

# Advanced Machine Learning and Autonomous Agents

## Reinforcement Learning I



Jesse Read

# Goals for Today

**Week 1:** Deciding on actions given a knowledge representation

**Week 2:** Learning the knowledge representation, and complex inference

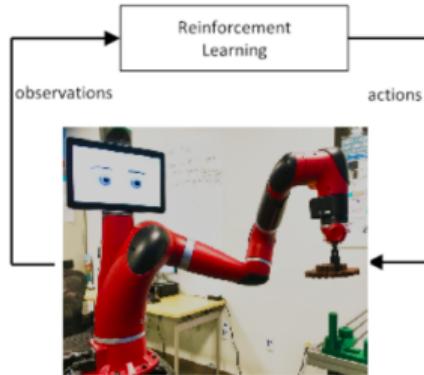
**Week 3:** Sequential Decision Making (Exploration-Exploitation tradeoff)

**Today:** Our decisions change the state of the environment (Reinforcement Learning)

## Objectives:

- ① Understand the main concepts behind Reinforcement Learning
- ② Get a feel for setting up RL problems in practice (warning: can be more subtle)

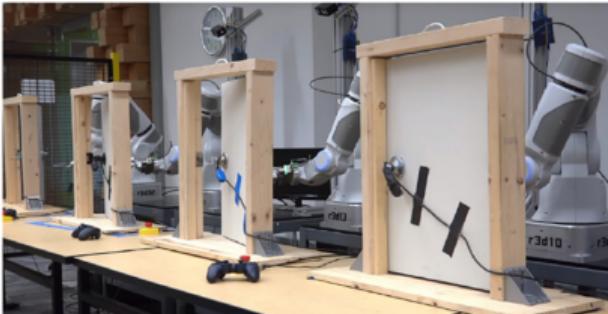
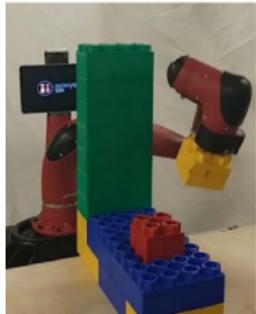
# Applications: Robotics and Autonomous Vehicles



► Helicopter

► Robot

► More Robots



# Applications: Logistics, Finance and Business

- Trading/finance/manage investment portfolio 
- In recommendation engines, email advertising
- Marketing and advertising; real-time bidding  



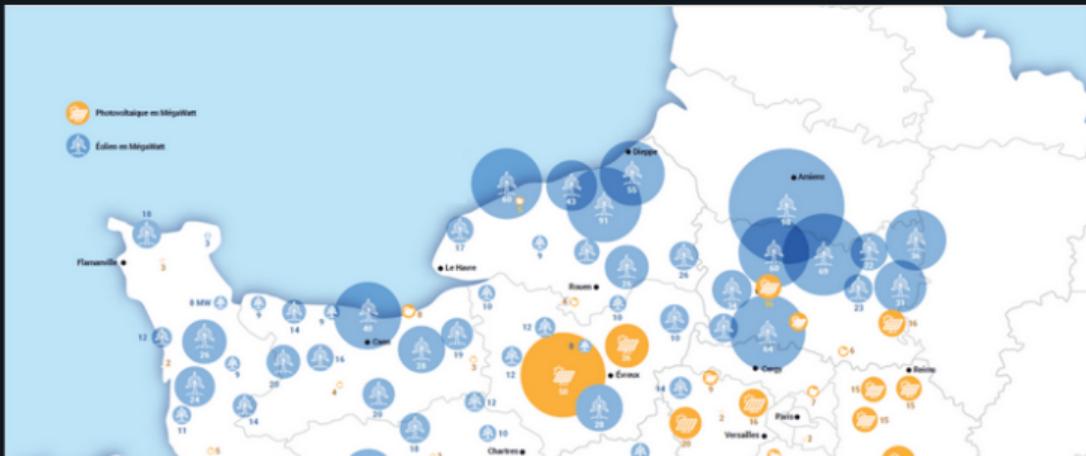
# Application: Energy Systems

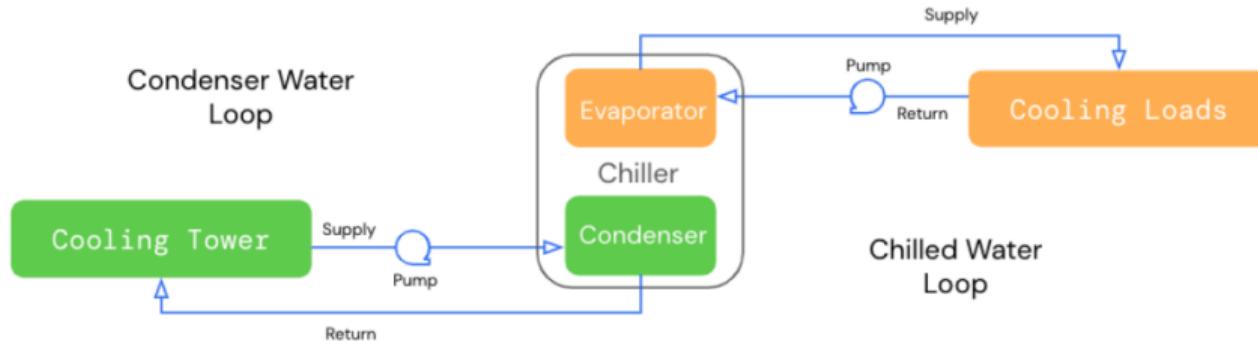
## La Région Ile-de-France et RTE lancent le Challenge IA pour la Transition Energétique

16.05.2023

TRANSITION ÉNERGÉTIQUE

Pour assurer l'alimentation électrique du territoire à chaque seconde, RTE surveille et pilote les flux d'électricité 24 heures sur 24, 7 jours sur 7, 365 jours par an. Le développement des énergies renouvelables, indispensable à la transition énergétique, nécessite d'adapter la gestion du système électrique. En effet, ces énergies sont variables en fonction des conditions météorologiques et réparties sur tout le territoire. Cela implique un pilotage en temps réel encore plus performant, avec des données plus nombreuses à traiter pour anticiper au mieux les situations. L'Intelligence Artificielle est une solution pour accompagner les dispatchers de RTE dans leurs missions.





Luo et al., Controlling Commercial Cooling Systems Using Reinforcement Learning, 2022. N.B. Also Energy

Technology

## Software update for world's wind farms could power millions more homes

An AI that predicts wind changes could boost wind turbine efficiency by 0.3 per cent, which globally would amount to enough extra electricity to keep a country running

By Matthew Sparkes

21 May 2023

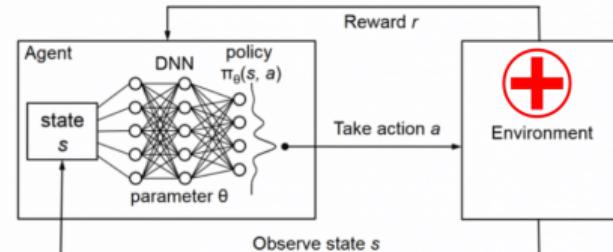


▲ A small efficiency gain for every wind turbine would lead to a big boost globally  
Justin Paget/Corbis/VCG via Getty Images

Current algorithms track wind patterns and adjust the turbine blades to anticipate changes, but [Alban Puech](#) and [Jesse Read](#) at the Polytechnic Institute of Paris believe artificial intelligence can do better. They have trained a reinforcement-learning algorithm to monitor wind patterns and develop its own strategy to maintain the turbine at the correct angle. And because

# Applications: Education, Healthcare, Agriculture, Politics, LLMs . . .

- Healthcare: Diagnosis, manage hospital resources, . . . ▶ Health
- Education: Personalised/auto-generated curriculum, . . .
- Farming and Agriculture: Managing resources, . . .
- Sports
- Politics: Obtaining a policy ▶ Politics



- Large Language Models (and other deep neural networks); HRL etc.

