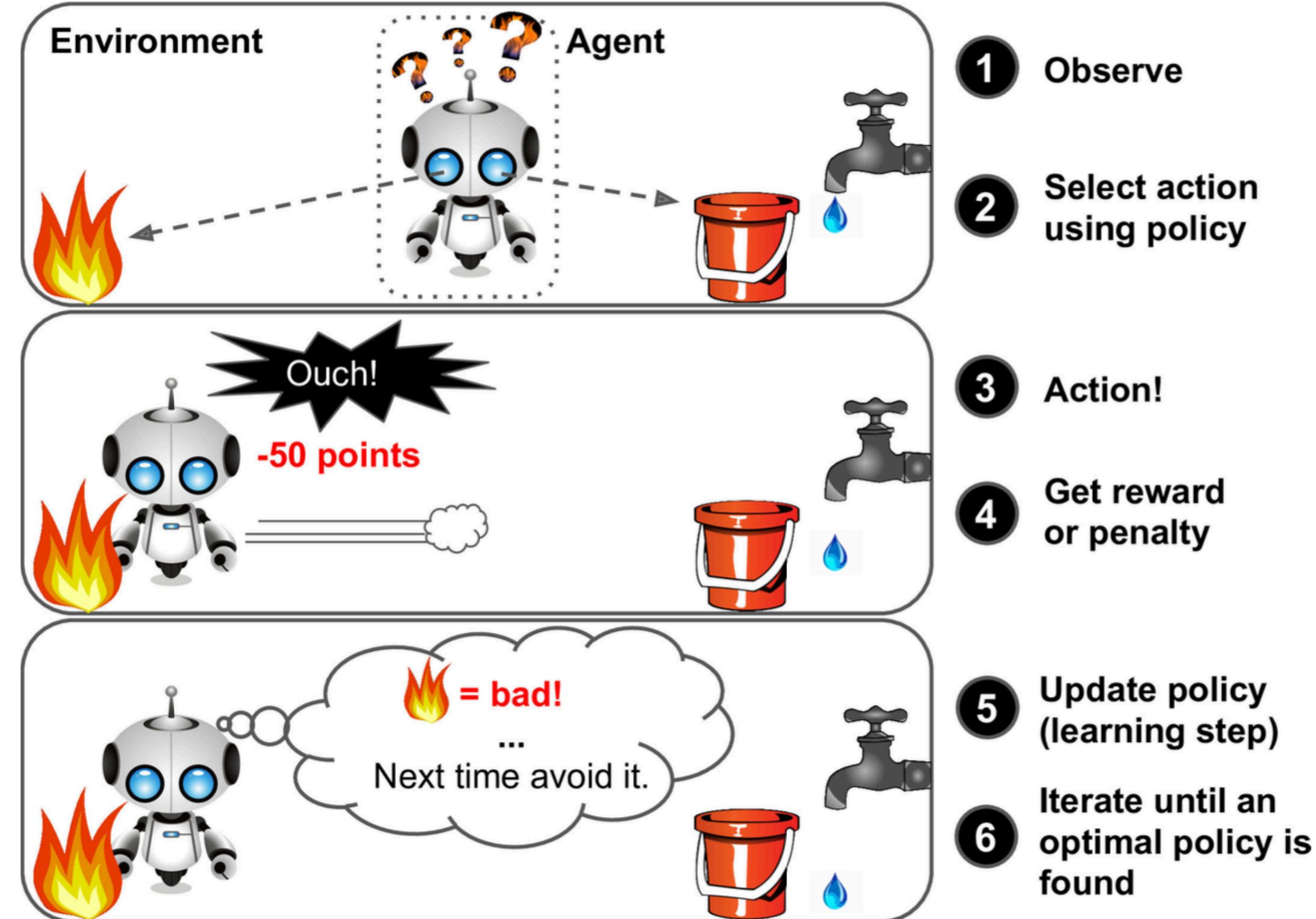
ATARI Game



Running Bot



Basic Concept



Problem Formulation

Agent (Decision maker) interacts with *environment* to complete a task.

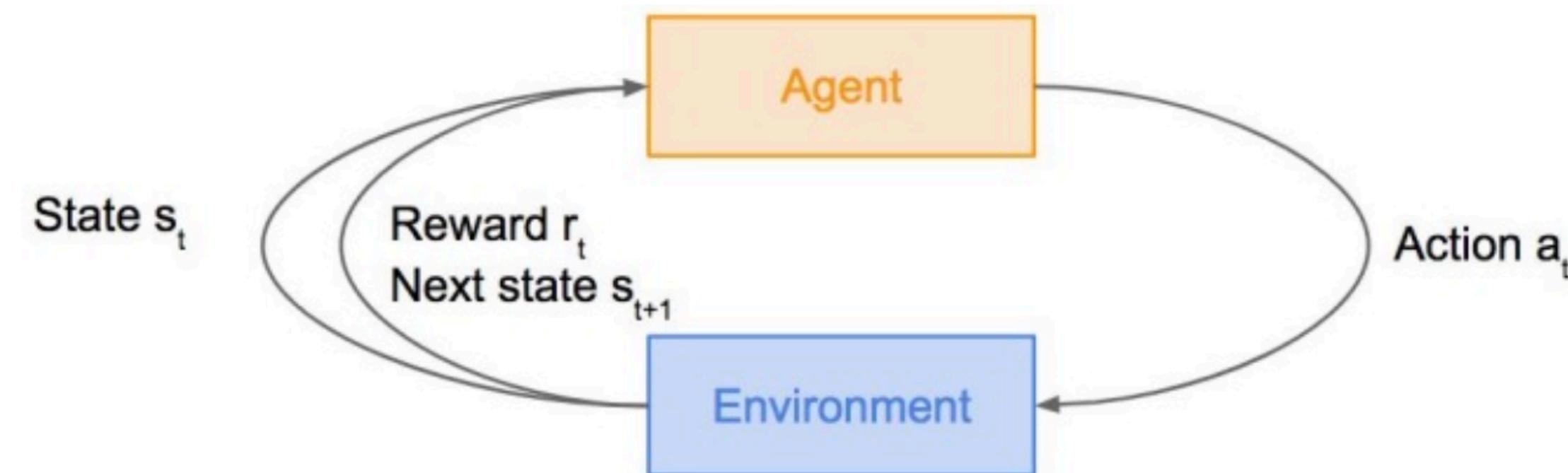
At any point in time, the agent

- ◆ Observes the current *state* $s_t \in S$ of the environment.
- ◆ Selects an *action* $a_t \in A$ to perform, resulting in change of state to s_{t+1} .
- ◆ Receives an immediate *reward* (feedback) r_t from the environment.

Each attempt at the task is referred to as an *episode* or *trial*.

Goal: Complete the task as successfully as possible -- maximizing the *cumulative reward across an episode*.

$$(s_1, a_1, r_1), (s_2, a_2, r_2), (s_3, a_3, r_3), \dots, (s_t, a_t, r_t), \dots$$



The *reward function* $r(\cdot)$ and *state transition function* $\delta(\cdot)$ are parts of the environment, referred to as *model of the world/environment*.

Model possibilities:

- ◆ Known/Not known to the agent
- ◆ Independent or dependent on past states and actions
- ◆ Deterministic or Non-deterministic (Stochastic)

Ex: Atari's Arcade Game

State = Positions of remaining raw pixels

Action = Game controls - Left, Right

Reward = Score increase/decrease at each time step

Task = Complete the game with least number of actions.



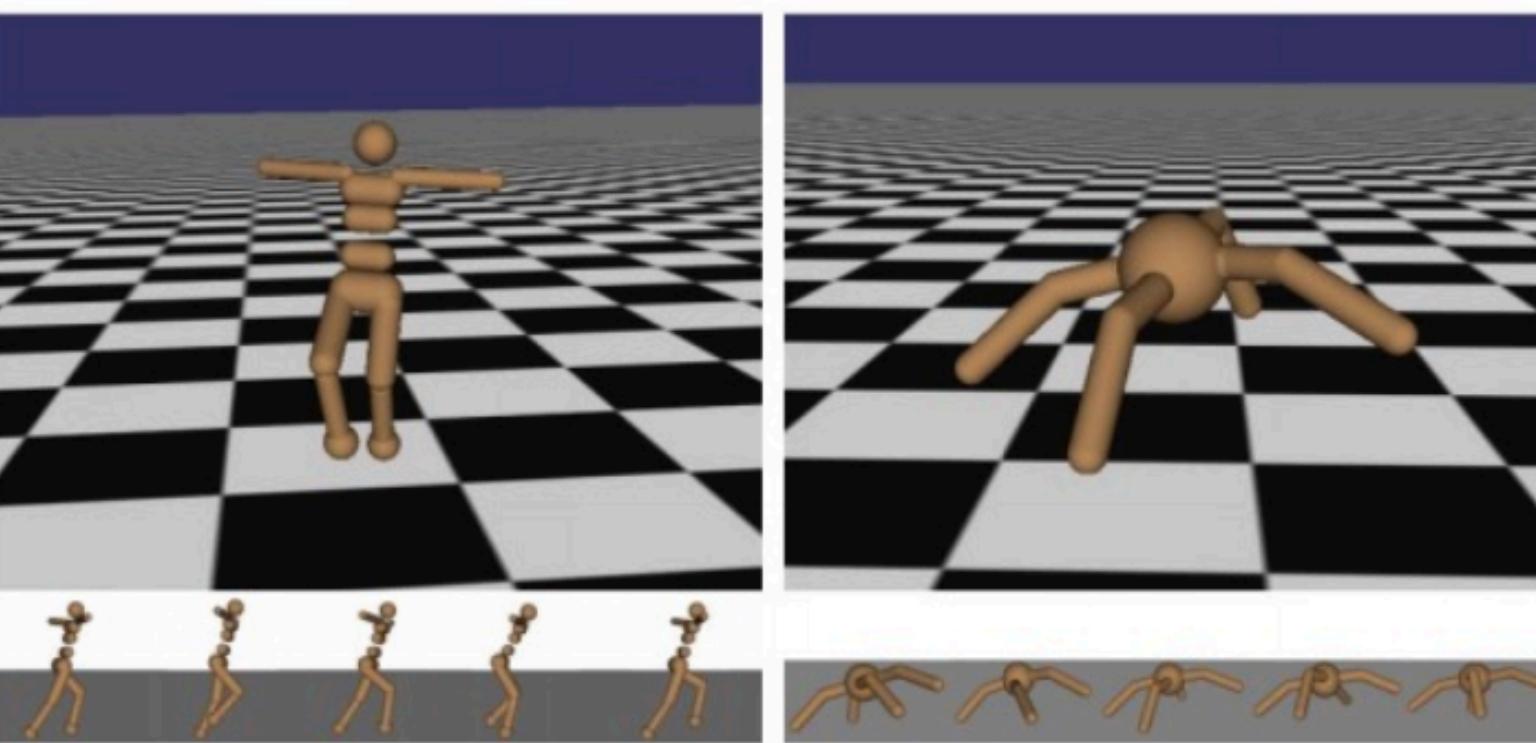
Ex: Robot Locomotion

State = Leg/body positions

Action = Leg/body movement

Reward = +1 if not falling down, 0 otherwise

Task = Make the robot move forward without falling down to the destination.

