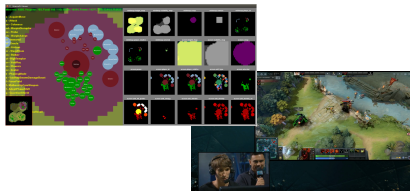
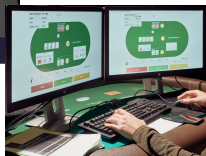


# Intro to deep RL, the concept of generalization and the importance of representation learning

Vincent François-Lavet

19 juin 2024

# Motivation : Overview



# Outline

## Introduction

# Outline

Introduction

Generalization in deep RL

Using self-supervised learning and abstract representations

# Outline

Introduction

Generalization in deep RL

Using self-supervised learning and abstract representations

Abstract representations for reasoning, exploration and transfer learning

# Outline

## Introduction

## Generalization in deep RL

Using self-supervised learning and abstract representations

Abstract representations for reasoning, exploration and transfer learning

A few other challenges for RL (disentanglement of controllable and uncontrollable feature + causal representations)

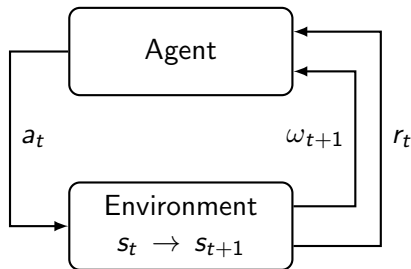
Model-based methods (planning-based techniques)

Model-free techniques

# Introduction

# Objective

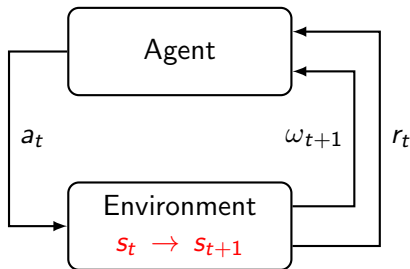
**From experience** in an environment,  
an artificial agent  
should be able to **learn** a sequential decision making task  
in order **to achieve goals**.





# Objective

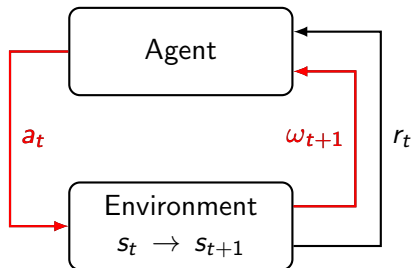
**From experience** in an environment,  
an artificial agent  
should be able to **learn** a sequential decision making task  
in order **to achieve goals**.



transitions  
are usually  
stochastic

# Objective

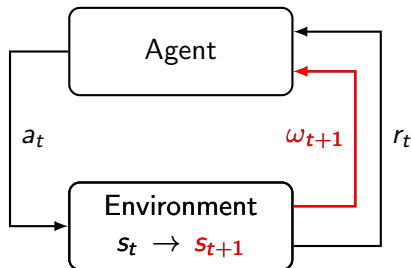
**From experience** in an environment,  
an artificial agent  
should be able to **learn** a sequential decision making task  
in order **to achieve goals**.



Observations and  
actions may be  
high dimensional

# Objective

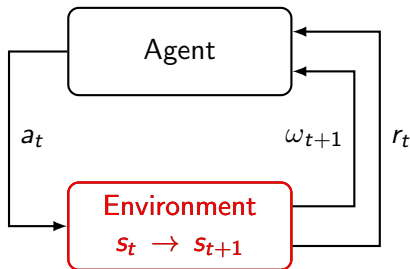
**From experience** in an environment,  
an artificial agent  
should be able to **learn** a sequential decision making task  
in order **to achieve goals**.



Observations may not  
provide full knowledge  
of the underlying  
state :  $\omega_t \neq s_t$

# Objective

**From experience** in an environment,  
an artificial agent  
should be able to **learn** a sequential decision making task  
in order **to achieve goals**.

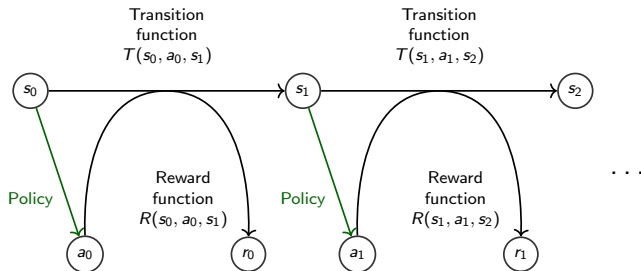


Experience may be constrained  
(e.g., not access to an accurate simulator or limited data)

# Definition of an MDP

An MDP can be defined as a 5-tuple  $(\mathcal{S}, \mathcal{A}, T, R, \gamma)$  where :

- ▶  $\mathcal{S}$  is a finite set of states  $\{1, \dots, N_S\}$ ,
- ▶  $\mathcal{A}$  is a finite set of actions  $\{1, \dots, N_A\}$ ,
- ▶  $T : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{P}(\mathcal{S})$  is the transition function (set of conditional transition probabilities between states),
- ▶  $R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathcal{R}$  is the reward function, where  $\mathcal{R}$  is a continuous set of possible rewards in a range  $R_{max} \in \mathbb{R}^+$  (e.g.,  $[0, R_{max}]$ ),
- ▶  $\gamma \in [0, 1)$  is the discount factor.



# Performance evaluation

In an MDP  $(\mathcal{S}, \mathcal{A}, T, R, \gamma)$ , the discounted expected return  $V^\pi(s) : \mathcal{S} \rightarrow \mathbb{R}$  ( $\pi \in \Pi$ , e.g.,  $\mathcal{S} \rightarrow \mathcal{A}$ ) is defined such that

$$V^\pi(s) = \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k} \mid s_t = s, \pi \right], \quad (1)$$

with  $\gamma \in [0, 1)$ .

From the definition of the (discounted) expected return, the optimal expected return can be defined as

$$V^*(s) = \max_{\pi \in \Pi} V^\pi(s). \quad (2)$$

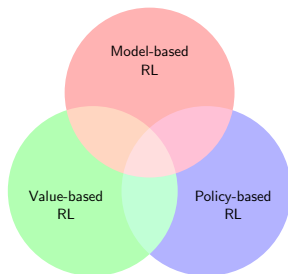
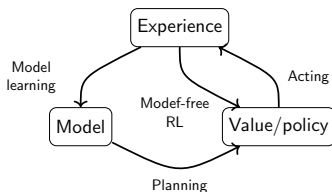
and the optimal policy can be defined as :

$$\pi^*(s) = \operatorname{argmax}_{\pi \in \Pi} V^\pi(s). \quad (3)$$

# Overview of the techniques used for finding the optimal policy $\pi^*$

In general, an RL agent may include one or more of the following components :

- ▶ a model of the environment in conjunction with a planning algorithm.
- ▶ a representation of a value function that provides a prediction of how good is each state or each couple state/action,
- ▶ a direct representation of the policy  $\pi(s)$  or  $\pi(s, a)$ , or



# Strengths and weaknesses of model-based methods

The respective strengths of the model-free versus model-based approaches depend on different factors.