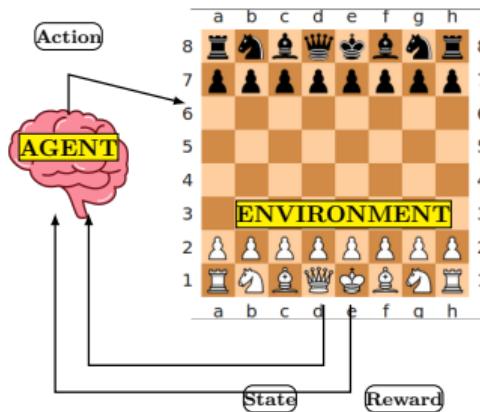


## General Setup: Agent and Environment; State, Action, Reward

- 1 General Setup: Agent and Environment; State, Action, Reward
- 2 Policy Search
- 3 Fully Observed Deterministic Environments
- 4 Stochastic Environments and MDPs
- 5 Learning: Return (Gain) and Value

# General Setup for Reinforcement Learning: Agent and the Environment

- ① The **agent** observes the **state** of the **environment**, and receives a **reward**
- ② The agent takes an available **action**
- ③ The next state is determined by **environment dynamics** ('rules of the game')

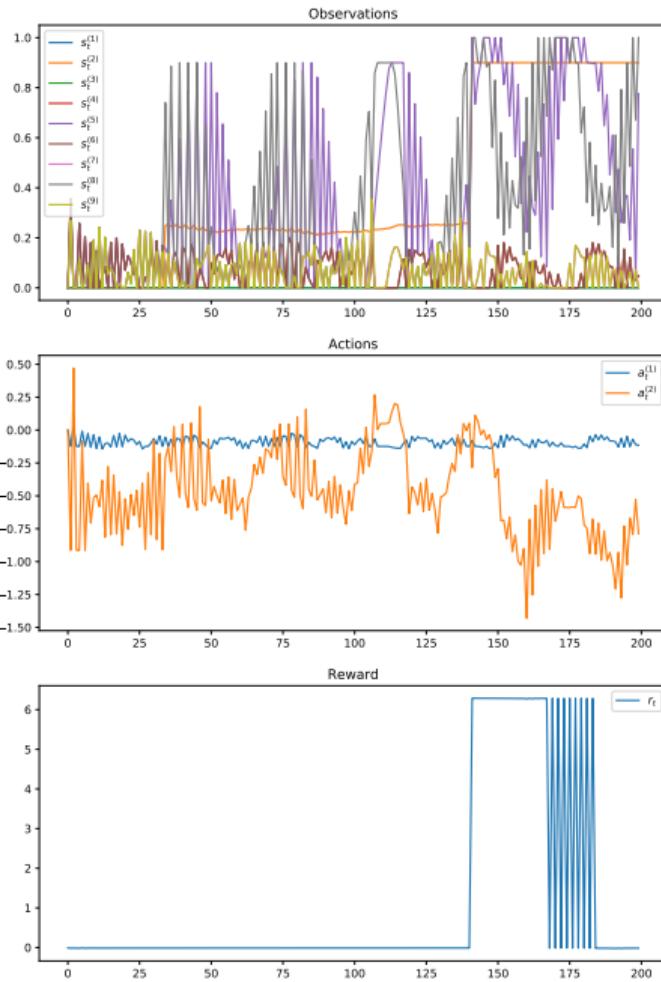


Important: The **agent may alter the environment** (and thus its own future state observations *and reward!*)

Environment-agent interaction (simplified):

```
env = Environment()
agent = Agent(env)
s = env.init()

while True:
    a = agent.act(s)
    s_next, r = env.step(a)
    agent.update((s,a,r,s_next))
    s = s_next
```



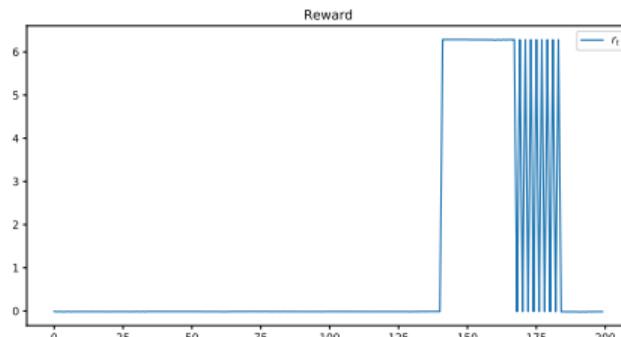
# The Reward Signal and Trajectories

At time  $t$  we observe  $s_t$ , take action  $a_t$ , obtain reward  $r_{t+1}$ , and advance to next state  $s_{t+1}$ . An episode (or trajectory,  $\tau$ ):

$$\tau = \{s_1, a_1, r_2, s_2, a_2, r_3, s_3, a_3, \dots, r_{T-1}, s_{T-1}, a_{T-1}, \dots, r_T, s_T\}$$

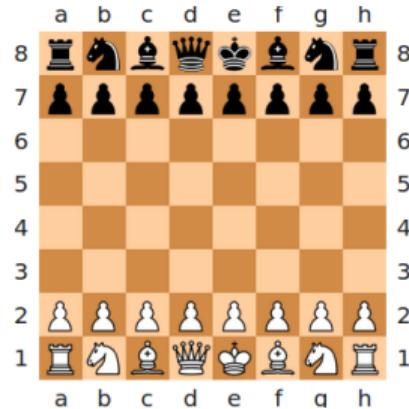
Challenges:

- Sparse rewards, e.g.,  $r_t = 0$  most of the time
- Weak rewards (limited impact of actions on reward)
- **Temporal credit assignment**: Which actions ( $a_{t-1}$ ?  $a_{t-100}$ ? both? none?) led to reward  $r_t$ ?



# The Environment

The **environment** provides interaction to the agent. Most importantly, it provides the environment dynamics, which determine how the agent can interact.



$$s' = f(s_t, a_t) \triangleright \text{Dynamic function}$$

$$r = r(s_t, a_t) \triangleright \text{Reward function}$$

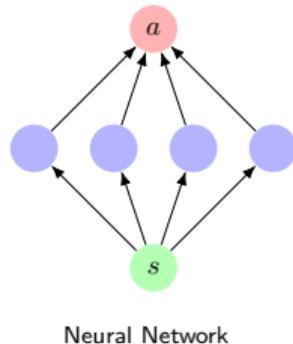
Note: often our intuition on real-world boundary between agent and its environment does not correspond to the modelisation of the problem.

# Agent and Policy

The policy defines the behaviour of the agent;

$$a_t = \pi(s_t)$$

indicates the action to take given the observation of the current state; i.e., the agent observes  $s_t$  and takes action  $a_t$  in response.



State $s_t$	Action $a_t$
0	+1
1	+1
2	-1
3	+1

Q-Table

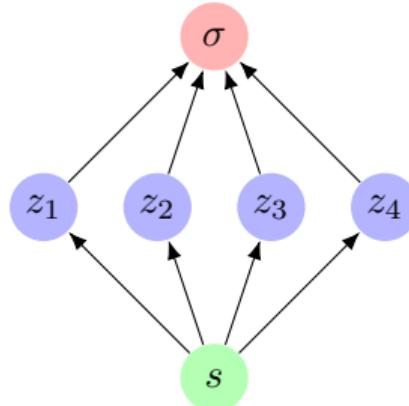
```
def act(s):  
    if s < 3:  
        # Go Right  
        return 1  
    else:  
        # Go Left  
        return -1
```

'Hand-coded' policy

A **stochastic policy** (provides a conditional probability distribution):

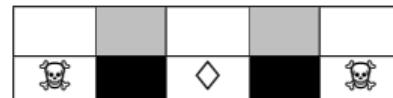
$$a_t \sim \pi(\cdot | s_t)$$

i.e.,  $P(A_t | S_t = s_t) = \pi(s_t)$ .



State $s_t$	$\pi(a_t = 1   s_t)$	$\pi(a_t = 0   s_t)$
	$\sigma$	$1 - \sigma$
0	0.1	0.9
1	0.5	0.5
2	1	0
3	0.2	0.8

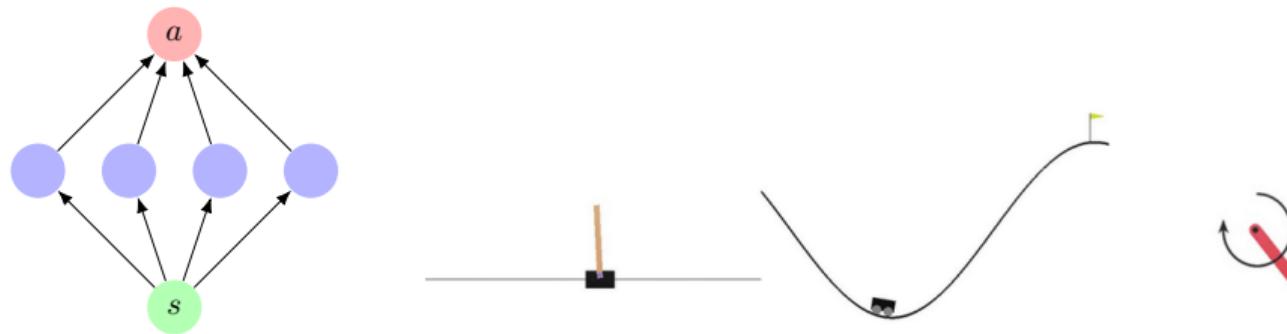
- Built-in exploration
- In competitive games, easier to defeat deterministic agents
- And **aliased states**: States that appear the same, but are not



# Policy Search

- 1 General Setup: Agent and Environment; State, Action, Reward
- 2 Policy Search
- 3 Fully Observed Deterministic Environments
- 4 Stochastic Environments and MDPs
- 5 Learning: Return (Gain) and Value

## Direct Policy Search (Black Box Optimization)



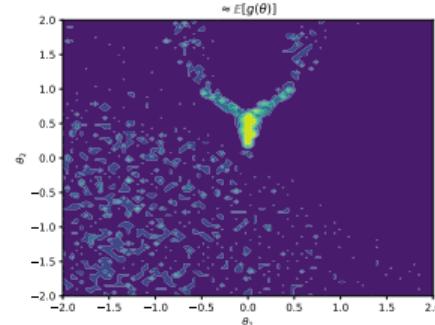
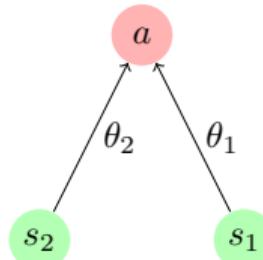
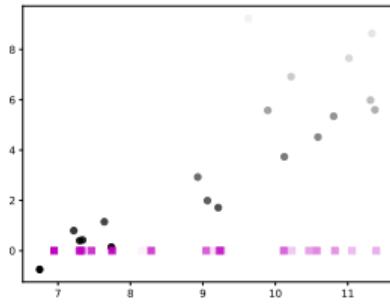
Suppose a **parametric policy**, parameters  $\theta$  (e.g., of a neural network agent); i.e.,

$$a_t = \pi_\theta(s_t)$$

we can search directly in parameter space, i.e., **policy search**:

$$\theta^* = \operatorname{argmax}_{\theta \in \Theta} g(\theta)$$

where  $g(\theta)$  indicates how good  $\pi_\theta$  is, e.g.,  $g(\theta) = 10 \Leftrightarrow$  agent survived 10 seconds / obtained score of 10 at end of episode.