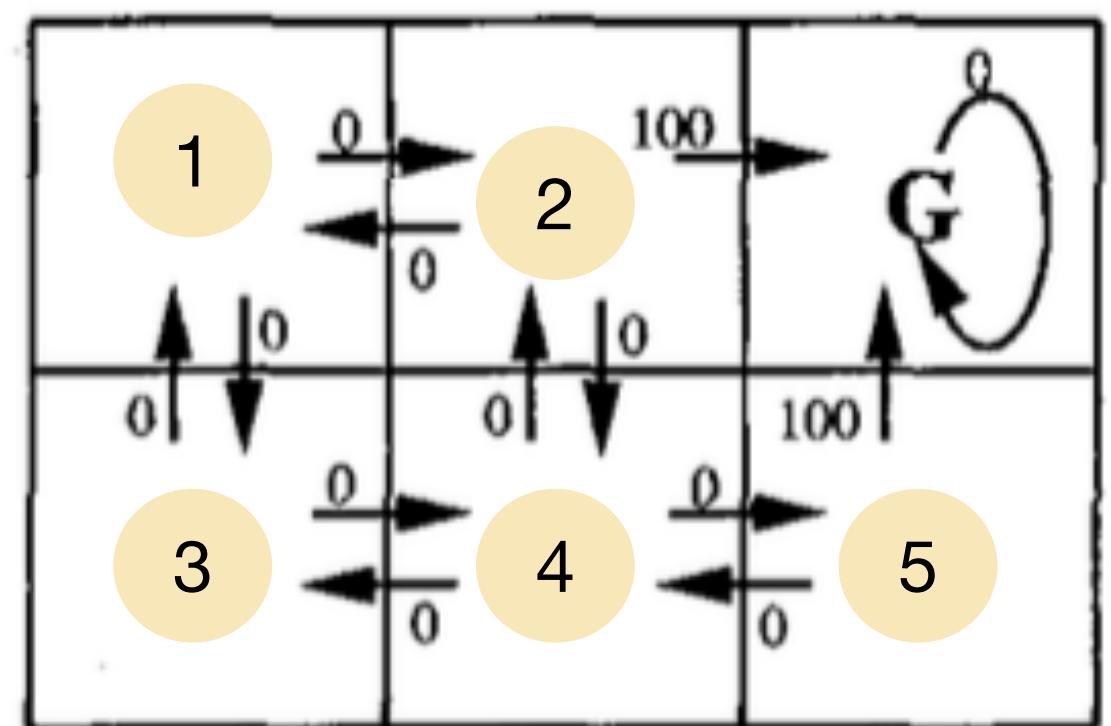


Ex: Six-grid-squares Environment

Assume $r(\cdot)$ and $\delta(\cdot)$ being deterministic and independent from past states and actions so that $r_t = r(s_t, a_t)$ and $s_{t+1} = \delta(s_t, a_t)$.

- ◆ What is the set of states ? $\{1, 2, 3, 4, 5, G\}$
- ◆ What is the set of actions ? $\{\text{Left}, \text{Right}, \text{Top}, \text{Down}, \text{Self}\}$
- ◆ What is the reward ? $\{0, 100\}$

Task = Go the cell G (Goal/Terminal/Absorbing state) with least number of moves from any initial cell.

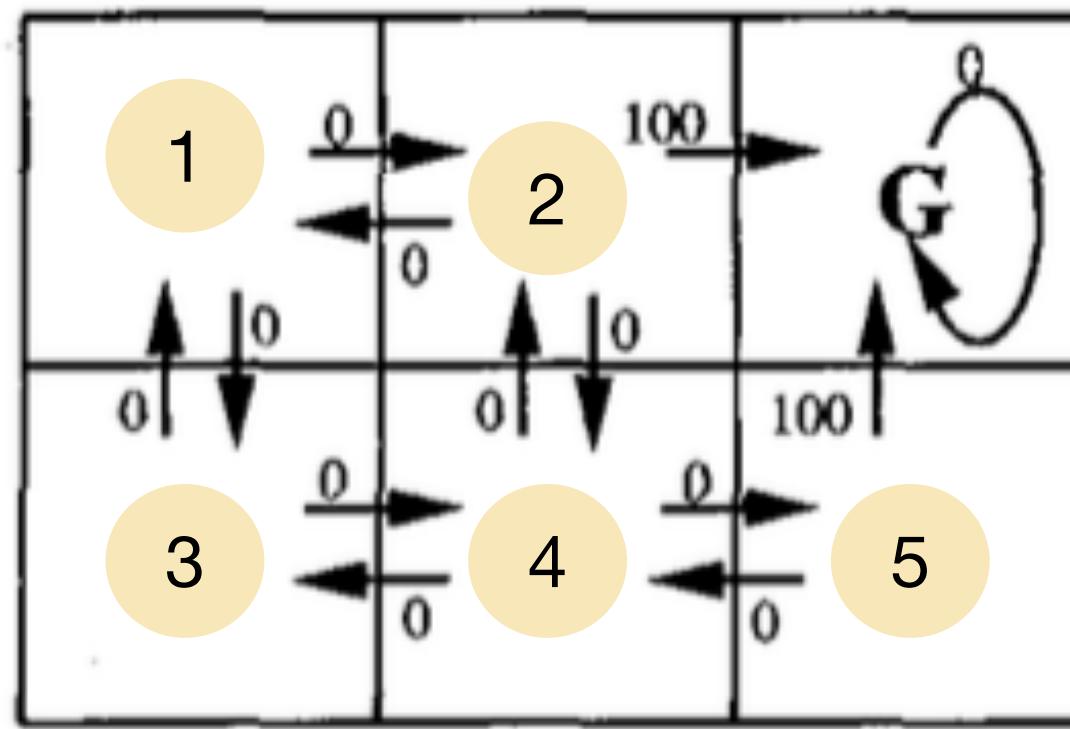


Immediate reward values $r(a, s)$

Policy

A policy (function) $\pi : S \rightarrow A$ is a *mapping* from $s_t \in S$ to $a_t \in A$, i.e., what action $a_t = \pi(s_t)$ to take in each state s_t .

Given $\pi(s)$, the action in state s is completely determined by s .



Ex: *Greedy action selection policy*

- Always take the action yielding the highest immediate reward.
- Fail to handle different characteristics of rewards, e.g., *delayed* or *contradictory*.

Need a different measure to decide an action in a given state.

Value Function – Discounted Reward

Defined as a *cumulative reward* expected to earn when the agent is in state s_t and follows the policy π .

For $0 \leq \gamma < 1$, the *discounted (cumulative) reward* obtained from using the policy π is defined by

discount factor

$$V^\pi(s_t) \triangleq r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots$$

$$= r_t + \gamma (r_{t+1} + \gamma r_{t+2} + \dots)$$

$$= r_t + \gamma V^\pi(s_{t+1})$$

$V^\pi(s)$ has the recurrence relation with that of next state s' .

What is the implication of discount factor γ ?

discount factor ស្ថិតិមានអំពីរក្សាបន្ទូរ future reward

Optimal Policy

The goal of RL is the agent learns the optimal policy -- Maximizing the cumulative reward in the long run.

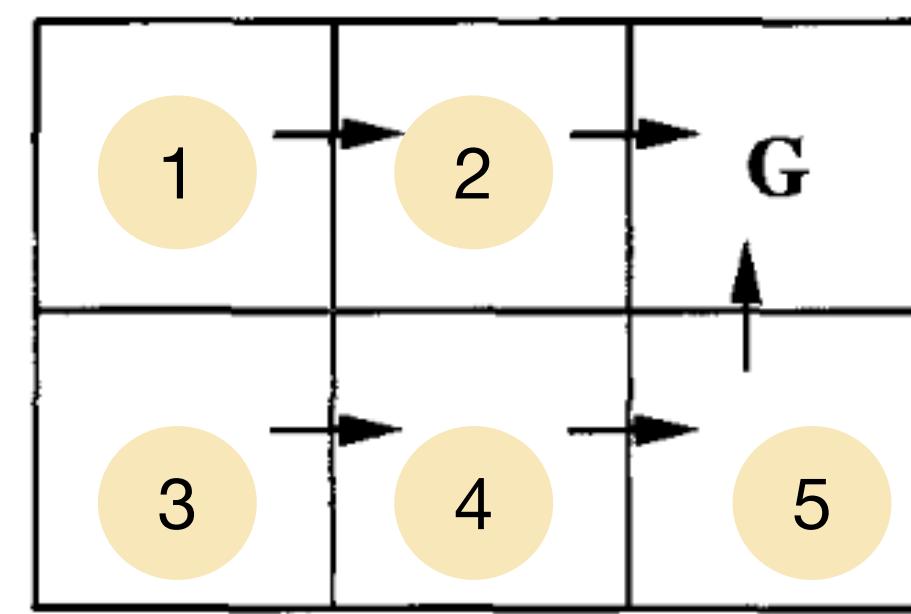
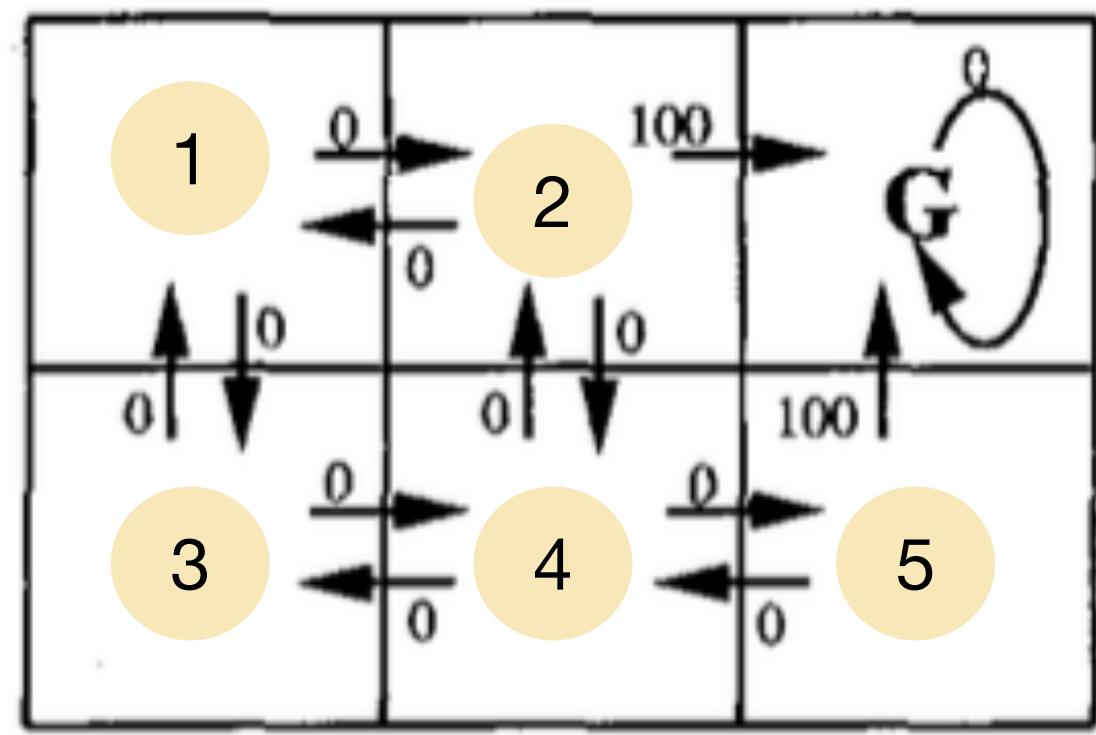
The optimal policy is defined as