- ✓ For some tasks, the model of the environment is available or can be learned efficiently due to the particular structure of the task.
- ✗ If the agent does not have access to a generative model of the environment, the learned model will have some inaccuracies.
- ✗ A model-based approach requires working in conjunction with a planning algorithm, which is computationally demanding.



# Generalization in deep RL

# Challenges of applying RL to real-world problems

In real-world scenarios, it is often not possible to let an agent interact freely and sufficiently in the actual environment :

- ▶ The agent may not be able to interact with the true environment but only with an inaccurate simulation of it. This is known as the **reality gap**.
- ▶ The agent might have access to only **limited data**. This can be due to safety constraints (robotics, medical trials, etc.), compute constraints or due to limited exogenous data (e.g., weather conditions, trading markets).

# Generalization

In an RL algorithm, *generalization* refers to either

- ▶ the capacity to achieve good performance in an environment where limited data has been gathered, or
- ▶ the capacity to obtain good performance in a related environment. This latter case can be tackled with specific *transfer learning* techniques.

# Overview

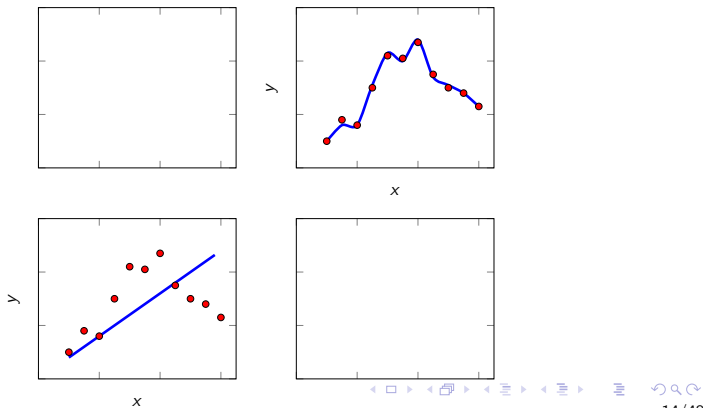
To understand generalization in RL from limited data, we will

- ▶ take inspiration from the bias-variance concept in supervised learning, and
- ▶ introduce the formulation in RL.

We'll then discuss how an agent can have a good generalization in RL (disclaimer : we'll see where deep RL comes in !)

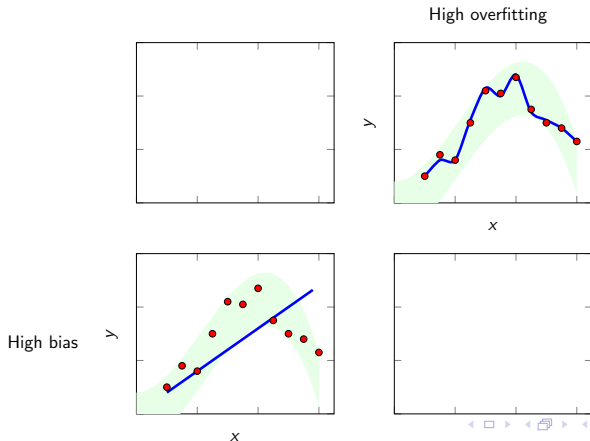
# Bias and overfitting in supervised learning

A supervised learning algorithm can be viewed as a mapping from a dataset  $D_{LS}$  of learning samples  $(x, y) \stackrel{\text{i.i.d.}}{\sim} (X, Y)$  into a predictive model  $f(x | D_{LS})$ .



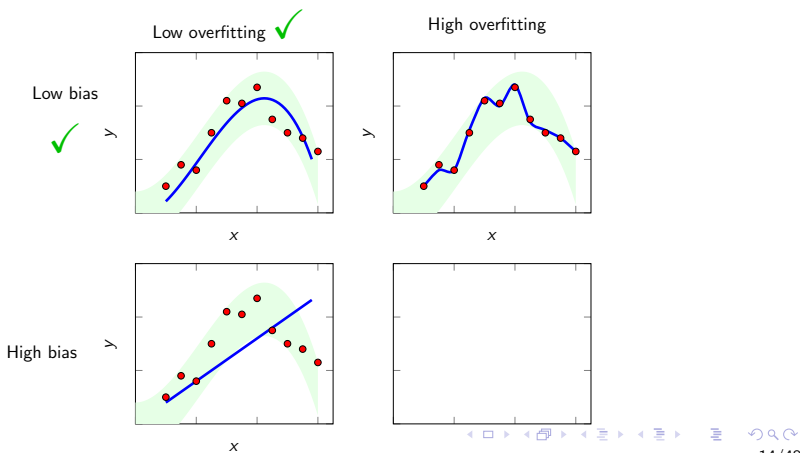
# Bias and overfitting in supervised learning

A supervised learning algorithm can be viewed as a mapping from a dataset  $D_{LS}$  of learning samples  $(x, y) \stackrel{\text{i.i.d.}}{\sim} (X, Y)$  into a predictive model  $f(x | D_{LS})$ .



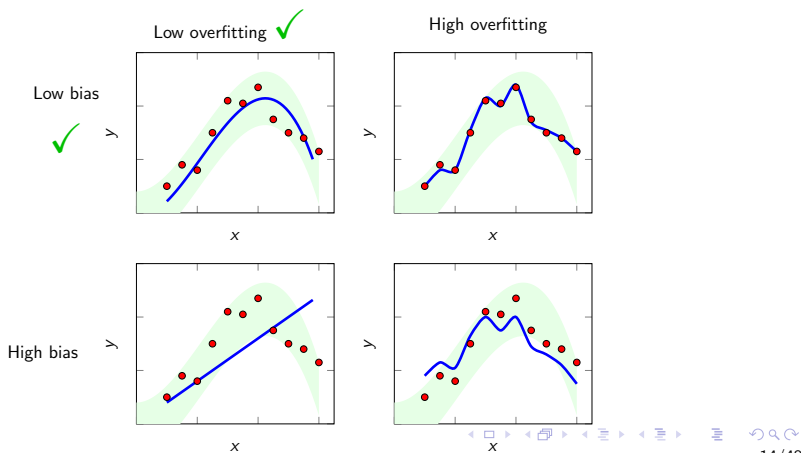
# Bias and overfitting in supervised learning

A supervised learning algorithm can be viewed as a mapping from a dataset  $D_{LS}$  of learning samples  $(x, y) \stackrel{\text{i.i.d.}}{\sim} (X, Y)$  into a predictive model  $f(x | D_{LS})$ .