Move the purple paddle left ( $a_t < 0$ ) or right ( $a_t > 0$ ) to catch a falling black ball (coordinates  $s_t = [s_1, s_2]$ )

- Many off-the-shelf search & optimization tools are available
- Works ‘embarrassingly well’ on some problems
- Handles continuous state/action spaces naturally
- No need to study reinforcement learning

But

- **Expensive** to evaluate  $g(\theta)$  (run simulations/play games)
- **Noisy** (large variance); randomness from the environment!
- **Curse of dimensionality**: optimization in high dimensions ( $> 100$ ) is very hard – it’s why we don’t train classifiers with black-box optimization!

# Fully Observed Deterministic Environments

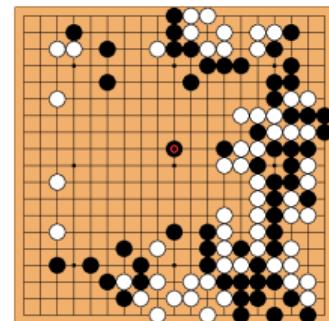
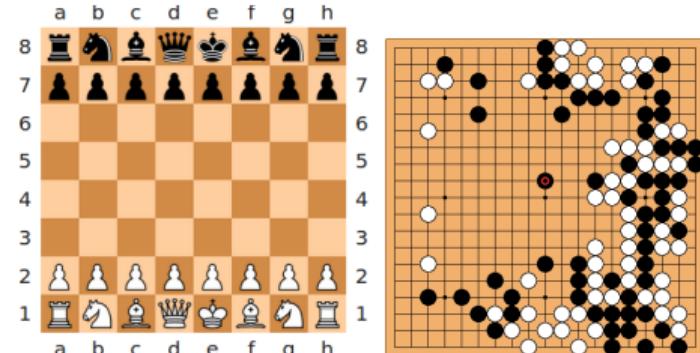
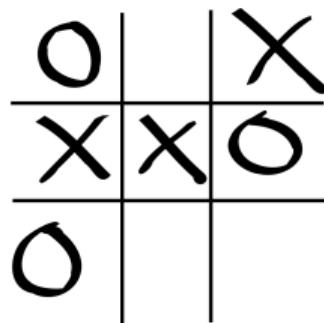
- 1 General Setup: Agent and Environment; State, Action, Reward
- 2 Policy Search
- 3 Fully Observed Deterministic Environments
- 4 Stochastic Environments and MDPs
- 5 Learning: Return (Gain) and Value

# Fully Observed Deterministic Environments

In **deterministic** environments: for any given action, the next state is known for certain.

In **fully observed** environments: complete view of the environment at time  $t$ .

Examples: 8-Puzzle, Noughts & Crosses (Tic Tac Toe), Chess, ...



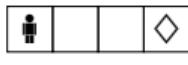
# Environment-Agent Interaction in Deterministic Environments

```
a0 ← π(s0)                                ▷ Initial action
for t = 1, 2, ... do
    st ← f(st-1, at-1)                ▷ Act, observe
    rt ← r(st-1, at-1)
    at ← π(st)                                ▷ Decide action
end for
```

- fully-observed case: full knowledge of  $f$ , and full view of  $s_t$ !
- reward function  $r$  also known

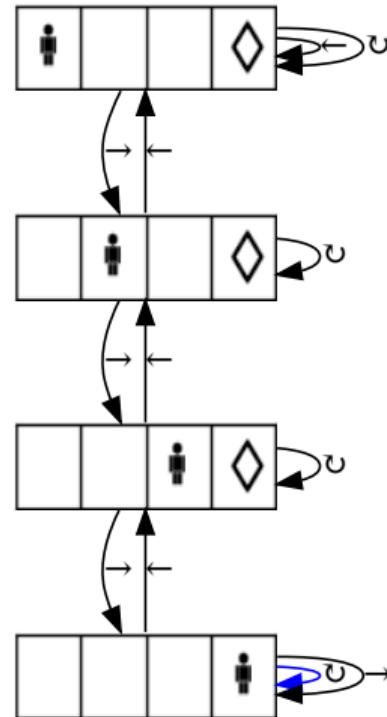
## Toy Example 1

- state  $s \in \{1, 2, 3, 4\}$
- action  $a \in \{\leftarrow, \circlearrowright, \rightarrow\}$
- reward function  $r(4, \circlearrowright) = 1$   
else  $r(s, a) = 0$  for other values of  $s, a$
- dynamics fn  $s' = f(s, a)$  depicted by  
*transition graph* (right)

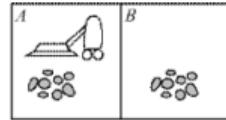
e.g.,  $s_t =$  

$$\square \boxed{\text{person}} \square \boxed{\diamond} = f(\boxed{\text{person}} \square \emptyset \boxed{\diamond}, \rightarrow)$$

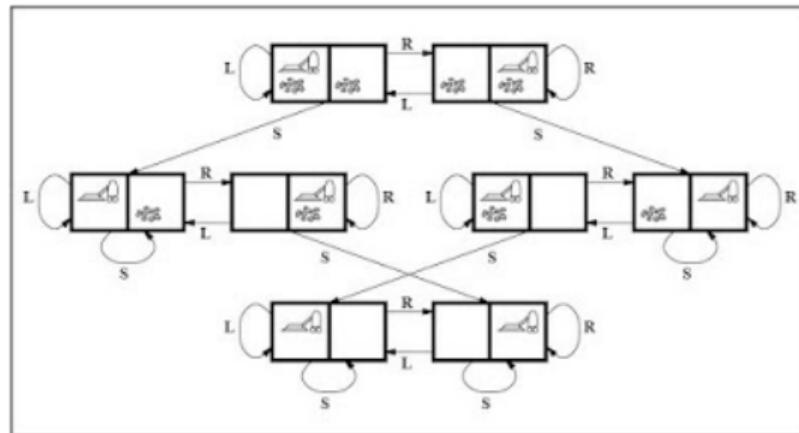
And end of episode at  $t$  when  $\circlearrowright = \pi(4)$



## Vacuum Cleaner World:

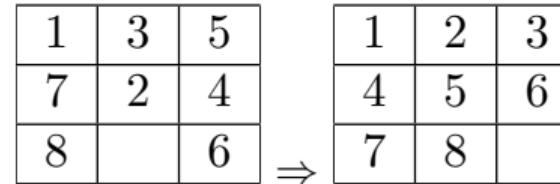


There are 8 possible states:



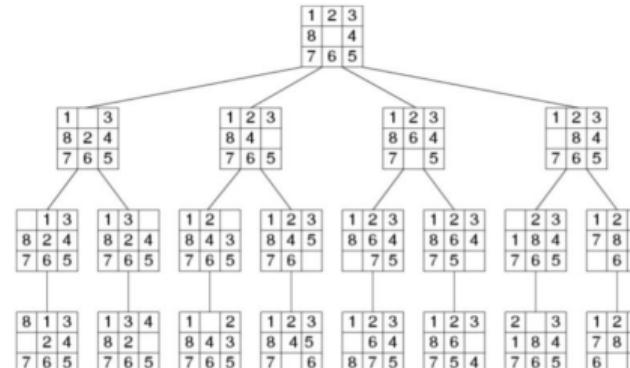
From Russel and Norvig, *Artificial Intelligence: A Modern Approach* (2nd ed.) 2002

## The 8-Puzzle:

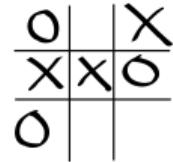


It's just a **search** (each **node** is a **state**):