$$\pi^*(s) \triangleq \arg \max_{\pi} V^{\pi}(s), \forall s$$

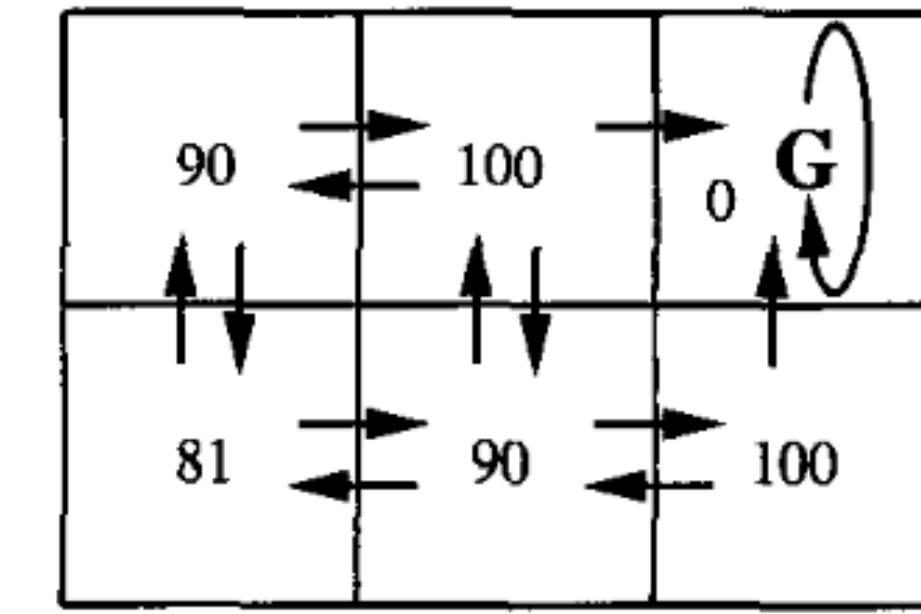
with the optimal policy value $V^{\pi^*}(s) \triangleq V^*(s)$

So, the optimal policy always yields the highest policy value for any starting state.

Example of Optimal Policy



One optimal policy



$V^*(s)$ values

Optimal Policy Value

The value of the optimal policy is

$$V^*(s_t) = \max_{a_t, a_{t+1}, a_{t+2}, \dots} \left\{ r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \gamma^2 r(s_{t+2}, a_{t+2}) + \dots \right\}$$

Since the immediate reward $r_t = r(s_t, a_t)$ only depends on s_t and a_t

$$\begin{aligned} V^*(s_t) &= \max_{a_t} [r(s_t, a_t) + \gamma \max_{a_{t+1}, a_{t+2}, \dots} (r(s_{t+1}, a_{t+1}) + \gamma r(s_{t+2}, a_{t+2}) + \dots)] \\ &= \max_{a_t} [r(s_t, a_t) + \gamma V^*(s_{t+1})] \\ V^*(s) &= \max_a [r(s, a) + \gamma V^*(\delta(s, a))] \end{aligned} \tag{1}$$

Recurrent relation

If $r(s, a)$ and $\delta(s, a)$ are known beforehand, $V^*(s)$ can be iteratively approximated **off-line**, known as **off-line learning** or **model-based learning**.

Off-line Learning: Value Iteration Approximation

For known $r(\cdot)$ and $\delta(\cdot)$, $V^*(s)$ can be iteratively approximated by using

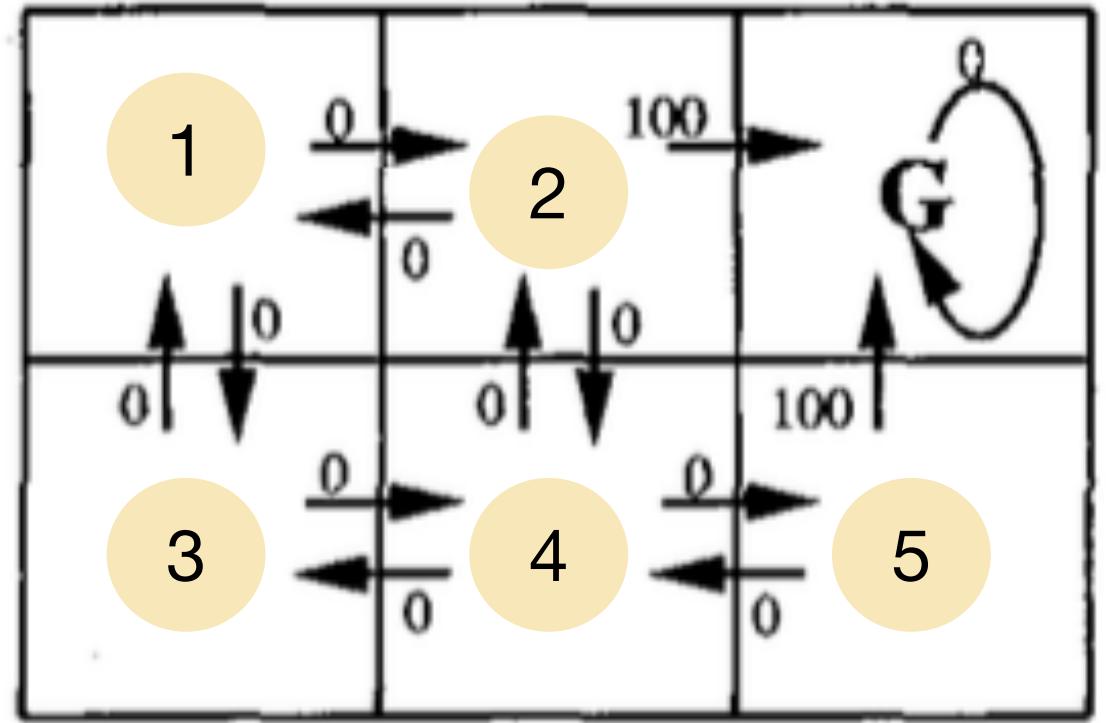
$$\hat{V}(s) \leftarrow \max_a [r(s, a) + \gamma \hat{V}(\delta(s, a))] \quad (2)$$

to reward w/ max score

Stop when $\hat{V}(s)$ converges.

From (1), once $V^*(s)$ has been determined, the optimal policy at each state s is obtained from

$$\pi^*(s) = \arg \max_a [r(s, a) + \gamma V^*(\delta(s, a))] \quad (3)$$



$$\hat{V}(s) \leftarrow \max_a [r(s, a) + \gamma \hat{V}(\delta(s, a))]$$

Initial

s	1	2	3	4	5	G
$V(s)$	0	0	0	0	0	0

1st Iteration:

$$V(1) = \max (0 + 0.9 \times 0, 0 + 0.9 \times 0) = 0$$

$$V(2) = \max (0 + 0.9 \times 0, 0 + 0.9 \times 0, 100 + 0.9 \times 0) = 100$$

$$V(3) = 0$$

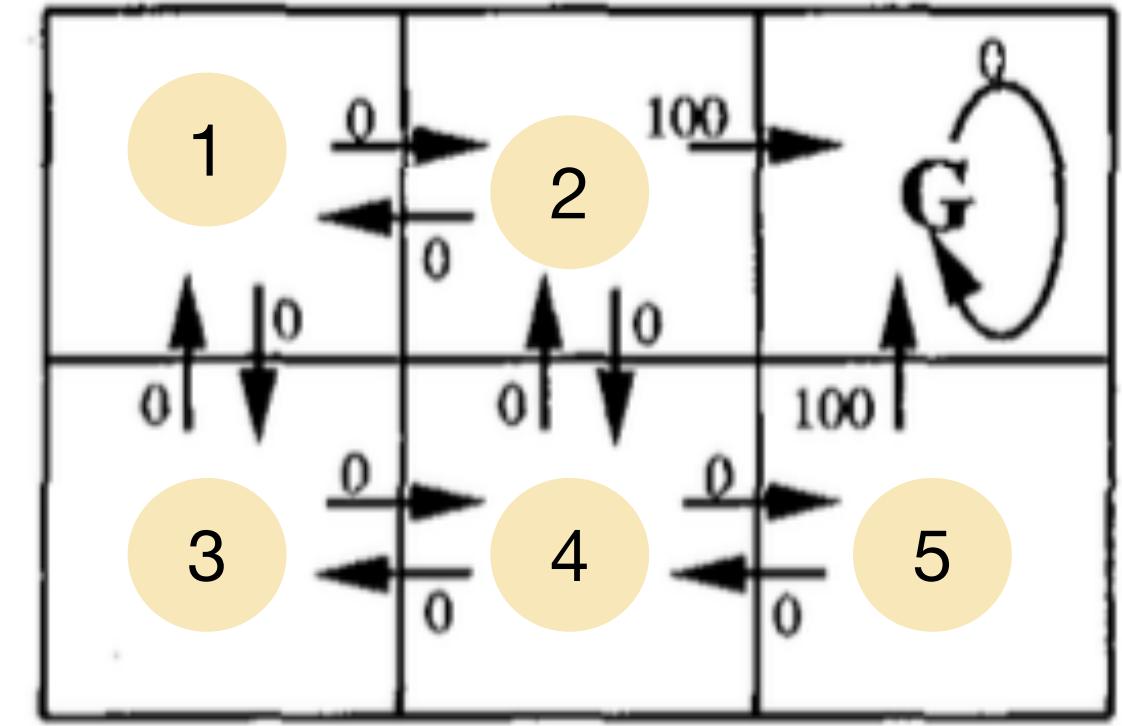
$$V(4) = 0$$

$$V(5) = 100$$

$$V(G) = 0$$

converges

s	1	2	3	4	5	G
$V(s)$	0	100	0	0	100	0



ສົດເບີທີ່ເດີປິກຈາ

$$\hat{V}(s) \leftarrow \max_a [r(s, a) + \gamma \hat{V}(\delta(s, a))]$$

s	1	2	3	4	5	G
$V(s)$	0	100	0	0	100	0

2nd Iteration:

$$V(1) = \max(0 + 0.9 \times 0, 0 + 0.9 \times 100) = 90$$

$$V(2) = 100$$

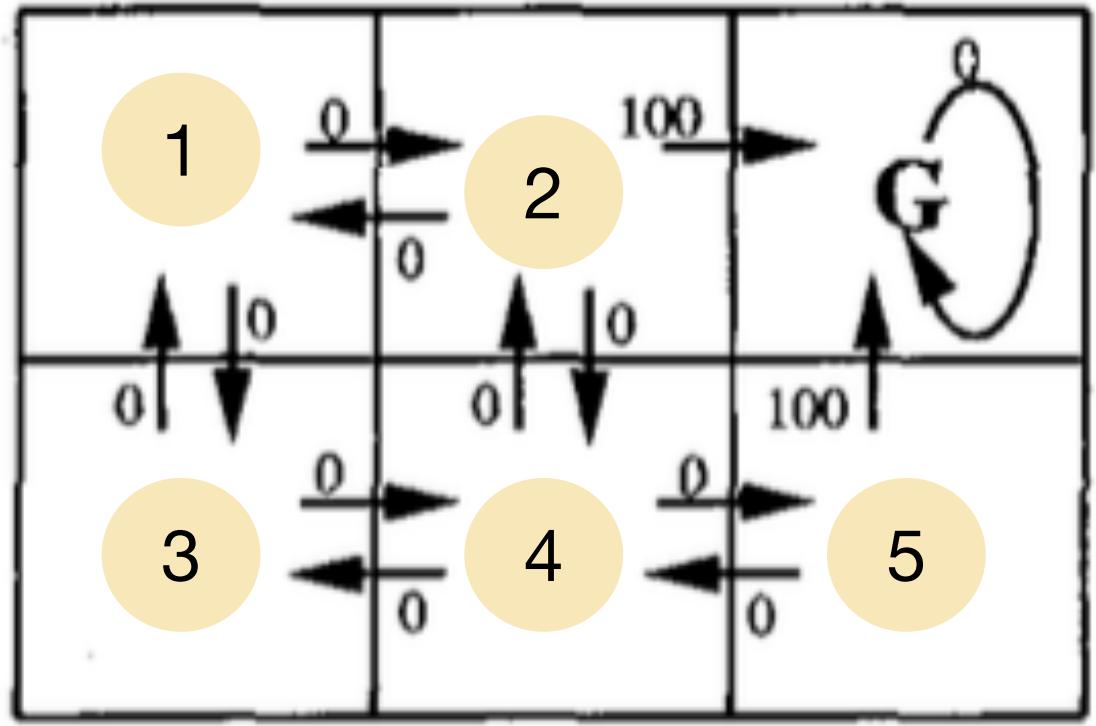
$$V(3) = \max(0 + 0.9 \times 90, 0 + 0.9 \times 0) = 81$$

$$V(4) = \max(0 + 0.9(81), 0 + 0.9(100), 0 + 0.9(100)) = 90$$

$$V(5) = 100$$

$$V(G) = 0$$

s	1	2	3	4	5	G
$V(s)$	90	100	81	90	100	0



$$\hat{V}(s) \leftarrow \max_a [r(s, a) + \gamma \hat{V}(\delta(s, a))]$$

discount factor

s	1	2	3	4	5	G
$V(s)$	90	100	81	90	100	0

3rd Iteration:

$$V(1) = 90$$

$$V(2) = 100$$

$$V(3) = 81$$

$$V(4) = 90$$

$$V(5) = 100$$

$$V(G) = 0$$

s	1	2	3	4	5	G
$V^*(s)$	90	100	81	90	100	0
$\pi^*(s)$	R	R	R, U	R, U	U	Self

$$\pi^*(s) = \arg \max [r(s, a) + \gamma V^*(\delta(s, a))]$$

state 1: $\max[0 + 0.9(100), 0 + 0.9(81)] = 90$ ↗ 90%

state 2: $\max(0 + 0.9(90), 0 + 0.9(90), 100 + 0.9(0)) = 100 - 75\%$

Action-Value or Quality (Q)-Function

In practice, the agent can only learn by interacting with the environment, i.e., **on-line learning**:

- ◆ Execute a certain action a_t in state s_t
- ◆ Observe the reward r_t and succeeding state s_{t+1} .

If we define the action-value or **Q function**:

$$Q(s_t, a_t) = \max_{a_{t+1}, a_{t+2}, \dots} \left\{ r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \gamma^2 r(s_{t+2}, a_{t+2}) + \dots \right\} \quad (4)$$

Immediate reward max future discounted cumulative reward

Then, the optimal policy and value can be determined from

$$\begin{aligned} \pi^*(s_t) &= \arg \max_{a_t} Q(s_t, a_t) \\ \pi^*(s) &= \arg \max_a Q(s, a) \end{aligned} \quad (5a)$$

$$V^*(s) = \max_a Q(s, a) \quad (5b)$$

One must express $Q(s_t, a_t)$ in the **recurrence relation** form.

The agent can iteratively approximate the Q function based on *observed rewards and next states* of each action taken without the knowledge of $r(s, a)$ and $\delta(s, a)$.

$$Q(s_t, a_t) = \max_{a_{t+1}, a_{t+2}, \dots} \left\{ r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \gamma^2 r(s_{t+2}, a_{t+2}) + \dots \right\}$$
$$Q(s, a) = r(s, a) + \gamma \max_{a'} Q(\delta(s, a), a')$$

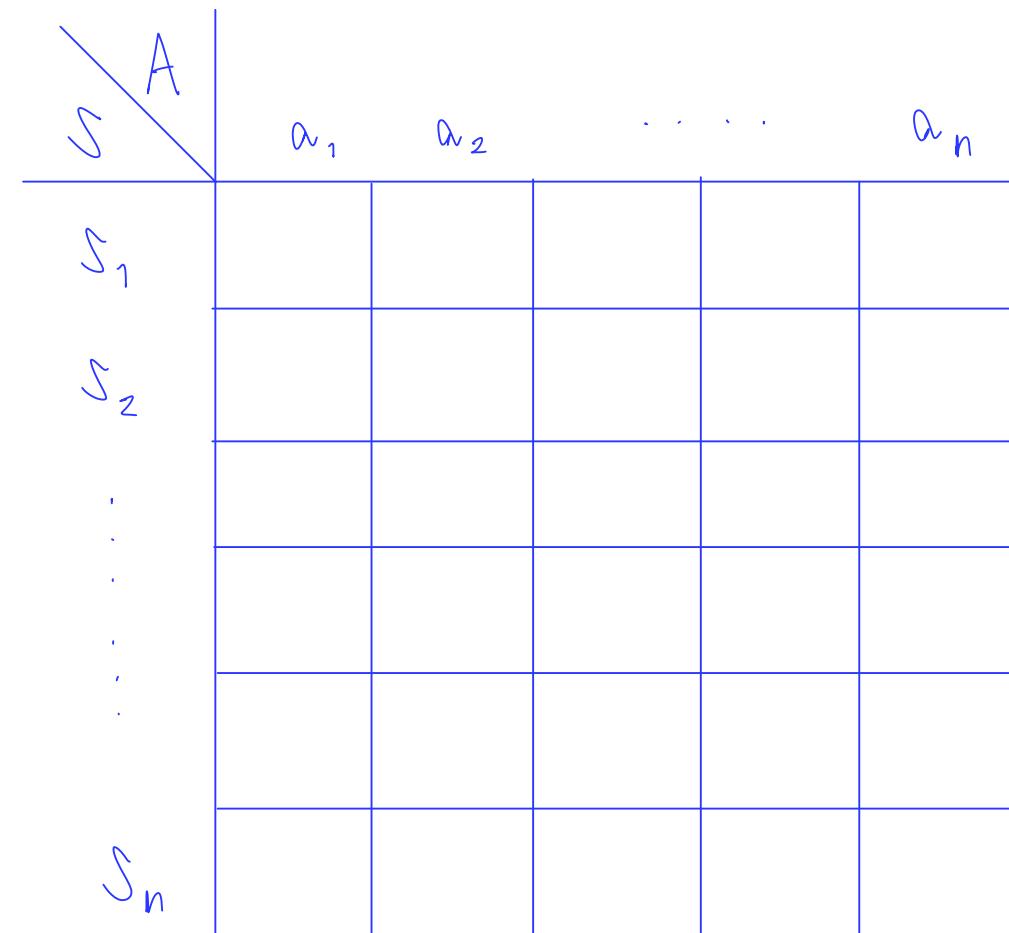
(6)

Q-Learning Algorithm

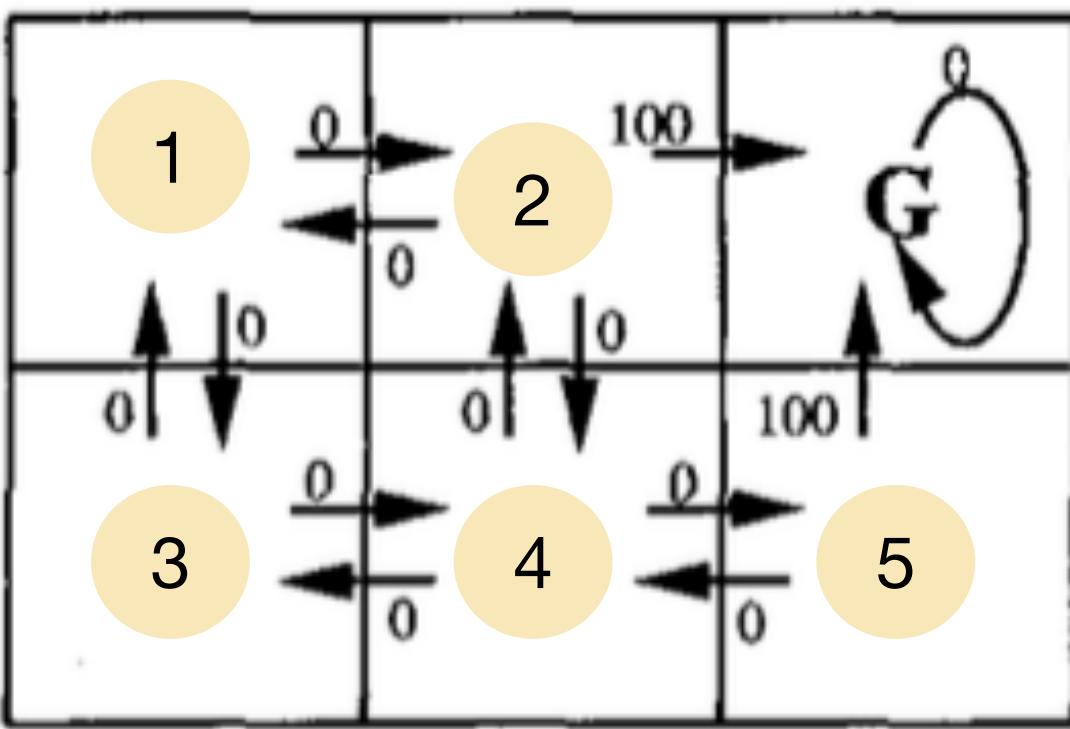
From (6), by observing reward $r = r(s, a)$ and next state $s' = \delta(s, a)$, the learner can iteratively approximate $Q(s, a)$ with:

$$\hat{Q}(s, a) \leftarrow r + \gamma \max_{a'} \hat{Q}(s', a') \quad (7)$$
$$s \leftarrow s'$$

The optimal policy is determined from (5a).