# Bias and overfitting in supervised learning

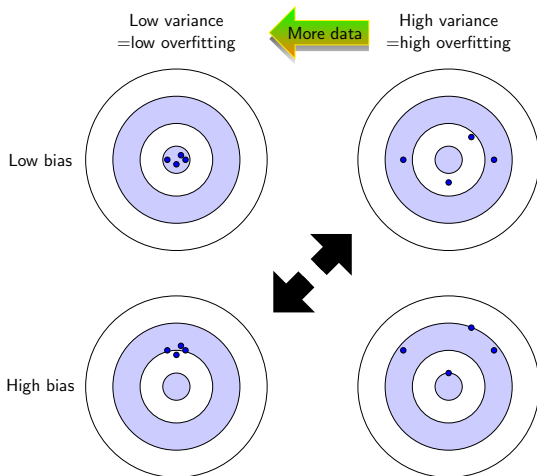
A supervised learning algorithm can be viewed as a mapping from a dataset  $D_{LS}$  of learning samples  $(x, y) \stackrel{\text{i.i.d.}}{\sim} (X, Y)$  into a predictive model  $f(x | D_{LS})$ .





## Bias and overfitting in supervised learning

There are many choices to optimize the learning algorithm and there is usually a tradeoff between the bias and the overfitting terms to reach best solution.



## Bias and overfitting in supervised learning

Assuming a random sampling scheme  $D_{LS} \sim \mathcal{D}_{LS}$ ,  $f(x | D_{LS})$  is a random variable, and so is its average error over the input space.

The expected value of this quantity is given by :

$$I[f] = \mathbb{E}_X \mathbb{E}_{D_{LS}} \mathbb{E}_{Y|X} L(Y, f(X | D_{LS})), \quad (4)$$

where  $L(\cdot, \cdot)$  is the loss function.

## Bias and overfitting in supervised learning

Assuming a random sampling scheme  $D_{LS} \sim \mathcal{D}_{LS}$ ,  $f(x | D_{LS})$  is a random variable, and so is its average error over the input space. The expected value of this quantity is given by :

$$I[f] = \mathbb{E}_X \mathbb{E}_{D_{LS}} \mathbb{E}_{Y|X} L(Y, f(X | D_{LS})), \quad (4)$$

where  $L(\cdot, \cdot)$  is the loss function. If  $L(y, \hat{y}) = (y - \hat{y})^2$ , the error naturally gives the **bias-variance decomposition** :

$$\mathbb{E}_{D_{LS}} \mathbb{E}_{Y|X} (Y - f(X | D_{LS}))^2 = \sigma^2(x) + \text{bias}^2(x), \quad (5)$$

where

$$\text{bias}^2(x) \triangleq (\mathbb{E}_{Y|X}(Y) - \mathbb{E}_{D_{LS}} f(x | D_{LS}))^2,$$

$$\sigma^2(x) \triangleq \underbrace{\mathbb{E}_{Y|X} (Y - \mathbb{E}_{Y|X}(Y))^2}_{\text{Internal variance}} + \underbrace{\mathbb{E}_{D_{LS}} (f(x | D_{LS}) - \mathbb{E}_{D_{LS}} f(x | D_{LS}))^2}_{\text{Parametric variance = overfitting}}.$$

# Bias and overfitting in reinforcement learning

This bias-variance decomposition highlights a tradeoff between

- ▶ an error directly introduced by the learning algorithm (the bias) and
- ▶ an error due to the limited amount of data available (the parametric variance).

# Bias and overfitting in reinforcement learning

Since there is no direct bias-variance decomposition for loss functions other than  $L_2$  loss in supervised learning, there is not an actual “bias-variance” tradeoff in RL.

However, there is still a tradeoff between a sufficiently rich learning algorithm (to reduce the model bias, which is present even when the amount of data would be unlimited) and a learning algorithm not too complex (so as to avoid overfitting to the limited amount of data).

# Bias and overfitting in RL

The *batch* or *offline* algorithm in RL can be seen as mapping a dataset  $D \sim \mathcal{D}$  into a policy  $\pi_D$  (independently of whether the policy comes from a model-based or a model-free approach) :

$$D \rightarrow \pi_D.$$

In an MDP, the suboptimality of the expected return can be decomposed as follows :

$$\begin{aligned} \mathbb{E}_{D \sim \mathcal{D}} [V^{\pi^*}(s) - V^{\pi_D}(s)] &= \underbrace{(V^{\pi^*}(s) - V^{\pi_{D_\infty}}(s))}_{\text{asymptotic bias}} \\ &+ \underbrace{\mathbb{E}_{D \sim \mathcal{D}} [(V^{\pi_{D_\infty}}(s) - V^{\pi_D}(s))]}_{\text{error due to finite size of the dataset } D_s \text{ referred to as } \textit{overfitting}}]. \end{aligned}$$

# How to improve generalization ?

We can improve generalization of RL thanks to the following elements :

- an **abstract representation** that discards non-essential features,
- the **objective function** (e.g., reward shaping, tuning the training discount factor) and
- the **learning algorithm** (type of function approximator and model-free vs model-based).

And of course, if possible :

- improve the dataset (exploration/exploitation dilemma in an online setting)

Questions so far ?

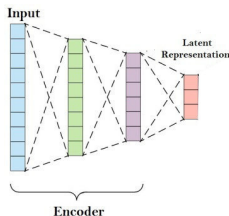


# Using self-supervised learning and abstract representations

# Foreword

## Vocabulary

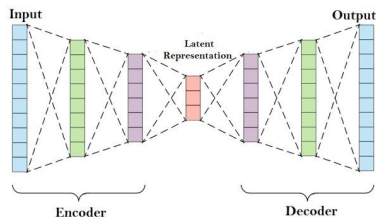
- ▶ An **encoder** is a specific deep learning component that transforms the input (to reduce the dimensionality).
- ▶ An **abstract representation** or **latent representation** is the representation obtained after the input goes through the encoder.