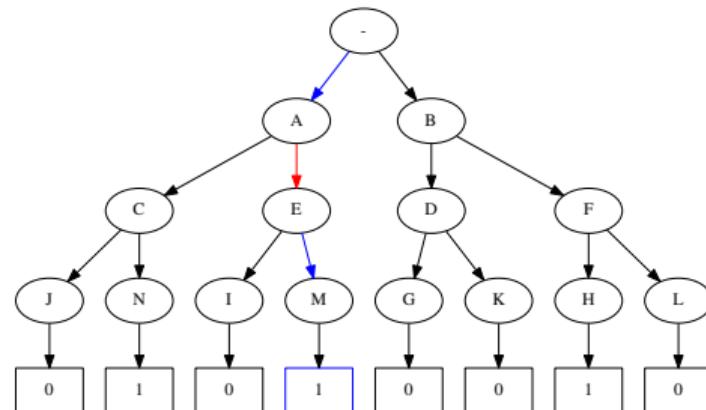
- ① Generate the search tree to goal state
- ② Apply payoff to leaf
- ③ Backup
- ④ Choose the best branch



# Adversaries



What about the other player (**adversary**)? Can be made deterministic via **minimax algorithm**,



where the adversary is the '**min**'. Assumptions:

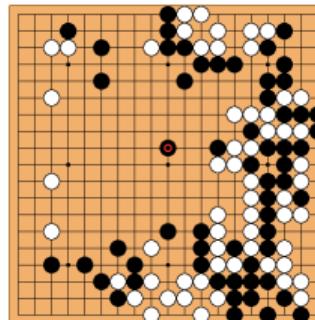
# Scaling Up

What if the branching factor is **too large** and/or long path to goal?



Use advanced/approximate-search techniques.

What if branching factor is **really large**?



# Stochastic Environments and MDPs

- 1 General Setup: Agent and Environment; State, Action, Reward
- 2 Policy Search
- 3 Fully Observed Deterministic Environments
- 4 Stochastic Environments and MDPs
- 5 Learning: Return (Gain) and Value

## Stochastic Environments

Deterministic environments: if agent chooses to step forward (action), the agent move forwards (next state). **Stochastic environments**, e.g., real world:

"Oh u got your nice hot coffee and feelin good? Here let me fuck that up for you."  
-The Universe



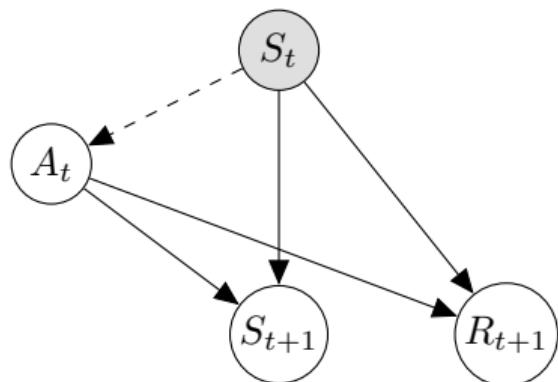
Image credits: the internet

The agent takes action  $a$  from state  $s$  and arrives to state  $s'$  with probability  $s' \sim p(\cdot | s, a)$ . A stochastic adversary makes your environment stochastic too!

# A Markov Decision Process (MDP): Model of the Environment

In stochastic environments, we cannot choose the next state deterministically.

- $\mathcal{S}$  state space,  $s \in \mathcal{S}$
- $\mathcal{A}$  action space,  $a \in \mathcal{A}(s)$
- $\mathcal{R}$  reward space,  $r(s, a) \in \mathcal{R}$
- $r(s, a)$  reward function
- $p(s' | s, a)$  transition dynamics
- $\pi(a | s)$  policy of the agent



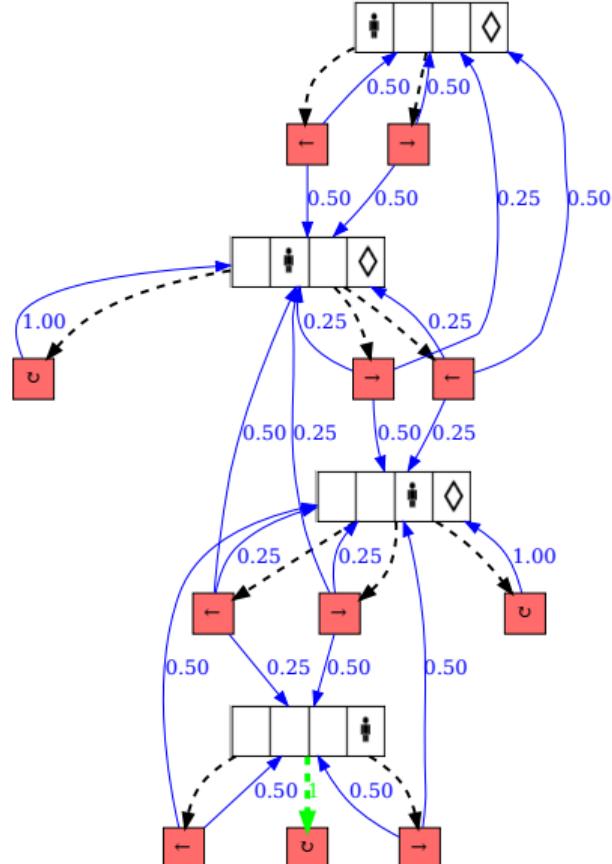
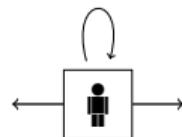
In Reinforcement Learning we do not necessarily know this model.

# Slippery Floor Environment

- state  $s \in \{1, 2, 3, 4\}$
- action  $a \in \{\leftarrow, \circlearrowright, \rightarrow\}$
- reward function  $r(4, \circlearrowright) = 1$   
else  $r(s, a) = 0$  for other values of  $s, a$
- transition dynamics  $s' \sim p(\cdot | s, a)$

e.g.,  $p(\boxed{\text{---} \mid \text{---} \mid \diamond} \mid \boxed{\text{---} \mid \text{---} \mid \diamond}, \leftarrow) = 0.5$

The floor is slippery!! And the lights are off, agent only sees 's':



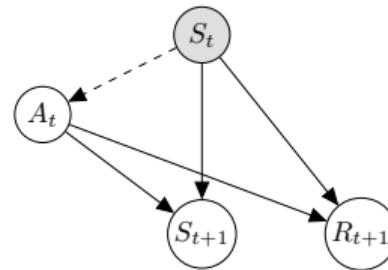
# The Markov Property

MDPs provide a framework for modeling **sequential decision making** in **stochastic** environments. Most reinforcement learning problems can be framed as an MDP.

- M The **Markov** property
- D We add **decisions** (motivated by rewards)
- P A [Markov] **process**

**Markov property:** The effects of action  $a_t$  from state  $s_t$  depend only on that action and state; i.e.,

$$p(s_{t+1} | s_t, a_t) = p(s_{t+1} | s_1, \dots, s_t, a_1, \dots, a_t)$$

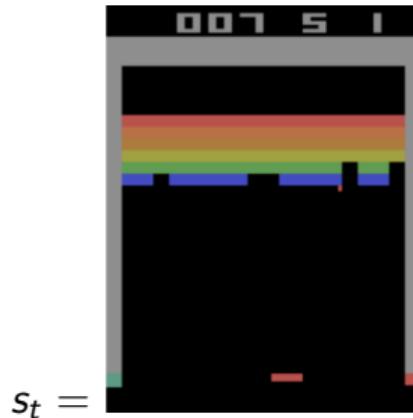


Implications: if Markov, it's possible for the agent to act optimally knowing (observing) state  $s_t$  only.

**Markov** Decision Process:

Does your environment satisfy the **Markov property**?

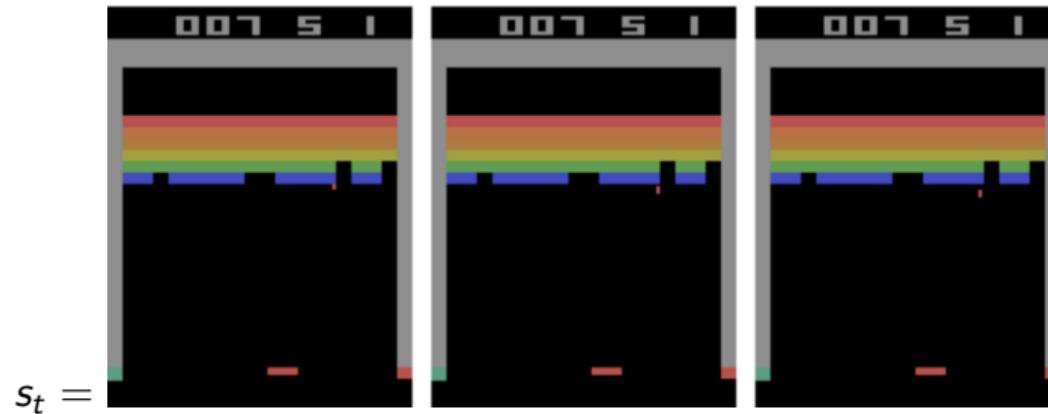
(Is it possible to perform optimally from any state  $s_t$ ?)



$s_t =$

Important consideration! Much of Reinforcement Learning is built around the Markov assumption!