Q-table at the start of episode 1



	Left	Right	Up	Down
1	0	0	0	0
2	0	0	0	0
3	0	0	0	0
4	0	0	0	0
5	0	0	0	0
G	0	0	0	0

Episode 1 (Initial state 4) เลือก action ที่ให้ค่า Q สูงสุด
หรือ random otherwise

t=0, s = 4, a = right, s' = 5

$$Q(s, \text{up}) = 100 + 0.9(0) = 100$$

t=1, s = 5, a = up, s' = G

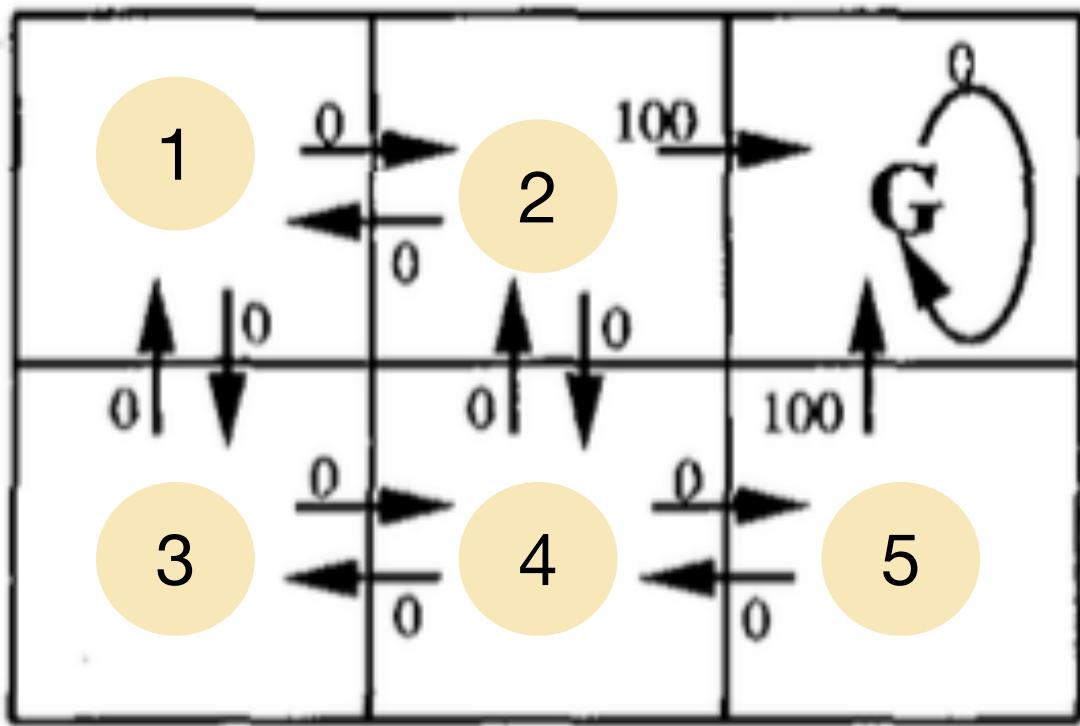
	Left	Right	Up	Down
1				
2				
3				
4	0			
5				
G				

$$Q(4, 2) \leftarrow 0 + 0.9(0) = 0$$

	Left	Right	Up	Down
1				
2				
3				
4	0			
5			100	
G				

$$Q(5, \text{up}) \leftarrow 100 + 0.9(0)$$

Q-table at the start of episode 2



	Left	Right	Up	Down
1				
2				
3				
4		0		
5			100	
G				

Episode 2 (Initial state 4)

$$Q(4, \text{Right}) = 0 + 0.9 \times 100$$

t=0: s = 4, a = right, s' = 5

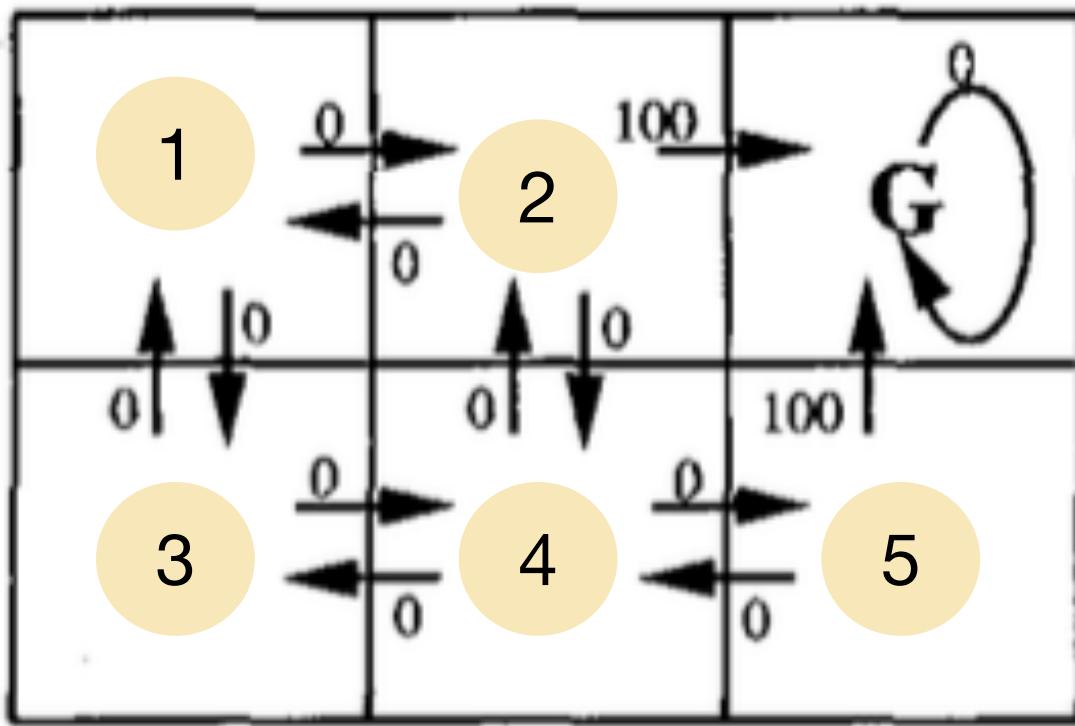
t=1: s = 5, a = up, s' = G

	Left	Right	Up	Down
1				
2				
3				
4		90		
5			100	
G				

	Left	Right	Up	Down
1				
2				
3				
4			90	
5				100
G				

$$Q(s, \text{up}) \leftarrow 100 + 0.9 \times 0$$

Q-table at the start of episode 3



	Left	Right	Up	Down
1				
2				
3				
4			90	
5			100	
G				

Episode 3 (Initial state 2)

$t=0: s = 2, a = \text{down}, s' = 4$

	Left	Right	Up	Down
1				
2				81
3				
4		90		
5			100	
G				

$$Q(2, \text{D}) = 0 + 0.9(90) \\ = 81$$

$t=1: s = 4, a = \text{right}, s' = 5$

	Left	Right	Up	Down
1				
2				81
3				
4		90		
5			100	
G				

$$Q(4, \text{R}) = 0 + 0.9(100) = 90$$

$t=2: s = 5, a = \text{up}, s' = G$

	Left	Right	Up	Down
1				
2				81
3				
4			90	
5			100	
G				

$$Q(5, \text{U}) = 100 + 0.9(0) = 100$$

Behavior Policy - Experimentation Strategies

The agent tries to learn the *optimal/target* policy from its history of interaction with the environment by using the *behavior policy*.

$$(s_1, a_1, r_1), (s_2, a_2, r_2), (s_3, a_3, r_3), \dots, (s_t, a_t, r_t), \dots$$

Behavior policy has the following two important strategies:

Exploitation Choose action that maximizes \hat{Q} .

Exploration Try different actions other than the current best, which could improve \hat{Q} in a long run.

Examples of Behavior Policy

Greedy action selection policy

Random action selection policy

ϵ -greedy selection policy $\xrightarrow{\text{Greed with prob. } 1 - \epsilon}$

Gibbs or **Boltzmann** action selection policy

$$\mathbb{P}\{a_i \mid s\} = \frac{e^{\hat{Q}(s, a_i)/\tau}}{\sum_j e^{\hat{Q}(s, a_j)/\tau}}, \quad \tau > 0$$

Math Framework of Reinforcement Learning - Markov Decision Process (MDP)

In general, the reward and state transition functions can have **non-deterministic** outcomes.

RL is mathematically formalized under the framework of **Markov Decision Process (MDP)**.