# Foreword

## Vocabulary

- ▶ An **encoder** is a specific deep learning component that transforms the input (to reduce the dimensionality).
- ▶ An **abstract representation** or **latent representation** is the representation obtained after the input goes through the encoder.



# Catcher

This environment has only a few important features

- (i) the position of the paddle and
- (ii) the position of the blocks.

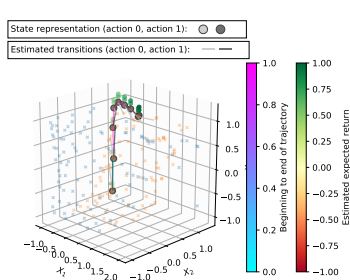
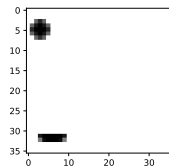


FIGURE –  
Without interpretability loss.

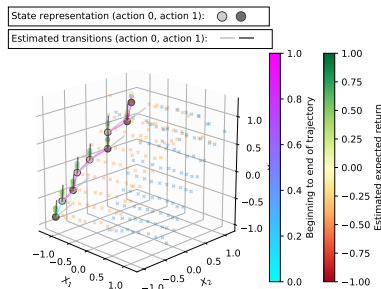


FIGURE – With interpretability loss :  
 $v(a^{(1)}) = (1, 1)$  and  $v(a^{(2)}) = (-1, 1)$ .

# Abstract representations for reasoning, exploration and transfer learning

# Combining model-based and model-free via abstract representations

We are interested to learn both the model and the value function through one abstract representation :

- ▶ it can enforce a good generalization (information bottleneck),
- ▶ planning is computationally efficient,
- ▶ it facilitates interpretation of the decisions taken by the agent,

# Information bottleneck

As opposed to auto encoders, we seek to preserve only *relevant* information and we apply the *Information Bottleneck* (IB) principle to representation learning of state. This gives us the functional that we minimize

$$\mathcal{L} = I[S; X] - \beta I[(X); \{X^*, A^*\}]$$

This corresponds to a trade-off between

- ▶ minimizing the encoding rate  $I[S; X]$  and
- ▶ maximizing the mutual information between the abstract state  $X$  (and reward) and the tuple previous abstract state, previous action  $(X^*, A^*)$ .

# Simple labyrinth

Abstract representation of states for a labyrinth task (without any reward).

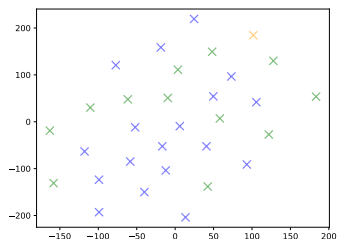
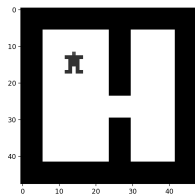


FIGURE – 2D representation using t-SNE (blue represents states where the agent is on the left part, green on the right part and orange in the junction).

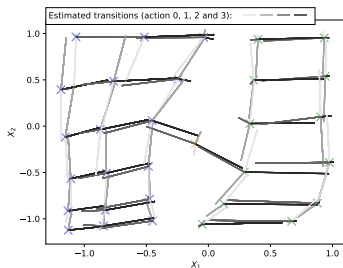


FIGURE – The CRAR agent is able to reconstruct a sensible representation of its environment in 2 dimensions.



# Combined Reinforcement Learning via Abstract Representations (CRAR)

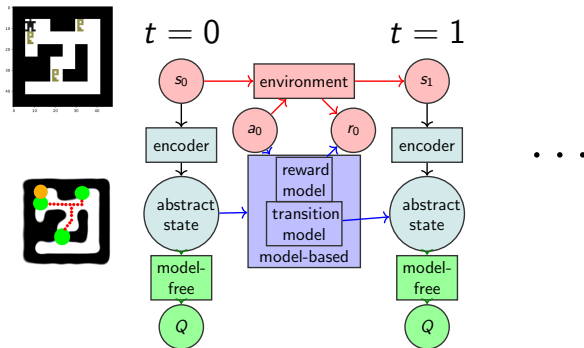


FIGURE – Illustration of the integration of model-based and model-free RL in the CRAR architecture.

The value function and the model are trained via the abstract representation.

# Learning the model

Here is how we learn the internal model :

$$\mathcal{L}_\tau(\theta_e, \theta_\tau) = | (e(s; \theta_e) + \tau(e(s; \theta_e), a; \theta_\tau) - e(s'; \theta_e)) |^2,$$

# Learning the model

Here is how we learn the internal model :

$$\mathcal{L}_\tau(\theta_e, \theta_\tau) = | (e(s; \theta_e) + \tau(e(s; \theta_e), a; \theta_\tau) - e(s'; \theta_e)) |^2,$$

$$\mathcal{L}_\rho(\theta_e, \theta_\rho) = | r - \rho(e(s; \theta_e), a; \theta_\rho) |^2,$$

$$\mathcal{L}_g(\theta_e, \theta_g) = | \gamma - g(e(s; \theta_e), a; \theta_g) |^2 .$$

These losses train the weights of both the encoder and the model-based components.

# Learning the model

Here is how we learn the internal model :

$$\mathcal{L}_\tau(\theta_e, \theta_\tau) = | (e(s; \theta_e) + \tau(e(s; \theta_e), a; \theta_\tau) - e(s'; \theta_e)) |^2,$$

$$\mathcal{L}_\rho(\theta_e, \theta_\rho) = | r - \rho(e(s; \theta_e), a; \theta_\rho) |^2,$$

$$\mathcal{L}_g(\theta_e, \theta_g) = | \gamma - g(e(s; \theta_e), a; \theta_g) |^2 .$$

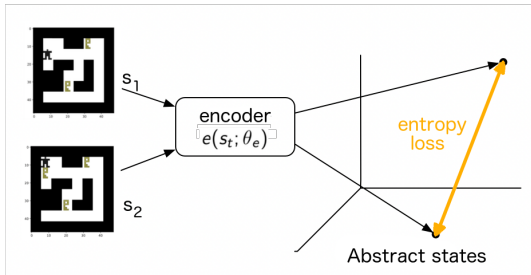
These losses train the weights of both the encoder and the model-based components.