Markovian state:



Other solutions: LSTMs, track the difference across frames, . . . .

i.e., if your decision process is not Markov, you can make it one (not necessarily easily)!

## Learning: Return (Gain) and Value

- 1 General Setup: Agent and Environment; State, Action, Reward
- 2 Policy Search
- 3 Fully Observed Deterministic Environments
- 4 Stochastic Environments and MDPs
- 5 Learning: Return (Gain) and Value

## Reinforcement Learning: Produce a Policy for an Environment

The output of reinforcement learning is a **policy**

$$a_t = \pi(s_t)$$

(for any state  $s_t \in \mathcal{S}$ ) – but *what* policy?

## Reinforcement Learning: Produce a Policy for an Environment

The output of reinforcement learning is a **policy**

$$a_t = \pi(s_t)$$

(for any state  $s_t \in \mathcal{S}$ ) – but **what** policy?

The best policy should take the **best action** from the **current state**;

$$a_t^* = \pi(s_t)$$

i.e., the action to **to optimize** (**this is the important question!**) ??? ...

to optimize what?



## The Gain (Finite Scenario)

The **gain** (aka **return**<sup>1</sup> or sum of future rewards) at step  $t$  is

$$\begin{aligned} G_t &= \sum_{i=t}^T R_{i+1} \\ &= R_{t+1} + R_{t+2} + \dots + R_T \end{aligned}$$

i.e., the **sum of rewards** of the episode; it indicates the **value at current time  $t$** .



---

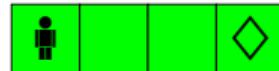
<sup>1</sup>We call it  $G$  the gain to avoid confusion with  $R$  the reward

## The Gain (Finite Scenario)

The **gain** (aka **return**<sup>1</sup> or sum of future rewards) at step  $t$  is

$$\begin{aligned} G_t &= \sum_{i=t}^T R_{i+1} \\ &= R_{t+1} + R_{t+2} + \dots + R_T \end{aligned}$$

i.e., the **sum of rewards** of the **episode**; it indicates the **value at current time  $t$** .



(green for gain). But: when  $T = \infty$ ? And  $r_{t+1} = 1$  vs  $r_{t+1000} = 1$ ?

---

<sup>1</sup>We call it  $G$  the gain to avoid confusion with  $R$  the reward

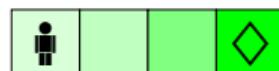
## The Return/Gain (Infinite Scenario)

The **return** (aka **gain**) at step  $t$  is

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (1)$$

$$= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \quad (2)$$

with **discount factor**  $\gamma \in (0, 1)$  which indicates the **relative value of closer rewards**.



Gain  $G$  from each possible  $s$  of optimal agent

Task: The agent should take actions  $A_t$  to maximize  $G_t$ ! **The policy is implicit!**

## The Return/Gain (Infinite Scenario)

The **return** (aka **gain**) at step  $t$  is

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (1)$$

$$= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \quad (2)$$

with **discount factor**  $\gamma \in (0, 1)$  which indicates the **relative value of closer rewards**.



Gain  $G$  from each possible  $s$  of optimal agent

Task: The agent should take actions  $A_t$  to maximize  $G_t$ ! **The policy is implicit!**

But

- $G_t$  is from the future! The future is uncertain!
- $G_t$  inherits the randomness (**uncertainty**) from environment *and* agent!

## The Return/Gain (Infinite Scenario)

The **return** (aka **gain**) at step  $t$  is

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (1)$$

$$= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \quad (2)$$

with **discount factor**  $\gamma \in (0, 1)$  which indicates the **relative value of closer rewards**.



Gain  $G$  from each possible  $s$  of optimal agent

Task: The agent should take actions  $A_t$  to maximize  $G_t$ ! **The policy is implicit!**

But

- $G_t$  is from the future! The future is uncertain!
- $G_t$  inherits the randomness (**uncertainty**) from environment *and* agent!

Therefore: We should maximize expected gain  $\mathbb{E}[G_t]$ , also called *value*.

## Value = Expected Gain

Reinforcement learning is about maximising value:

$$\begin{aligned}\pi_* &= \max_{\pi} \mathbb{E}[G_t] \\ a_t^* &= \pi_*(s_t)\end{aligned}$$

= minimising expected loss = the general setting of machine learning!

Remarks:

- gain  $G_t$  calculated according to discount factor  $\gamma$  and from time  $t$
- expectation (randomness)  $\mathbb{E}$  from our policy  $\pi$  and the environment dynamics  $p$

## Motivation and Role of the Discount Factor $\gamma$

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots + \gamma^{T-1} R_T$$

- Reward today is better than reward tomorrow
- The future is increasingly uncertain (difficult to represent)
- The agent may die at any time
- $\gamma \in (0, 1)$  determines the **present value of future rewards**
- The reward received  $k$  time steps in the future is worth only  $\gamma^{k-1}$  times what it would be worth if it were received now
- If  $\gamma$  is close to 0, the agent is ‘myopic’: maximize immediate rewards
- If  $\gamma$  approaches 1, the agent is more ‘far-sighted’, takes future rewards into account more strongly (even at the cost of missing out on short-term rewards)
- In practice:  $\gamma$  can have a significant impact on agent performance

# The Value Function

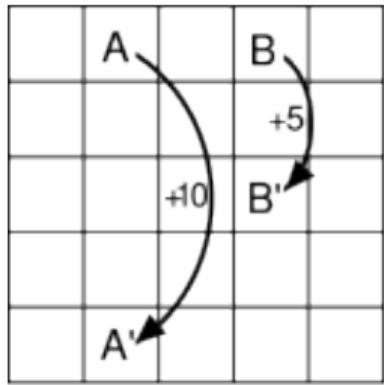
The value function (aka state-value function),

$$V^\pi(s) = \mathbb{E}[G_t | S_t = s]$$

maps a state to a value; which is the expected return from that state following policy  $\pi$  interacting with the environment.

We may think of a vector of  $|\mathcal{S}|$  values; if  $\mathcal{S} = \{s_1, s_2, s_3, s_4\}$ :

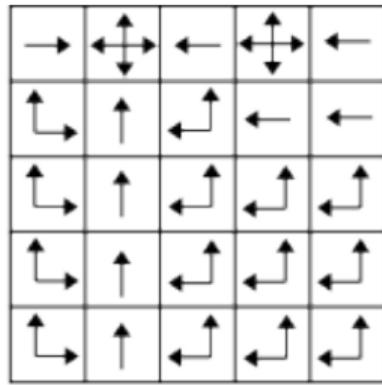
	$s_1$	$s_2$	$s_3$	$s_4$
$V^\pi$				



a) gridworld

22.0	24.4	22.0	19.4	17.5
19.8	22.0	19.8	17.8	16.0
17.8	19.8	17.8	16.0	14.4
16.0	17.8	16.0	14.4	13.0
14.4	16.0	14.4	13.0	11.7

b)  $V^*$



c)  $\pi^*$

Ref: Sutton and Barto. *Reinforcement Learning: An Introduction*, 2018.

# A V-function and Policy for the FROZEN LAKE Environment:

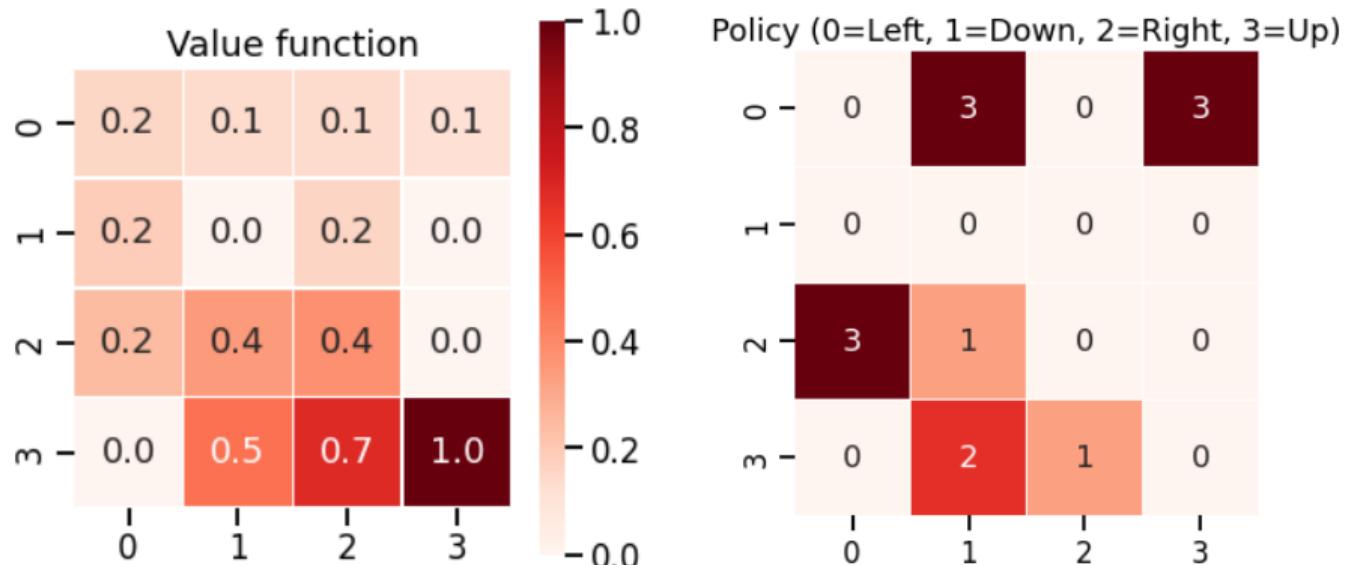
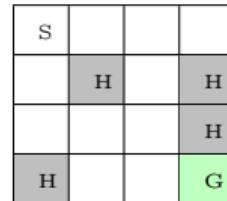


Image credits: Jérémie Decock

## The Action-Value Function

The **action-value function**,

$$Q^\pi(s, a) = \mathbb{E}[G_t \mid S_t = s, A_t = a]$$

maps a state *and action* to a value.