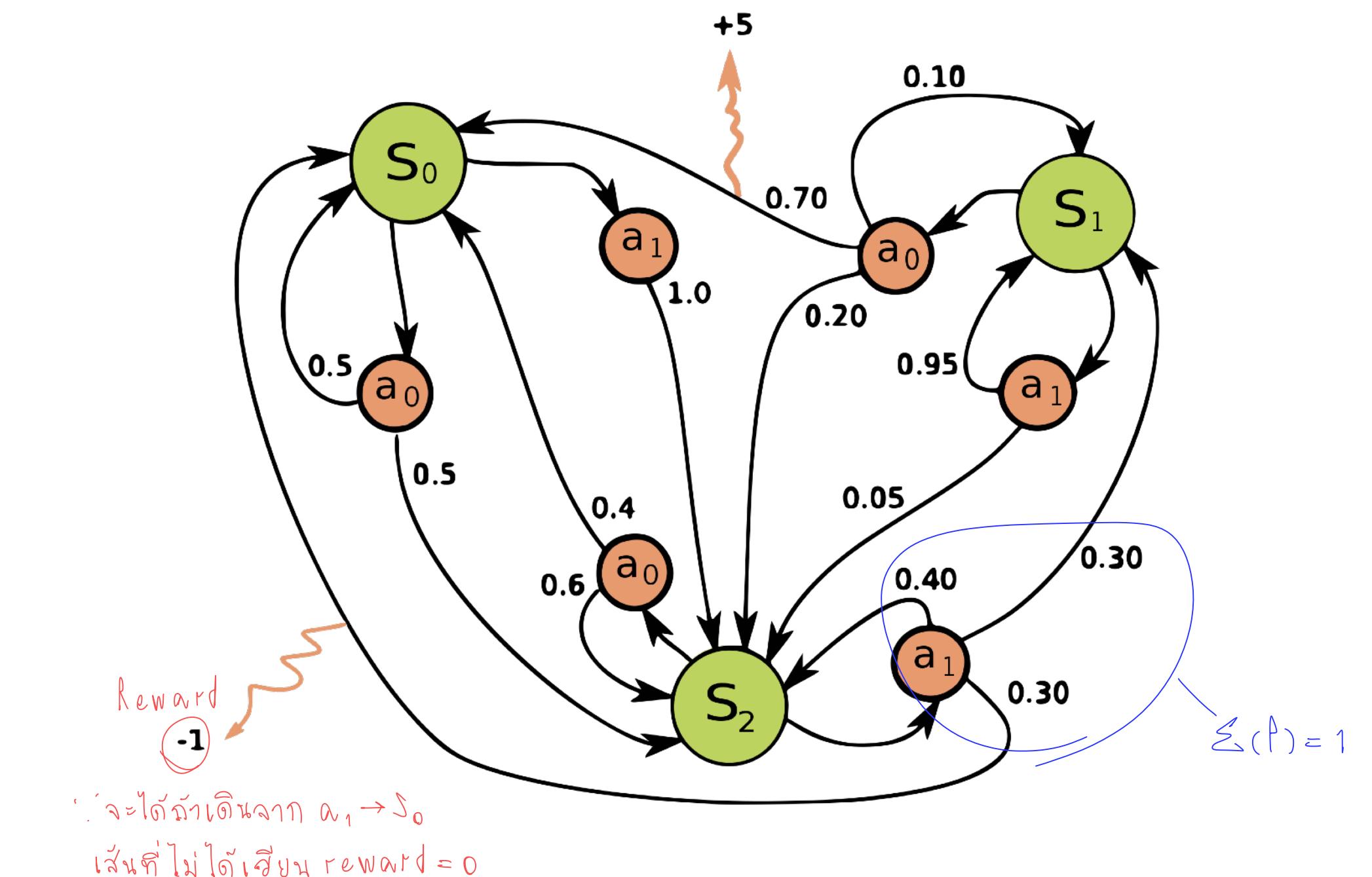
Four elements of MDP:

- ◆ State space S
- ◆ Action space A
- ◆ State-Action transition probabilities
- ◆ State-Action reward probabilities

$$\mathbb{P}\{s' \mid s, a\} \triangleq \mathbb{P}\{s_{t+1} = s' \mid s_t = s, a_t = a\}$$

$$\mathbb{E}\{r_t(s_t, a_t)\} = \sum_r r \mathbb{P}\{r_t = r \mid s_t = s, a_t = a\}$$

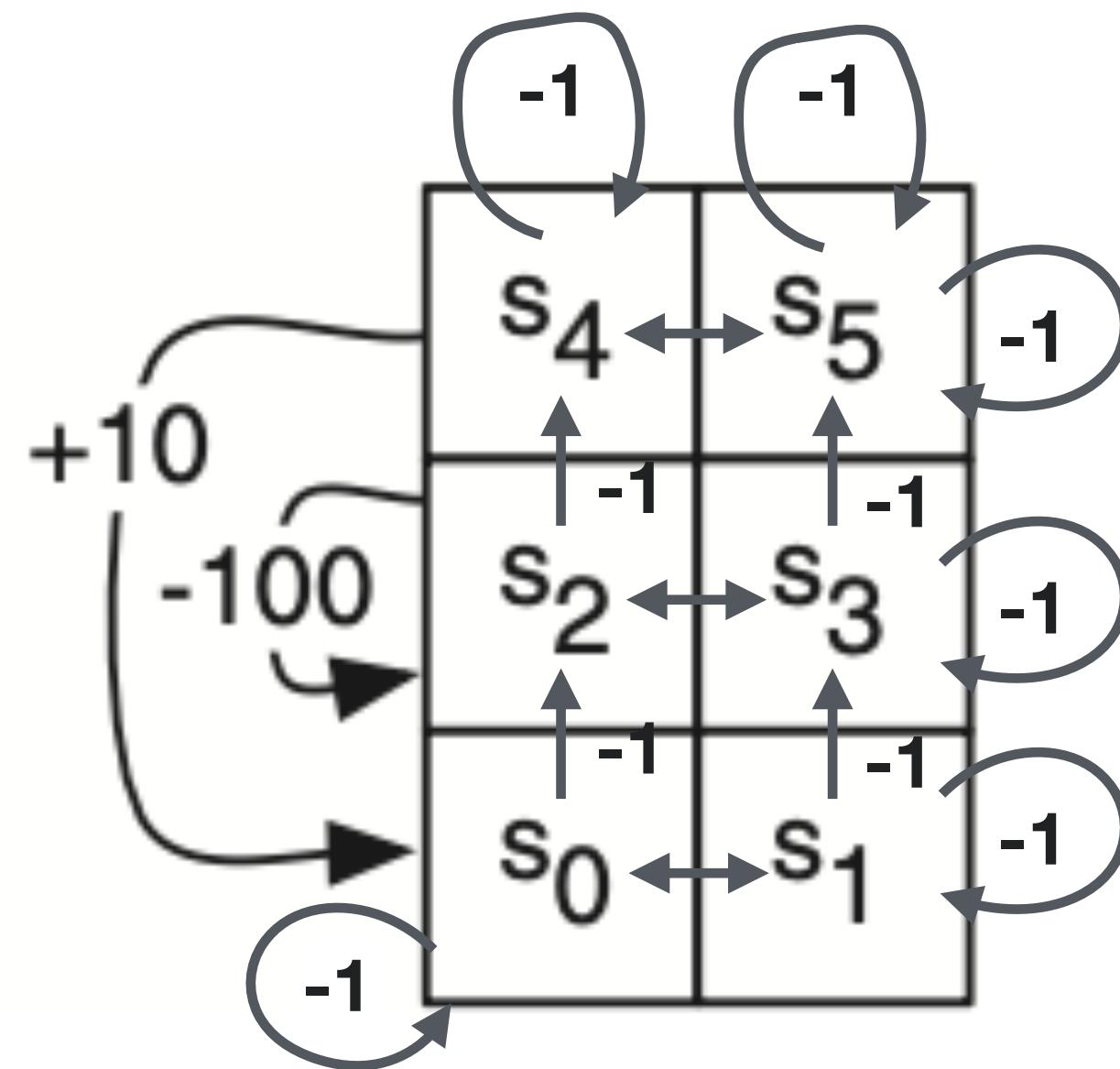
Globally optimal policy is analytically solvable under the perfect knowledge of reward and transition structures.



Example of a simple MDP with three states (green circles) and two actions (orange circles), with two rewards (orange arrows)

https://en.wikipedia.org/wiki/Markov_decision_process

Example



Set of actions and rewards

upC (for “up carefully”) The agent goes up, except in states s_4 and s_5 , where the agent stays still, and has a reward of -1 .

right The agent moves to the right in states s_0, s_2, s_4 with a reward of 0 and stays still in the other states, with a reward of -1 .

left The agent moves one state to the left in states s_1, s_3, s_5 . In state s_0 , it stays in state s_0 and has a reward of -1 . In state s_2 , it has a reward of -100 and stays in state s_2 . In state s_4 , it gets a reward of 10 and moves to state s_0 .

up With a probability of 0.8 it acts like *upC*, except the reward is 0. With probability 0.1 it acts as a *left*, and with probability 0.1 it acts as *right*.

* Determine the policy to maximize the long-term cumulative reward.