Training of the value function is done with DDQN

When learning the transition function there is a pressure to decrease the amount of information being represented.  
In our model, we introduce an **entropy loss** :

$$\mathcal{L}_{d1}(\theta_e) = \exp(-C_d \|e(s_1; \theta_e) - e(s_2; \theta_e)\|_2),$$

where  $s_1$  and  $s_2$  are random past states of the agent and  $C_d$  is a constant.



## Interpretability

Interpretability can mean that some features of the state representation are distinctly affected by some actions. The following optional loss makes the predicted abstract state change aligned with the chosen embedding vector  $v(a)$  :

$$\mathcal{L}_{\text{interpr}}(\theta_e, \theta_\tau) = -\cos\left(\tau(e(s; \theta_e), a; \theta_\tau)_{0:n}, v(a)\right),$$

where  $\cos$  stands for the cosine similarity.

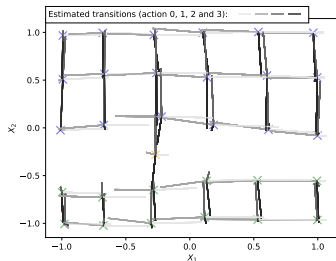


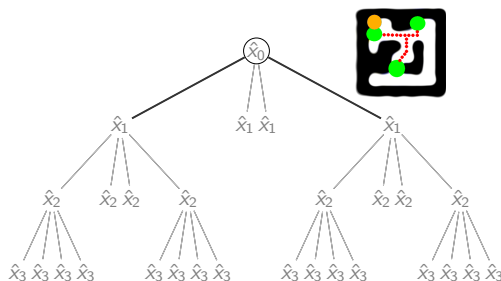
FIGURE – With enforcing  $\mathcal{L}_{\text{interpr}}$  and  $v(a_0) = [1, 0]$

## Planning

The trajectories for some sequence of actions are estimated recursively as follows for any  $t'$  :

$$\hat{x}_{t'} = \begin{cases} e(s_t; \theta_e), & \text{if } t' = t \\ \hat{x}_{t'-1} + \tau(\hat{x}_{t'-1}, a_{t'-1}; \theta_\tau), & \text{if } t' > t \end{cases}$$

A set  $\mathcal{A}^*$  of best potential actions is considered based on  $Q(\hat{x}_t, a; \theta_Q)$  ( $\mathcal{A}^* \subseteq \mathcal{A}$ ).



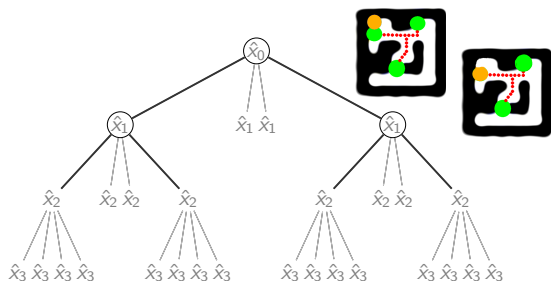
Expansion from  $\hat{x}_0$  until a certain depth.

## Planning

The trajectories for some sequence of actions are estimated recursively as follows for any  $t'$  :

$$\hat{x}_{t'} = \begin{cases} e(s_t; \theta_e), & \text{if } t' = t \\ \hat{x}_{t'-1} + \tau(\hat{x}_{t'-1}, a_{t'-1}; \theta_\tau), & \text{if } t' > t \end{cases}$$

A set  $\mathcal{A}^*$  of best potential actions is considered based on  $Q(\hat{x}_t, a; \theta_Q)$  ( $\mathcal{A}^* \subseteq \mathcal{A}$ ).



Expansion from  $\hat{x}_0$  until a certain depth.

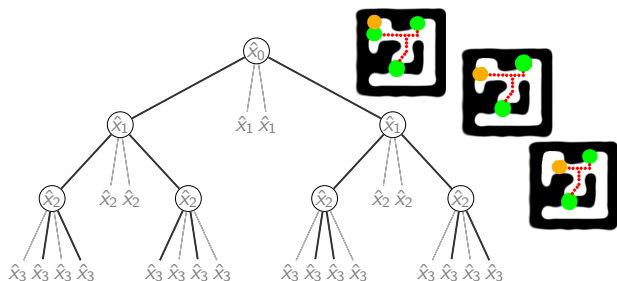


## Planning

The trajectories for some sequence of actions are estimated recursively as follows for any  $t'$  :

$$\hat{x}_{t'} = \begin{cases} e(s_t; \theta_e), & \text{if } t' = t \\ \hat{x}_{t'-1} + \tau(\hat{x}_{t'-1}, a_{t'-1}; \theta_\tau), & \text{if } t' > t \end{cases}$$

A set  $\mathcal{A}^*$  of best potential actions is considered based on  $Q(\hat{x}_t, a; \theta_Q)$  ( $\mathcal{A}^* \subseteq \mathcal{A}$ ).



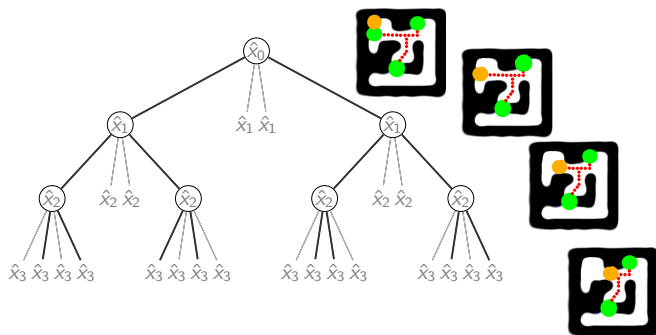
Expansion from  $\hat{x}_0$  until a certain depth.

## Planning

The trajectories for some sequence of actions are estimated recursively as follows for any  $t'$  :

$$\hat{x}_{t'} = \begin{cases} e(s_t; \theta_e), & \text{if } t' = t \\ \hat{x}_{t'-1} + \tau(\hat{x}_{t'-1}, a_{t'-1}; \theta_\tau), & \text{if } t' > t \end{cases}$$

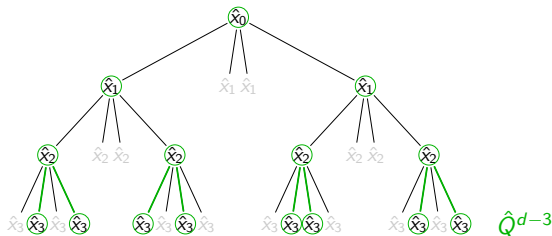
A set  $\mathcal{A}^*$  of best potential actions is considered based on  $Q(\hat{x}_t, a; \theta_Q)$  ( $\mathcal{A}^* \subseteq \mathcal{A}$ ).