We may think of a table of  $|\mathcal{S}| \times |\mathcal{A}|$  values; with  $\mathcal{A} = \{a_1, a_2\}$ :

$Q^\pi$	$a_1$	$a_2$
$s_1$		
$s_2$		
$s_3$		
$s_4$		

The advantage of  $Q$  over  $V$  is that it can tell us what actions are available from  $s$ .  
The **greedy policy** is explicit:

$$a_t = \pi(s_t) = \operatorname{argmax}_{a \in \mathcal{A}} Q(s_t, a)$$

## Bellman Equation for $V^\pi$ (Preliminaries)

Recall (from Slide 41) the **value function**:

$$V^\pi(s) = \mathbb{E}[G_t | S_t = s] \quad (3)$$

A recursion trick we will need:

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots \\ &= R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots) \\ &= R_{t+1} + \gamma G_{t+1} \quad \text{▷ Recursion!} \end{aligned} \quad (4)$$

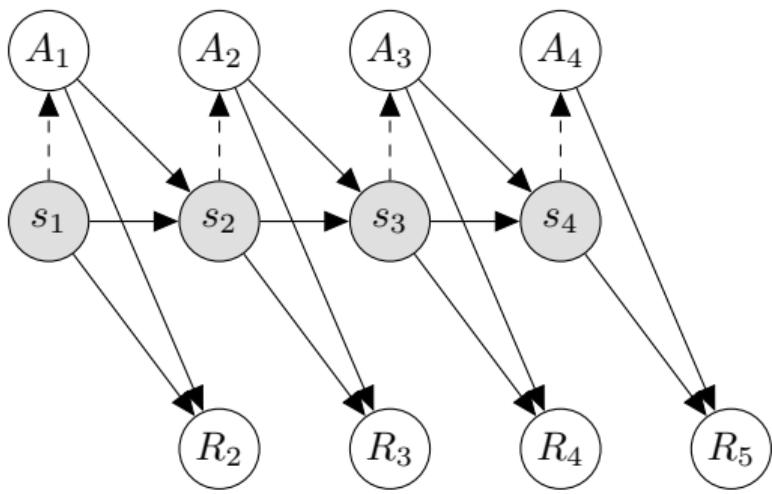
A reminder on conditional probability and expectations and marginalizing out:

$$\mathbb{E}_{A,R,S'}[g(A, R, S') | S = s] = \sum_{a,r,s'} g(a, r, s') p(a, r, s' | s)$$

The expectation is with respect to  $A, S$ .

Total law of expectation:

$$\mathbb{E}[G] = \mathbb{E}_S[\mathbb{E}[G|S]]$$



## Bellman Equation $V^\pi$ (Derivation)

From an MDP  $(\mathcal{S}, \mathcal{A}, \mathcal{R}, p, \gamma)$ :  $\forall s \in \mathcal{S}$ , under policy  $\pi$ :

$$V^\pi(s) = \mathbb{E}[G_t | S_t = s] \quad \triangleright \text{Cond. expectation - cf. Eq. (3)} \quad (5)$$

$$= \mathbb{E}[R_{t+1} + \gamma G_{t+1} | S_t = s] \quad \triangleright \text{Recursion; from Eq. (4)}$$

$$= \mathbb{E}[R_{t+1} | S_t = s] + \gamma \mathbb{E}[G_{t+1} | S_{t+1} = s']$$

$$= \mathbb{E}_{a_t \sim \pi(\cdot | s_t)} [r(s_t, a_t)] + \gamma \mathbb{E}[G_{t+1} | S_{t+1} = s']$$

$$= \sum_{a, s'} \underbrace{p(a | s)p(s' | s, a)}_{p(a, s' | s)} [r(s, a) + \gamma \mathbb{E}_\pi [G_{t+1} | S_{t+1} = s']] \quad \triangleright \text{Marginalise out}$$

$$= \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} p(s' | s, a) \left[ r(s, a) + \gamma \underbrace{\mathbb{E}_\pi [G_{t+1} | S_{t+1} = s']}_{V^\pi(s')} \right] \quad (6)$$

Where Eq. (6) completes the recursion via Eq. (5).

**Why is this important?**

## Bellman Equation $V^\pi$ (Derivation)

From an MDP  $(\mathcal{S}, \mathcal{A}, \mathcal{R}, p, \gamma)$ :  $\forall s \in \mathcal{S}$ , under policy  $\pi$ :

$$V^\pi(s) = \mathbb{E}[G_t | S_t = s] \quad \triangleright \text{Cond. expectation - cf. Eq. (3)} \quad (5)$$

$$= \mathbb{E}[R_{t+1} + \gamma G_{t+1} | S_t = s] \quad \triangleright \text{Recursion; from Eq. (4)}$$

$$= \mathbb{E}[R_{t+1} | S_t = s] + \gamma \mathbb{E}[G_{t+1} | S_{t+1} = s']$$

$$= \mathbb{E}_{a_t \sim \pi(\cdot | s_t)} [r(s_t, a_t)] + \gamma \mathbb{E}[G_{t+1} | S_{t+1} = s']$$

$$= \sum_{a, s'} \underbrace{p(a | s)p(s' | s, a)}_{p(a, s' | s)} [r(s, a) + \gamma \mathbb{E}_\pi [G_{t+1} | S_{t+1} = s']] \quad \triangleright \text{Marginalise out}$$

$$= \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} p(s' | s, a) \left[ r(s, a) + \gamma \underbrace{\mathbb{E}_\pi [G_{t+1} | S_{t+1} = s']}_{V^\pi(s')} \right] \quad (6)$$

Where Eq. (6) completes the recursion via Eq. (5).

**Why is this important?**

A [closed form] solution for **expected return**/value!

## Bellman Equation $V^\pi$ (Derivation)

From an MDP  $(\mathcal{S}, \mathcal{A}, \mathcal{R}, p, \gamma)$ :  $\forall s \in \mathcal{S}$ , under policy  $\pi$ :

$$V^\pi(s) = \mathbb{E}[G_t | S_t = s] \quad \triangleright \text{Cond. expectation - cf. Eq. (3)} \quad (5)$$

$$= \mathbb{E}[R_{t+1} + \gamma G_{t+1} | S_t = s] \quad \triangleright \text{Recursion; from Eq. (4)}$$

$$= \mathbb{E}[R_{t+1} | S_t = s] + \gamma \mathbb{E}[G_{t+1} | S_{t+1} = s']$$

$$= \mathbb{E}_{a_t \sim \pi(\cdot | s_t)} [r(s_t, a_t)] + \gamma \mathbb{E}[G_{t+1} | S_{t+1} = s']$$

$$= \sum_{a, s'} \underbrace{p(a | s)p(s' | s, a)}_{p(a, s' | s)} [r(s, a) + \gamma \mathbb{E}_\pi [G_{t+1} | S_{t+1} = s']] \quad \triangleright \text{Marginalise out}$$

$$= \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} p(s' | s, a) \left[ r(s, a) + \gamma \underbrace{\mathbb{E}_\pi [G_{t+1} | S_{t+1} = s']}_{V^\pi(s')} \right] \quad (6)$$

Where Eq. (6) completes the recursion via Eq. (5).

**Why is this important?**

A [closed form] solution for **expected return**/value!

**However:** what we really want (goal of reinforcement learning) is  $\pi_*$ .

$$V \approx V^\pi$$

Recall:

	$s_1$	$s_2$	$s_3$	$s_4$
$V^\pi$				

But **what are the values?** Where do we get them from?

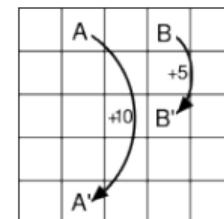
**We need to run policy  $\pi$  to find out!**

But **what is policy  $\pi$ ?** Where do we get it from?

**We need to calculate  $V^\pi$  to find out!**

Notation: If  $V^\pi$  provides  $\pi$ , then what we really want is  $v^*$ , it provides the *optimal* policy  $\pi^*$ .