$$P(s_0 | s_4, \text{Up}) = 0.1$$

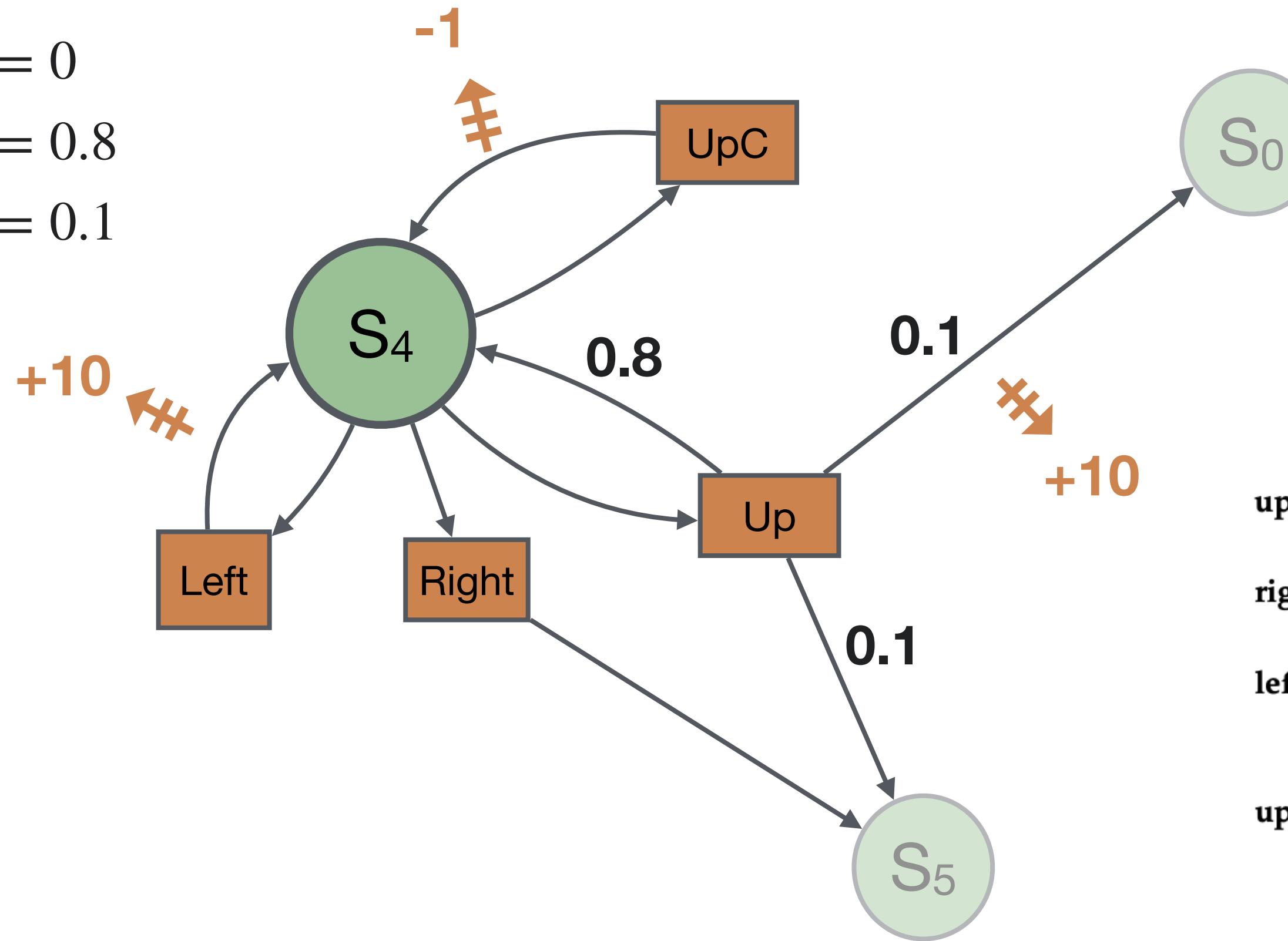
$$P(s_1 | s_4, \text{Up}) = 0$$

$$P(s_2 | s_4, \text{Up}) = 0$$

$$P(s_3 | s_4, \text{Up}) = 0$$

$$P(s_4 | s_4, \text{Up}) = 0.8$$

$$P(s_5 | s_4, \text{Up}) = 0.1$$



upC (for “up carefully”) The agent goes up, except in states s_4 and s_5 , where the agent stays still, and has a reward of -1.

right The agent moves to the right in states s_0, s_2, s_4 with a reward of 0 and stays still in the other states, with a reward of -1.

left The agent moves one state to the left in states s_1, s_3, s_5 . In state s_0 , it stays in state s_0 and has a reward of -1. In state s_2 , it has a reward of -100 and stays in state s_2 . In state s_4 , it gets a reward of 10 and moves to state s_0 .

up With a probability of 0.8 it acts like upC , except the reward is 0. With probability 0.1 it acts as a *left*, and with probability 0.1 it acts as *right*.

Temporal Difference (TD) Error

Consider sequence of observed numeric values v_1, v_2, v_3, \dots

Let A_k be an estimate of the expected value of (most recent) k data points. Then,

$$\begin{aligned} A_k &= \frac{v_1 + v_2 + \dots + v_k}{k} \\ &= A_{k-1} + \alpha_k \underbrace{(v_k - A_{k-1})}_{\text{previous estimate}}, \quad \alpha_k = \frac{1}{k} \end{aligned}$$

Temporal difference (TD) error

$$\begin{array}{l} v_k > A_{k-1} \\ v_k < A_{k-1} \end{array}$$

- ◆ How does A_k change with the latest observed value?
- ◆ How much A_k change with the latest observed value?

In RL, more recent values are more accurate and could be weighted more as the learning progress.

Q-Learning in a Non-Deterministic Environment

One form of *temporal difference learning* where the TD error is the change in Q value.

The agent learns from the history of interaction as a sequence of *experiences* $\langle s, a, r, s' \rangle$ using

$$\hat{Q}(s, a) \leftarrow \hat{Q}(s, a) + \underbrace{\alpha [r(s, a) + \gamma \max_{a'} \hat{Q}(s', a')]}_{\substack{\text{Actual return} = \text{Observed reward} \\ + \text{Discounted future reward}}} - \underbrace{\hat{Q}(s, a)}_{\text{Current estimated return}}, \quad 0 < \alpha < 1$$

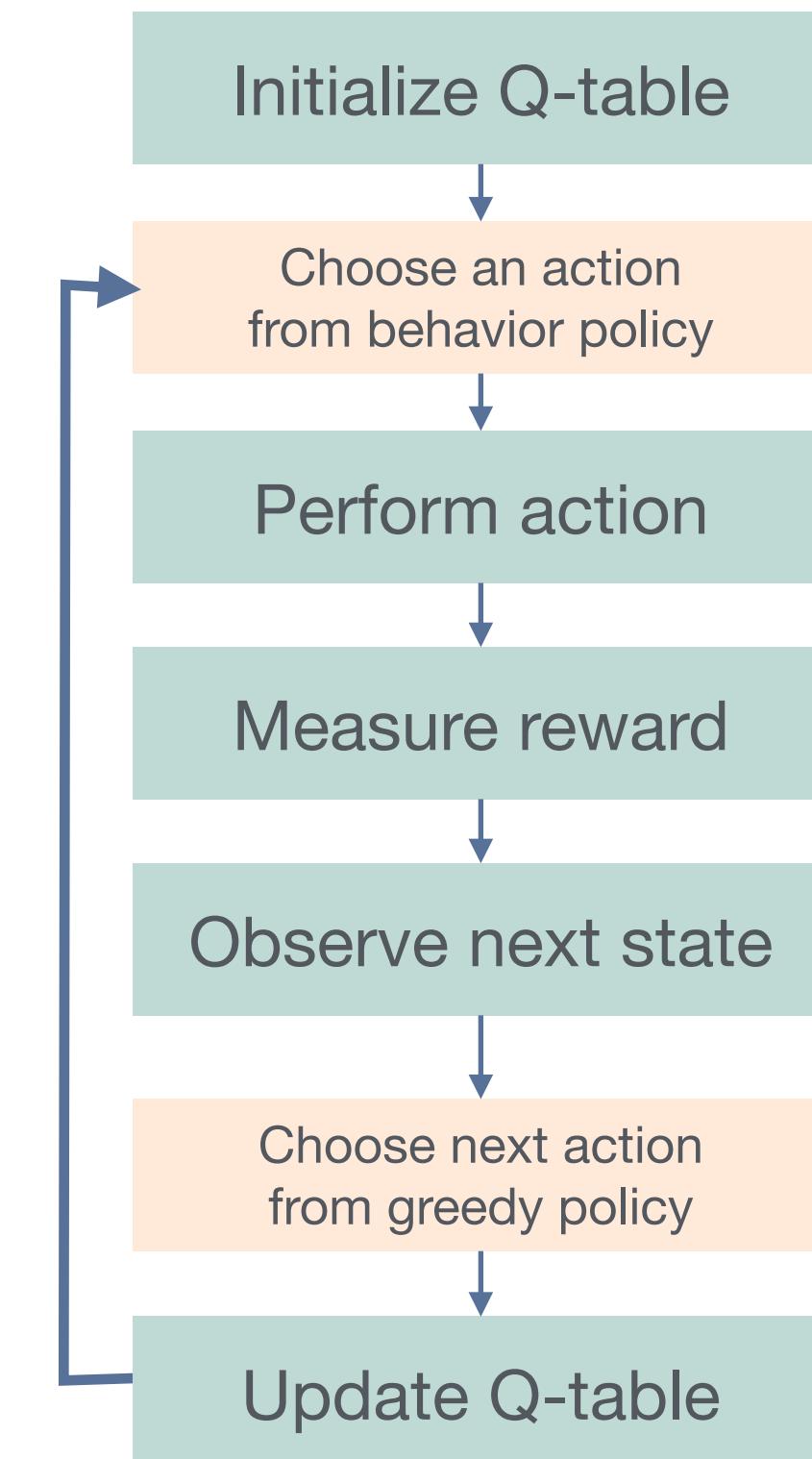
- ◆ Current action a in state s chosen by using the behavior policy.
- ◆ Next state action a' chosen to maximize the actual return (the greedy policy).

Regarded as an *off-policy* learning algorithm - Learn an optimal policy regardless of the behavior policy used during learning provided that it explores enough state-action pairs.

```

1: controller Q-learning( $S, A, \gamma, \alpha$ )
2:   Inputs
3:      $S$  is a set of states
4:      $A$  is a set of actions
5:      $\gamma$  the discount
6:      $\alpha$  is the step size
7:   Local
8:     real array  $Q[S, A]$ 
9:     previous state  $s$ 
10:    previous action  $a$ 
11:    initialize  $Q[S, A]$  arbitrarily
12:    observe current state  $s$ 
13:    repeat
14:      select and carry out an action  $a$ 
15:      observe reward  $r$  and state  $s'$ 
16:       $Q[s, a] \leftarrow Q[s, a] + \alpha (r + \gamma \max_{a'} Q[s', a'] - Q[s, a])$ 
17:       $s \leftarrow s'$ 
18:    until termination

```



$$\hat{Q}(s, a) \leftarrow \hat{Q}(s, a) + \alpha [r(s, a) + \gamma \max_{a'} \hat{Q}(s', a') - \hat{Q}(s, a)]$$

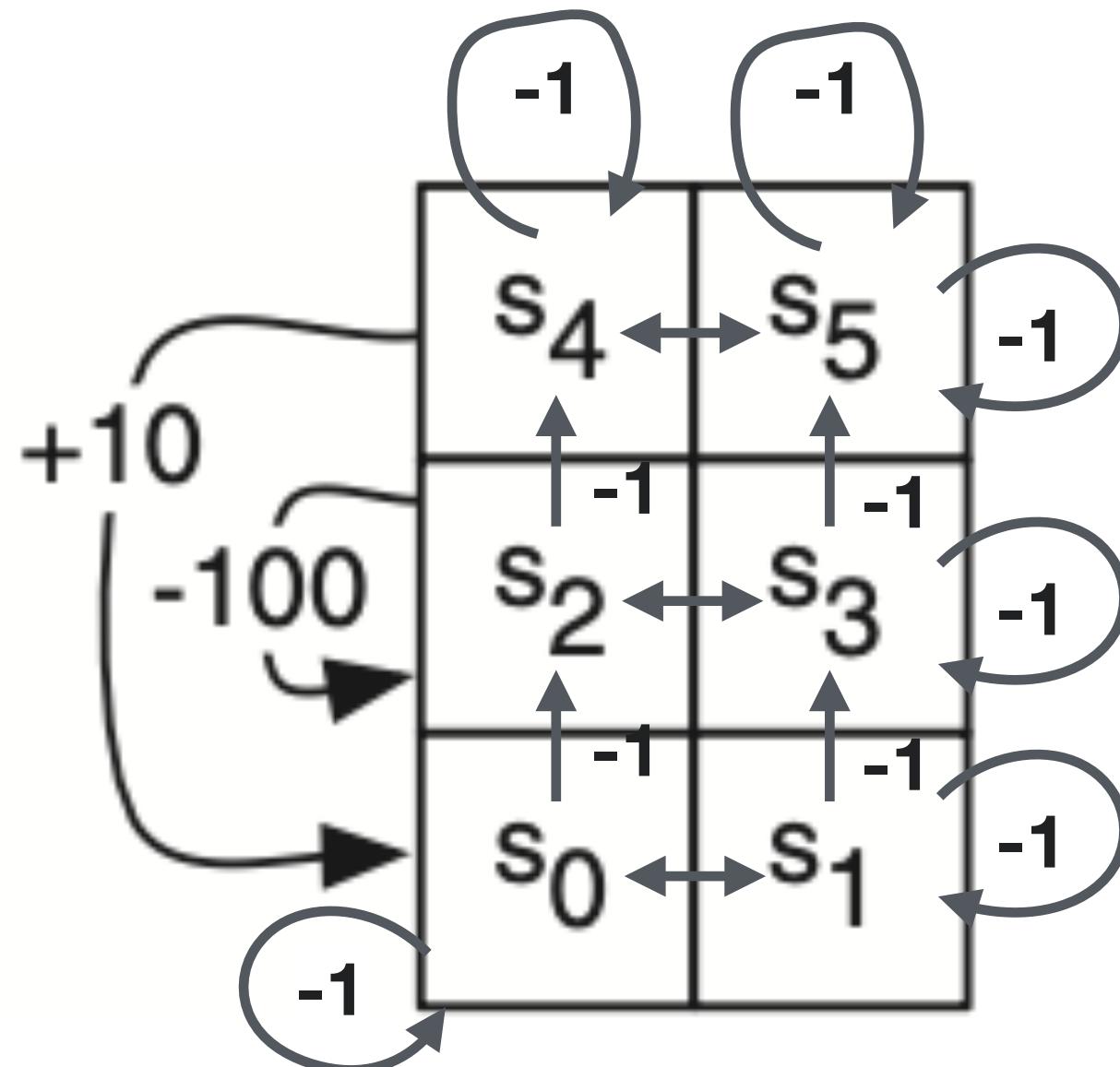
$$\alpha = 0.2, \gamma = 0.9$$

upC (for “up carefully”) The agent goes up, except in states s_4 and s_5 , where the agent stays still, and has a reward of -1 .

right The agent moves to the right in states s_0, s_2, s_4 with a reward of 0 and stays still in the other states, with a reward of -1 .

left The agent moves one state to the left in states s_1, s_3, s_5 . In state s_0 , it stays in state s_0 and has a reward of -1 . In state s_2 , it has a reward of -100 and stays in state s_2 . In state s_4 , it gets a reward of 10 and moves to state s_0 .

up With a probability of 0.8 it acts like *upC*, except the reward is 0 . With probability 0.1 it acts as a *left*, and with probability 0.1 it acts as *right*.



Initial Q-table

	UpC	Right	Left	Up
S0	0	0	0	0
S1	0	0	0	0
S2	0	0	0	0
S3	0	0	0	0
S4	0	0	0	0
S5	0	0	0	0

History of interaction

	<i>s</i>	<i>a</i>	<i>r</i>	<i>s'</i>
1	<i>s</i> ₀	<i>upC</i>	-1	<i>s</i> ₂
2	<i>s</i> ₂	<i>up</i>	0	<i>s</i> ₄
3	<i>s</i> ₄	<i>left</i>	10	<i>s</i> ₀
4	<i>s</i> ₀	<i>upC</i>	-1	<i>s</i> ₂

$$\hat{Q}(s, a) \leftarrow \hat{Q}(s, a) + \alpha [r(s, a) + \gamma \max_{a'} \hat{Q}(s', a') - \hat{Q}(s, a)]$$

step size

Q-Learning

$$\alpha = 0.2, \gamma = 0.9$$

$$\hat{Q}(S_2, Up) = 0 + 0.2(0 + 0.9(0) - 0) = 0$$

History of interaction

	<i>s</i>	<i>a</i>	<i>r</i>	<i>s'</i>
1	s_0	UpC	-1	s_2
2	s_2	Up	0	s_4
3	s_4	Left	10	s_0
4	s_0	UpC	-1	s_2

1st step
 $\hat{Q}(S_0, UpC) \leftarrow 0 + 0.2(-1 + 0.9(0) - 0) = -0.2$

	UpC	Right	Left	Up
S0	-0.2			
S1				
S2				
S3				
S4				
S5				

2nd step

	UpC	Right	Left	Up
S0				
S1				
S2				0
S3				
S4				
S5				

Initial Q-table

	UpC	Right	Left	Up
S0	0	0	0	0
S1	0	0	0	0
S2	0	0	0	0
S3	0	0	0	0
S4	0	0	0	0
S5	0	0	0	0

3rd step

	UpC	Right	Left	Up
S0	-0.2			
S1				
S2				
S3				
S4			2	
S5				

4th step

	UpC	Right	Left	Up
S0	-0.36			
S1				
S2				0
S3				
S4			2	
S5				

$$\begin{aligned} \hat{Q}(S_0, UpC) &= -0.2 + 0.2(-1 + 0.9(0) - (-0.2)) \\ &= -0.36 \end{aligned}$$

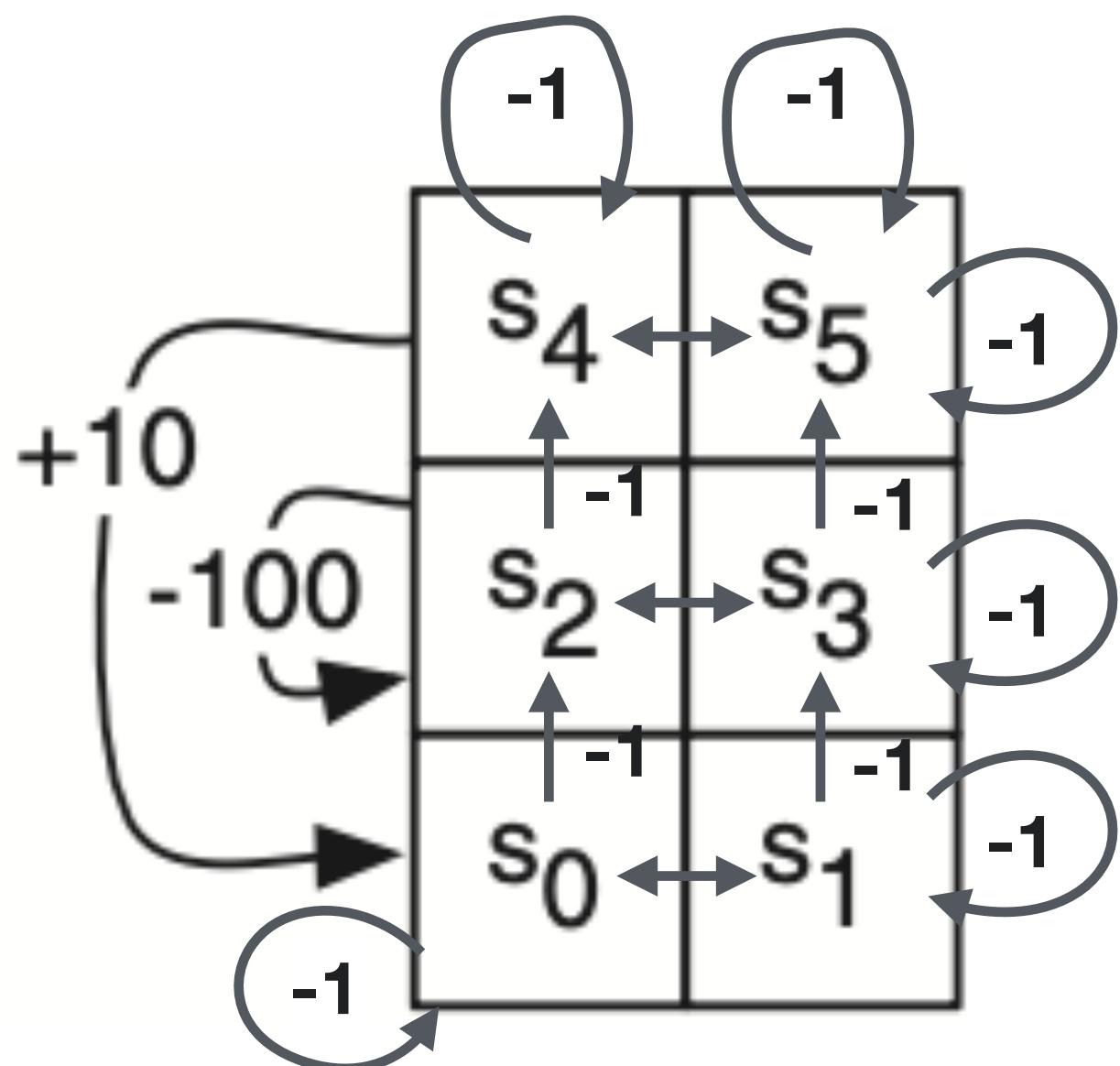
upC (for “up carefully”) The agent goes up, except in states s_4 and s_5 , where the agent stays still, and has a reward of -1 .

right The agent moves to the right in states s_0, s_2, s_4 with a reward of 0 and stays still in the other states, with a reward of -1 .

left The agent moves one state to the left in states s_1, s_3, s_5 . In state s_0 , it stays in state s_0 and has a reward of -1 . In state s_2 , it has a reward of -100 and stays in state s_2 . In state s_4 , it gets a reward of 10 and moves to state s_0 .

up With a probability of 0.8 it acts like *upC*, except the reward is 0. With probability 0.1 it acts as a *left*, and with probability 0.1 it acts as *right*.

	UpC	Right	Left	Up
S0				x
S1				x
S2	x			
S3				x
S4			x	
S5				x



Optimal from TD learning $\pi^*(s) = \arg \max_a Q(s, a)$

$$\pi^*(s_0) = \text{up}$$

$$\pi^*(s_1) = \text{up}$$

$$\pi^*(s_2) = \text{upC}$$

$$\pi^*(s_3) = \text{up}$$

$$\pi^*(s_4) = \text{left}$$

$$\pi^*(s_5) = \text{left}$$

Generalizing from Examples - Deep Q Learning

Q-learning requires that all state-action pairs are to be visited infinitely often for \hat{Q} to converge to actual Q .

Impossible to learn when the action-state space gets large (infinite or prohibitively large), the cost of execution actions is high.

Substitute a (deep) neural network for the lookup table and use each $\hat{Q}(s, a)$ update as a training example.

- ◆ One network with state and action as input.
- ◆ Train a separate network for each action with states as input.