Expansion from  $\hat{x}_0$  until a certain depth.

We define recursively the depth- $d$  estimated expected return as

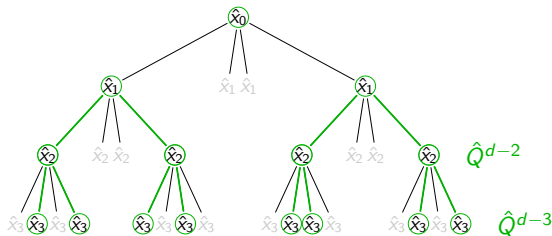
$$\hat{Q}^d(\hat{x}_t, a) = \begin{cases} \rho(\hat{x}_t, a; \theta_\rho) + g(\hat{x}_t, a; \theta_g) \max_{a' \in \mathcal{A}^*} \hat{Q}^{d-1}(\hat{x}_{t+1}, a'), & \text{if } d > 0 \\ Q(\hat{x}_t, a; \theta_k), & \text{if } d = 0 \end{cases}$$



Backup

We define recursively the depth- $d$  estimated expected return as

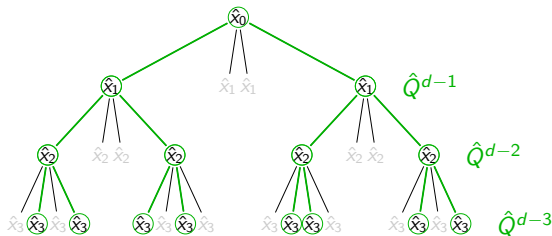
$$\hat{Q}^d(\hat{x}_t, a) = \begin{cases} \rho(\hat{x}_t, a; \theta_\rho) + g(\hat{x}_t, a; \theta_g) \max_{a' \in \mathcal{A}^*} \hat{Q}^{d-1}(\hat{x}_{t+1}, a'), & \text{if } d > 0 \\ Q(\hat{x}_t, a; \theta_k), & \text{if } d = 0 \end{cases}$$



Backup

We define recursively the depth- $d$  estimated expected return as

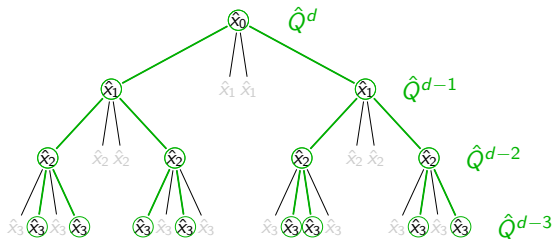
$$\hat{Q}^d(\hat{x}_t, a) = \begin{cases} \rho(\hat{x}_t, a; \theta_\rho) + g(\hat{x}_t, a; \theta_g) \max_{a' \in \mathcal{A}^*} \hat{Q}^{d-1}(\hat{x}_{t+1}, a'), & \text{if } d > 0 \\ Q(\hat{x}_t, a; \theta_k), & \text{if } d = 0 \end{cases}$$



Backup

We define recursively the depth- $d$  estimated expected return as

$$\hat{Q}^d(\hat{x}_t, a) = \begin{cases} \rho(\hat{x}_t, a; \theta_\rho) + g(\hat{x}_t, a; \theta_g) \max_{a' \in \mathcal{A}^*} \hat{Q}^{d-1}(\hat{x}_{t+1}, a'), & \text{if } d > 0 \\ Q(\hat{x}_t, a; \theta_k), & \text{if } d = 0 \end{cases}$$



Backup

## Planning - summary

$$\hat{x}_{t'} = \begin{cases} e(s_t; \theta_e), & \text{if } t' = t \\ \hat{x}_{t'-1} + \tau(\hat{x}_{t'-1}, a_{t'-1}; \theta_\tau), & \text{if } t' > t \end{cases}$$

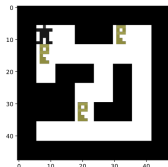
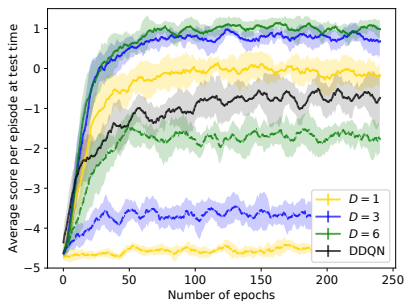
$$\hat{Q}^d(\hat{x}_t, a) = \begin{cases} \rho(\hat{x}_t, a; \theta_\rho) + g(\hat{x}_t, a; \theta_g) \max_{a' \in \mathcal{A}^*} \hat{Q}^{d-1}(\hat{x}_{t+1}, a'), & \text{if } d > 0 \\ Q(\hat{x}_t, a; \theta_k), & \text{if } d = 0 \end{cases}$$

To obtain the action selected at time  $t$ , we use a hyper-parameter  $D \in \mathbb{N}$  and use a simple sum of the Q-values obtained with planning up to a depth  $D$  :

$$Q_{plan}^D(\hat{x}_t, a) = \sum_{d=0}^D \hat{Q}^d(\hat{x}_t, a).$$

The optimal action is given by  $\operatorname{argmax}_{a \in \mathcal{A}} Q_{plan}^D(\hat{x}_t, a)$ .

# Generalization



**FIGURE** – Meta-learning score on a distribution of labyrinths where the training is done with a limited number of transitions obtained by a random policy.  $2 \times 10^5$  tuples,  $\sim 500$  labyrinths.

More details : *Combined Reinforcement Learning via Abstract Representations*, V. Francois-Lavet, Y. Bengio, D. Precup, J. Pineau (AAAI 2019).



## Another important challenge : exploration

- ▶ Undirected exploration (e.g.  $\epsilon$ -greedy)
- ▶ Directed exploration
  - ▶ When rewards are not sparse, a measure of the uncertainty on the value function can be used ;
  - ▶ If sparse rewards or no rewards, some exploration rewards have to be used.

# Exploration

Given a point  $x$  in representation space, we define a **reward function for novelty** that considers the *sparsity* of states around  $x$  - with the average distance between  $x$  and its  $k$ -nearest-neighbors in its visitation history buffer  $\mathcal{B}$  :

$$\hat{\rho}_x(x) = \frac{1}{k} \sum_{i=1}^k d(x, x_i), \quad (6)$$

where  $x$  is a given encoded state,  $k \in \mathbb{Z}^+$ ,  $d(\cdot, \cdot)$  is some distance metric in  $\mathbb{R}^{n_x}$  and  $x_i$  are the  $k$  nearest neighbors (by encoding states in  $\mathcal{B}$  to representational space).

More details : *Novelty Search in representational space for sample efficient exploration*, D. Tao, V. Francois-Lavet, J. Pineau (NeurIPS 2020).

# Exploration

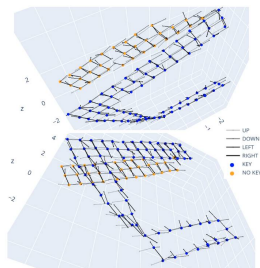
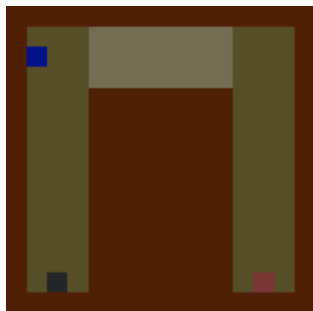


Figure : Multi step environment (left) and the abstract representations of states (right)

## Exploration

This technique can also be used for control tasks such as the double pendulum (acrobot), where only intrinsic rewards allows the agent to solve the task.

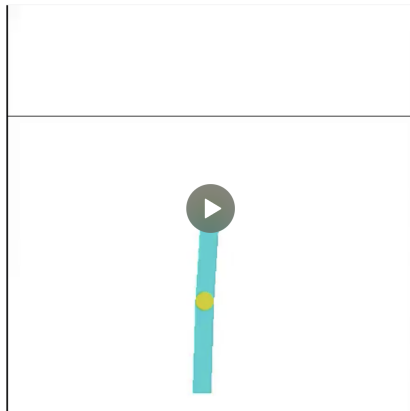
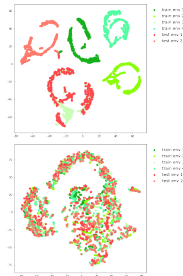
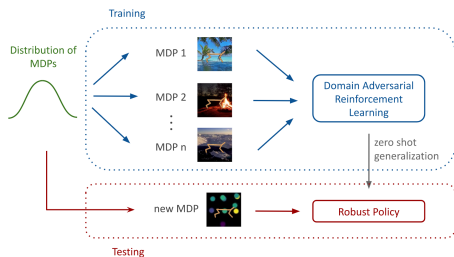


Figure : acrobot.

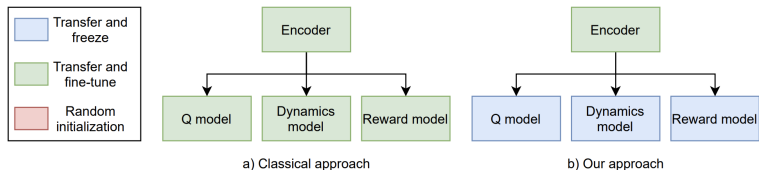
# Transfer learning



**FIGURE** – Set up : the agent is trained in a distribution of MDPs and evaluation is done in new domains with unknown backgrounds.

More details : *Domain adversarial reinforcement learning*, B. Li, V. Francois-Lavet, T. Doan, J. Pineau (2020).

# Component transfer learning

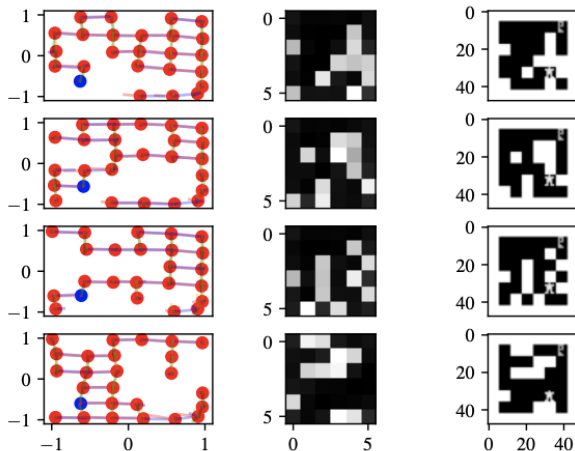


**FIGURE** – Set up : the agent is trained in a distribution of MDPs and evaluation is done in new domains with unknown backgrounds.

More details : *Component Transfer Learning for Deep RL Based on Abstract Representations*, Geoffrey van Driessell, V. Francois-Lavet (2021).

A few other challenges for RL  
(disentanglement of controllable  
and uncontrollable feature + causal  
representations)

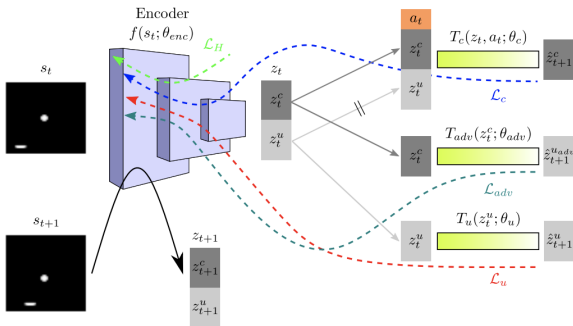
## Disentangled (un-)controllable features



**FIGURE** – The disentangled latent state representations of four different random maze observations. The left column represents the controllable latent representation. The middle column represents the uncontrollable latent representation and the right column is the original state.



# Disentangled (un-)controllable features



**FIGURE** – The disentangled latent state representations of four different random maze observations. The left column represents the controllable latent representation. The middle column represents the uncontrollable latent representation and the right column is the original state.

More details : *Disentangled (Un) Controllable Features*. JE Kooi, M Hoogendoorn, V François-Lavet (2022).

# Causality

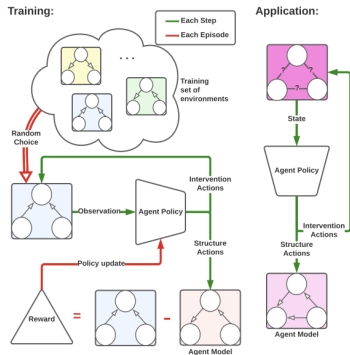


FIGURE – Learning to learn causal graphs.

More details : *A Meta-Reinforcement Learning Algorithm for Causal Discovery.*  
Andreas Sauter, Erman Acar, Vincent François-Lavet (2022)

# Conclusions

## Introduction

## Generalization in deep RL

Using self-supervised learning and abstract representations

Abstract representations for reasoning, exploration and transfer learning

A few other challenges for RL (disentanglement of controllable and uncontrollable feature + causal representations)

Model-based methods (planning-based techniques)

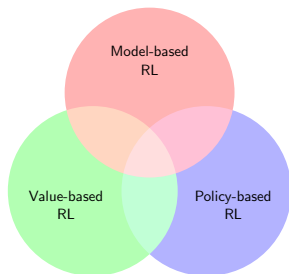
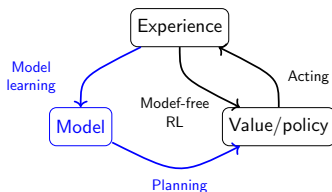
Model-free techniques

Questions ?

# Overview of the techniques used for finding the optimal policy $\pi^*$

In general, an RL agent may include one or more of the following components :

- a model of the environment in conjunction with a planning algorithm.
- ▶ a representation of a value function that provides a prediction of how good is each state or each couple state/action,
- ▶ a direct representation of the policy  $\pi(s)$  or  $\pi(s, a)$ , or



# Model-based methods (planning-based techniques)

# Motivation for planning with tree search

Given that you're playing the crosses, what would be your next move?

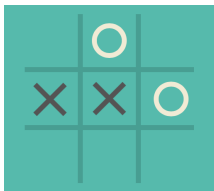


FIGURE – Illustration of a state in the tic-tac-toe game.

→ how did you come up with that choice?

## Monte-Carlo Tree Search methods

The overall idea is to estimate the action with the highest expected return.

$$V^*(s) = Q^*(s, a = \pi^*) = \mathbb{E}_{\pi^*}[r_0 + \gamma r_1 + \dots]$$

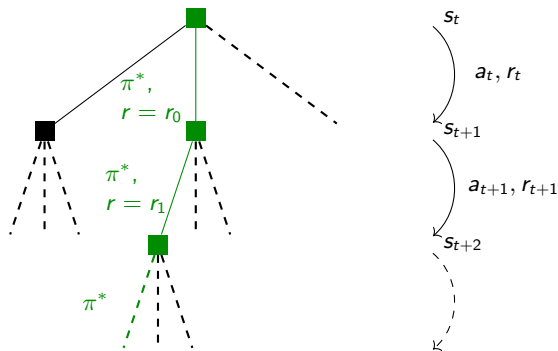


FIGURE – Illustration of model-based planning with tree search.



# Motivation

MCTS algorithms need (only) a generative model of the environment (i.e. model-based) :

$$s_{t+1}, r_t \sim G(s_t, a_t)$$

Advantages :

- ▶ it is possible to obtain samples without having the whole transition function for the model in an explicit form.
- ▶ it can learn a strong policy only where needed (from the current state  $s$ ).
- ▶ it is useful for a sequence of decisions.

# MCTS

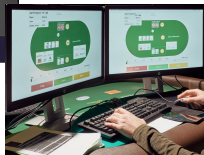
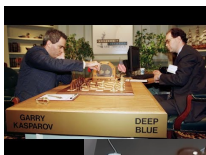
MCTS can converge to the optimal policy (finite action space, finite horizon) from any state  $s$  as long as the generative model is accurate.

However,

- The breath of search grows with the actions space.
- The depth of search grows with the horizon considered.

# Applications

- ▶ Tree search algorithms can be used along with different heuristics as well as model-free deep RL techniques.
- MCTS has been a key part of alpha Go for instance.



# Strengths and weaknesses of model-based methods

The respective strengths of the model-free versus model-based approaches depend on different factors.

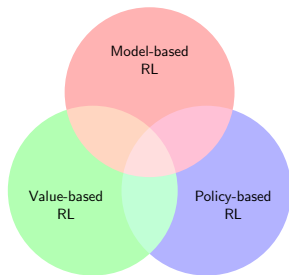
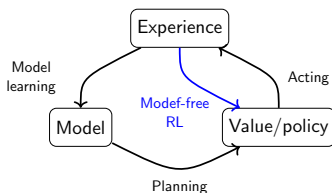
- ✓ For some tasks, the model of the environment is available or can be learned efficiently due to the particular structure of the task.
- ✗ If the agent does not have access to a generative model of the environment, the learned model will have some inaccuracies.
- ✗ A model-based approach requires working in conjunction with a planning algorithm, which is computationally demanding.

# Model-free techniques

# Overview of deep RL

In general, the learning algorithm in RL may include one or more of the following components :

- ▶ a model of the environment in conjunction with a planning algorithm.
- a value function that provides a prediction of how good is each state or each couple state/action (main focus), or
- a direct representation of the policy  $\pi(s)$  or  $\pi(s, a)$



Deep learning has brought its generalization capabilities to RL.

# Convergence Q-learning

Theorem : Given a finite MDP, the Q-learning algorithm given by the update rule

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha_t [r_t + \gamma \max_{a' \in \mathcal{A}} Q_t(s_{t+1}, a') - Q_t(s_t, a_t)],$$

converges w.p.1 to the optimal Q-function as long as

- ▶  $\sum_t \alpha_t = \infty$  and  $\sum_t \alpha_t^2 < \infty$ , and
- ▶ the exploration policy  $\pi$  is such that  $P_\pi[a_t = a | s_t = s] > 0, \forall (s, a)$ .

# Convergence Q-learning

Theorem : Given a **finite MDP**, the Q-learning algorithm given by the update rule

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha_t [r_t + \gamma \max_{a' \in \mathcal{A}} Q_t(s_{t+1}, a') - Q_t(s_t, a_t)],$$

converges w.p.1 to the optimal Q-function as long as

- ▷  $\sum_t \alpha_t = \infty$  and  $\sum_t \alpha_t^2 < \infty$ , and
- ▷ the exploration policy  $\pi$  is such that  $P_\pi[a_t = a | s_t = s] > 0, \forall (s, a)$ .



## Limitations of tabular approaches

A tabular approach fails for large scale problems due to the curse of dimensionality.

- ▶ Robot states with 10 features (e.g. position, speed, angle of joints) discretized into 100 bins  $\rightarrow 100^{10} = 10^{20}$  states.
- ▶ Chess :  $\approx 10^{120}$  states
- ▶ Go :  $\approx 10^{170}$  states

# Limitations of tabular approaches

A tabular approach fails for large scale problems due to the curse of dimensionality.

- ▶ Robot states with 10 features (e.g. position, speed, angle of joints) discretized into 100 bins  $\rightarrow 100^{10} = 10^{20}$  states.
- ▶ Chess :  $\approx 10^{120}$  states
- ▶ Go :  $\approx 10^{170}$  states

## Three problems

- Memory
- Compute time
- No generalization in the limited data context