	$s_1$	$s_2$	$s_3$	$s_4$
$v^*$				



a) gridworld

22.0	24.4	22.0	19.4	17.5
19.8	22.0	19.8	17.8	16.0
17.8	19.8	17.8	16.0	14.4
16.0	17.8	16.0	14.4	13.0
14.4	16.0	14.4	13.0	11.7

b)  $V^*$

→	↔	←	↔	←
↳	↑	↖	↗	↙
↑	↗	↖	↑	↖
↖	↑	↖	↗	↑
↑	↗	↖	↑	↖

c)  $\pi^*$

## Bellman Optimality Equation $v^*$ and $\pi^*$

So: we have an expression for  $V^\pi(s)$ .

The **optimal policy** (and the value function associated with it):

$$\pi^*(s) = \operatorname{argmax}_\pi V^\pi(s) \quad \forall s \in \mathcal{S}$$

$$v^*(s) = \max_\pi V^\pi(s) \quad \forall s \in \mathcal{S}$$

The value of state  $s$ , following policy  $\pi$  (note the **greedy policy**) is

$$V^\pi(s) = Q^\pi(s, \pi(s)) = \max_{a \in \mathcal{A}(s)} Q^\pi(s, a)$$

It means that (from Eq. (6),  $V^\pi(s)$ ):

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

The action-value function  $Q$  **shares the same optimal policy**  $\pi^*$ .

$$\begin{aligned}
v^*(s) &= \max_{a \in \mathcal{A}(s)} q^{\pi^*}(s, a) \\
&= \dots \\
&= \max_{a \in \mathcal{A}(s)} \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} p(s', r | s, a) [r + \gamma v^*(s')] \\
\pi^*(s) &= \operatorname{argmax}_{a \in \mathcal{A}(s)} q^{\pi^*}(s, a)
\end{aligned}$$

Note, again: the **max** and the **argmax**.

$$\begin{aligned}
v^*(s) &= \max_{a \in \mathcal{A}(s)} q^{\pi^*}(s, a) \\
&= \dots \\
&= \max_{a \in \mathcal{A}(s)} \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} p(s', r | s, a) [r + \gamma v^*(s')] \\
\pi^*(s) &= \operatorname{argmax}_{a \in \mathcal{A}(s)} q^{\pi^*}(s, a)
\end{aligned}$$

Note, again: the **max** and the **argmax**.

Recall our main objective: how to get  $\pi^*$ ? (policy which takes optimal actions)

# Generalized Policy Iteration

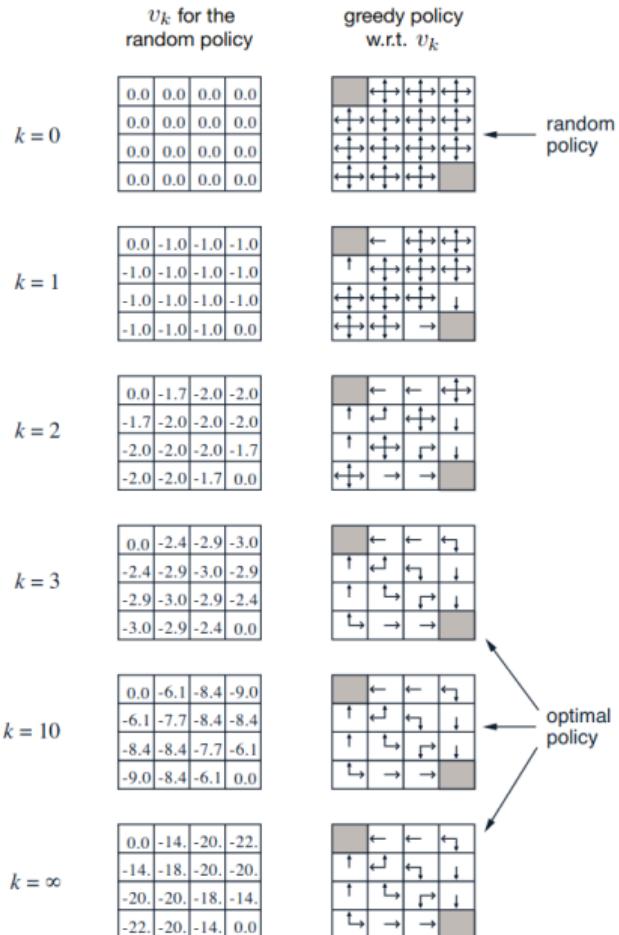
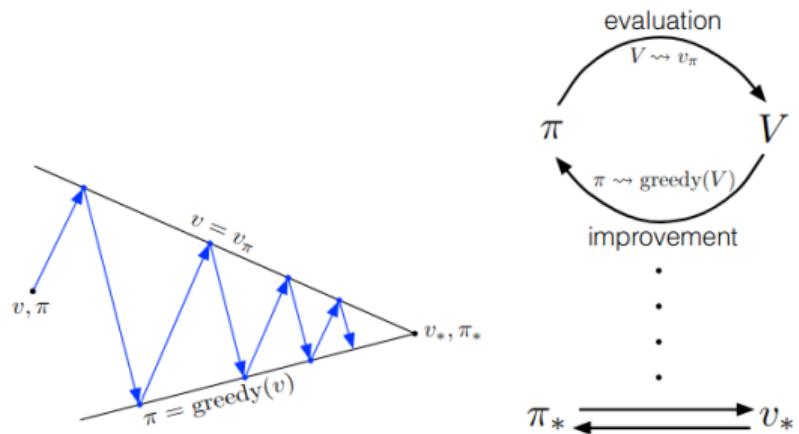
The bellman equation is typically solved iteratively.

Goals:

- ① make the value function consistent with the current policy, and
- ② make the policy **greedy** with respect to the current value function

Iterative solutions:

- Value iteration
- Policy iteration
- Q-learning
- Sarsa
- ...



# Exploration vs Exploitation

From any state  $s$ , at time  $t$ , we want:

$$\underset{\pi}{\operatorname{argmax}} \mathbb{E}_{s' \sim p(s), a \sim \pi(s)}[G_t]$$

In **reinforcement learning**, we **do not know** the underlying system dynamics ( $p$ ) of the MDP, we need to learn it through interaction with the environment.

**Easy:** run  $\pi$ , gather data  $(s, a, r', s')$ , calculate/approximate  $\mathbb{E}$ .

**Hard:** how to run  $\pi$ , i.e., which  $\pi$ ?

Exploration  $\neq$  Exploitation

## A Monte Carlo View

- ① Play many episodes with a policy

$$a_t \sim \pi(\cdot | s_t)$$

and record the **gain** from each  $(s_t, a_t)$ -pair

- ② Use these samples to approximate the expectation:

$$Q^\pi(s, a) = \frac{1}{n} \sum_{i=1}^n g_{a|s}^{(i)} \approx \mathbb{E}[G | s, a]$$

- ③ Repeat many times. Then employ the **greedy policy**  $\pi$  on  $Q$ :

$$a_t = \operatorname{argmax}_{a \in \mathcal{A}(s_t)} \pi(a | s_t) = \operatorname{argmax}_{a \in \mathcal{A}(s_t)} Q^\pi(s_t, a)$$

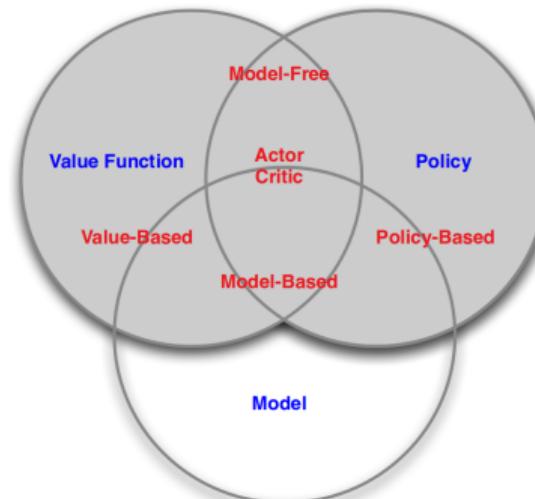
Trade-off: **Exploration vs exploitation.**

# Types of Reinforcement Learning Agents

Value-Based or Policy-Based ?

- Value Based (value function; policy is implicit)
- Policy Based (no value function)
- Actor Critic (both value function + policy)

Model Free vs Model Based: We can build a model of the environment (or not).



## Summary / Important Concepts

- Agent and Environment
- States, Actions, Reward
- Model of the environment (MDP)
  - Markov property
  - Deterministic vs stochastic
  - Unknown vs fully observable vs partially observable
- Return/Gain (finite horizon vs infinite horizon;  $\gamma$ )
- Policy (deterministic policy, stochastic policy, greedy policy)
- State Value and Action-Value functions
- Bellman Equations
- Exploration vs Exploitation trade-off

Reinforcement learning is difficult (compared to supervised learning), but there are many applications where supervised learning does not apply.

# Deep Reinforcement Learning



[2]

Action space is simple (move left, right, shoot, . . .), but **state space is huge!** Q-table impossible!

# A Deep Representation of the Q-table

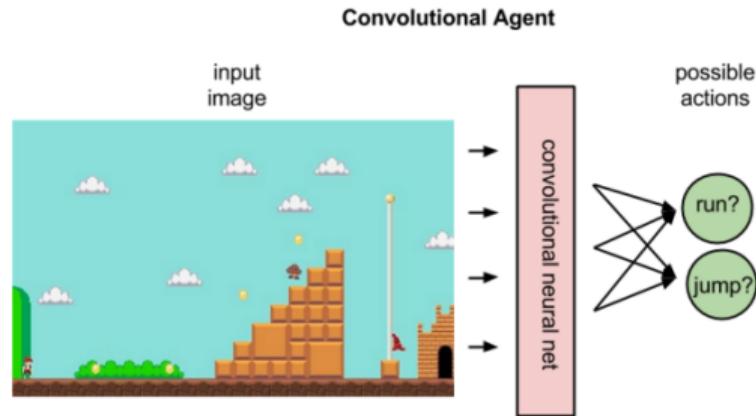
We replace the Q table with a neural network:

$$Q(s, a) = Q_w(s)_a$$

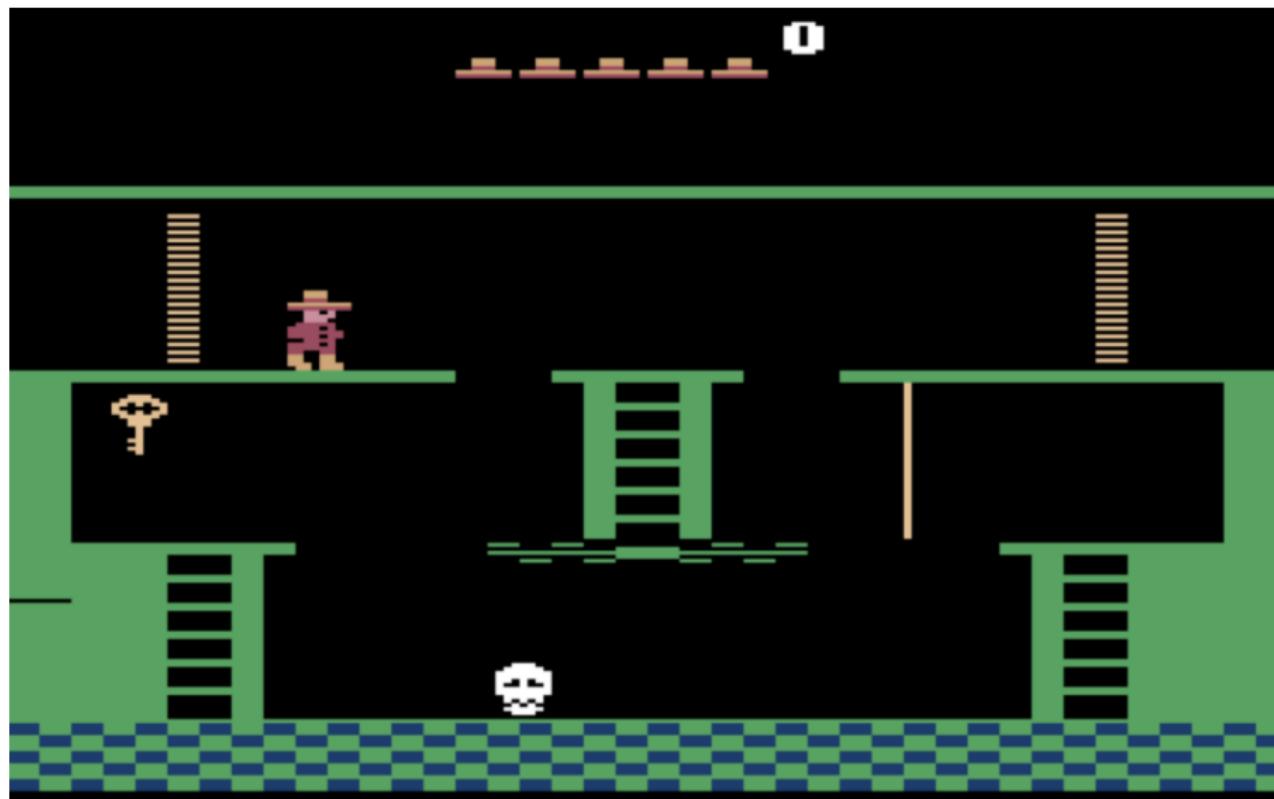
where  $Q_w$  is a (deep) neural network parametrized by  $w$ .

For example,

$$[Q(s, a_1), \dots, Q(s, a_4)] = Q_w(s)$$



# Montezuma's Revenge ('unsolved' by Reinforcement Learning)



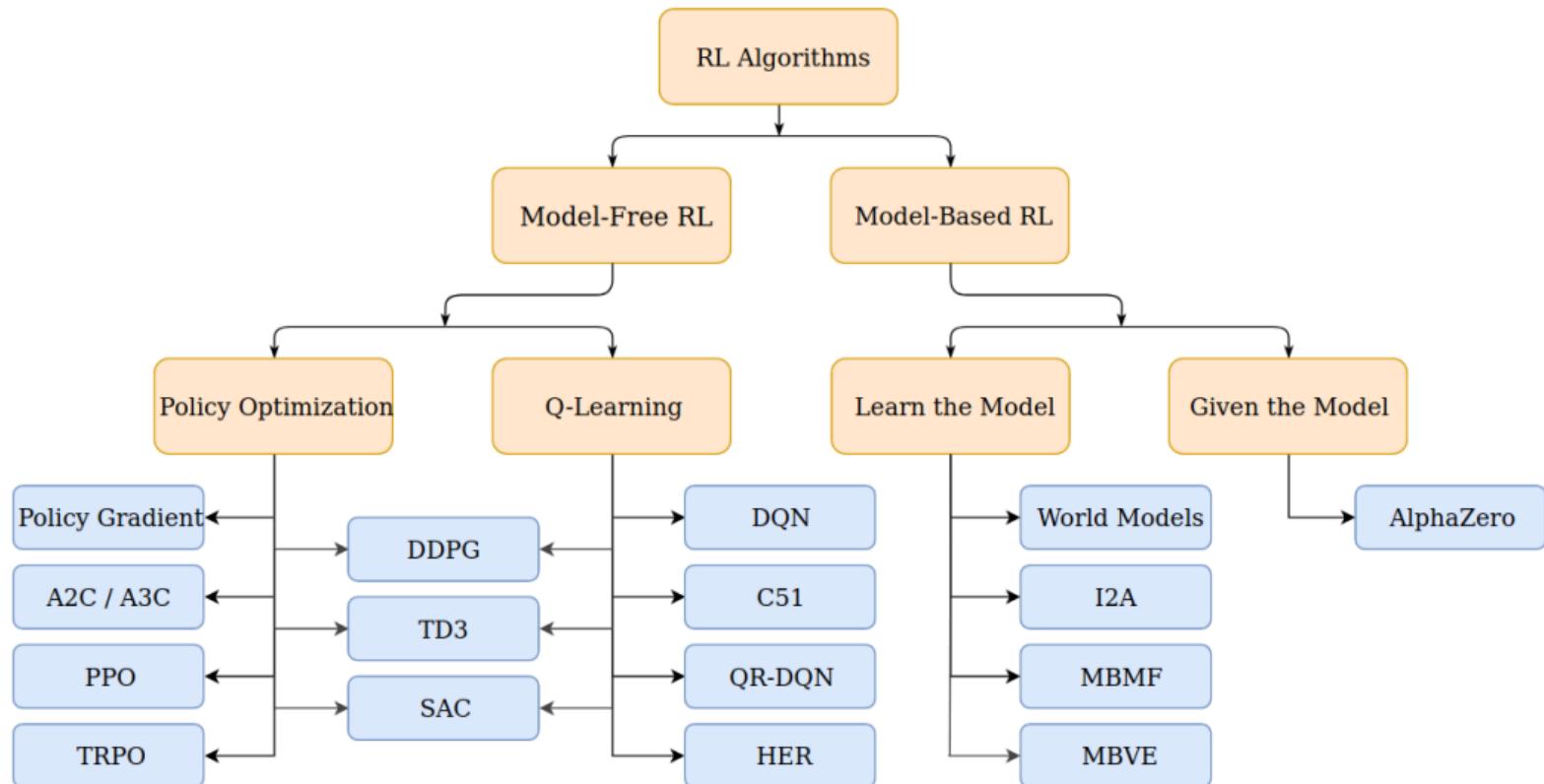
Montezuma's Revenge – Notoriously Difficult for RL Agents [4]

# Modern Challenges

- Learning quickly with **limited samples**; even when **high-dimensional** and **complex**
- **Generalization** and **transfer**: towards a foundation model
- Robustness and **safety** issues (**offline RL** and **model-based RL**)
- Explainable agents, **interpretable** actions
- Delayed/partial/weak observations, and handling the associated uncertainty
- Multi-objective reward fns, multi-agent, alignment
- Legal and **ethical concerns**, and **alignment** issues (agent's reward function vs ours)



# A Taxonomy of Methods



# Advanced Machine Learning and Autonomous Agents

## Reinforcement Learning I



Jesse Read

\Version: