

An outlook on deep reinforcement learning

Vincent François-Lavet

September 6th, 2019

Outline

Intro to reinforcement learning

Generalisation from limited data in reinforcement learning

How to improve generalization ?

Best practices in deep RL

Combining model-based and model-free via abstract representations

Conclusions

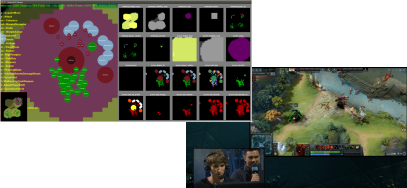
Intro to reinforcement learning

Machine learning relates to the capability of computers to **learn from examples** without following explicitly defined rules.

Three types of machine learning tasks can be described.

- ▶ Supervised learning is the task of inferring a classification or regression from labeled training data.
- ▶ Unsupervised learning is the task used to draw inferences from datasets consisting of input data without labeled responses.
- ▶ Reinforcement learning (RL) is the task concerned with how software agents ought to take actions in an environment in order to achieve some objectives.

Motivation



Motivation

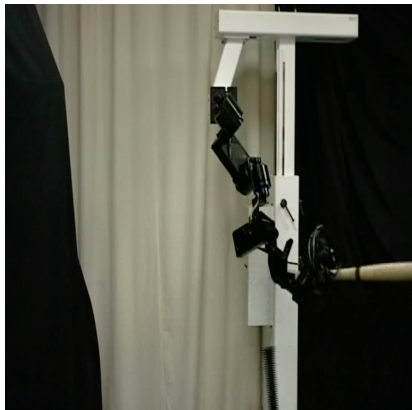
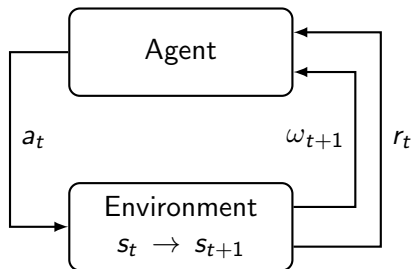


FIGURE – Application in robotics (credits : Jan Peters' team, Darmstadt)

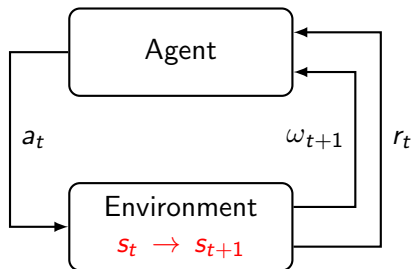
Objective

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



Objective

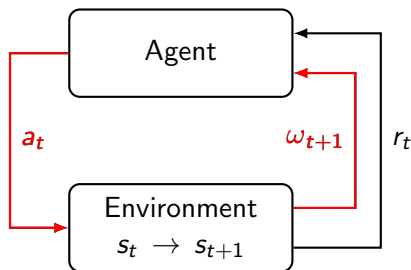
From experience in an environment,
an artificial agent
should be able to **learn** a sequential decision making task
in order **to achieve goals**.



transitions
are usually
stochastic

Objective

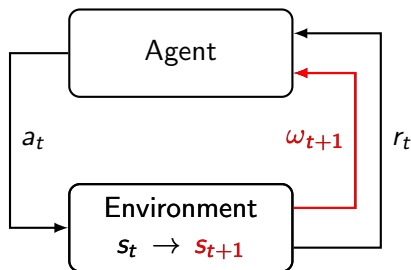
From experience in an environment,
an artificial agent
should be able to **learn** a sequential decision making task
in order **to achieve goals**.



Observations and
actions may be
high dimensional

Objective

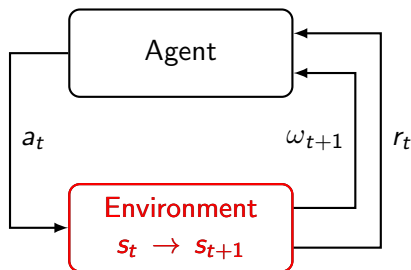
From experience in an environment,
an artificial agent
should be able to **learn** a sequential decision making task
in order **to achieve goals**.



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)

Introduction

- ▶ Experience is gathered in the form of sequences of observations $\omega \in \Omega$, actions $a \in \mathcal{A}$ and rewards $r \in \mathbb{R}$:

$$\omega_0, a_0, r_0, \dots, a_{t-1}, r_{t-1}, \omega_t$$

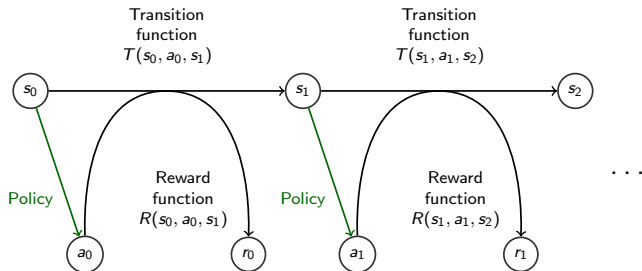
- ▶ In a fully observable environment, the state of the system $s_t \in \mathcal{S}$ is available to the agent.

$$s_t = \omega_t$$

Definition of an MDP

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

- ▶ \mathcal{S} is a finite set of states $\{1, \dots, N_{\mathcal{S}}\}$,
- ▶ \mathcal{A} is a finite set of actions $\{1, \dots, N_{\mathcal{A}}\}$,
- ▶ $T : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ 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 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 expected return, the optimal expected return can be defined as

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

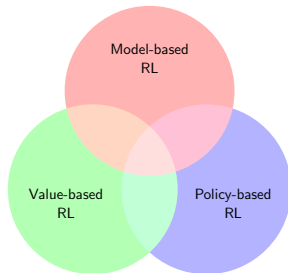
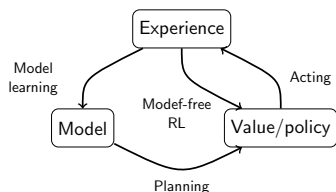
and the optimal policy can be defined as :

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

Overview of deep RL

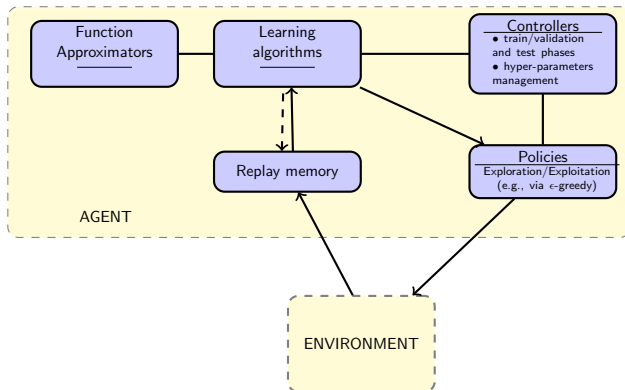
In general, an RL agent may include one or more of the following components :

- ▶ 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
- ▶ a model of the environment in conjunction with a planning algorithm.



Deep learning has brought its generalization capabilities to RL.

Overview



Implementation : <https://github.com/VinF/deer>

Generalisation from limited data in reinforcement learning

Overview

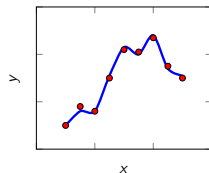
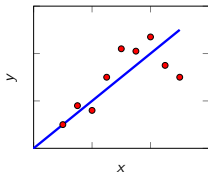
To understand generalization in RL from limited data, we will

- ▶ recall the concept in supervised learning, and
- ▶ introduce the formulation in RL.

We'll then discuss how an agent can have a good generalization in RL.

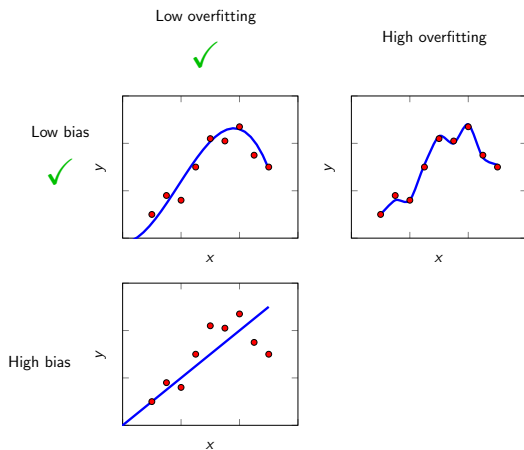
Bias and overfitting in supervised learning

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 $f(x | D_{LS})$.



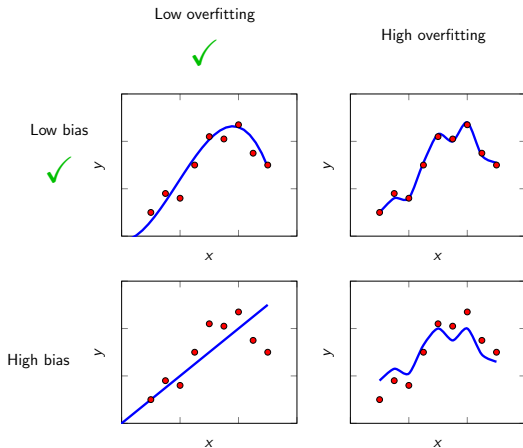
Bias and overfitting in supervised learning

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

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 supervised 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).

Note that there is no such direct bias-variance decomposition for loss functions other than the L_2 loss! It is however always possible to decompose the prediction error with a term related to the lack of expressivity of the model (the bias) and a term due to the limited amount of data (overfitting).

Bias and overfitting in RL

The *off-policy learning algorithm* in RL can be seen as mapping a dataset $D \sim \mathcal{D}$ into a policy π_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 obtain the best policy?



FIGURE – Schematic representation of the bias-overfitting tradeoff.

How to improve generalization ?

How to improve generalization ?

We can optimize the bias-overfitting tradeoff 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)

More details : V François-Lavet, et al. "*An introduction to deep reinforcement learning*". Foundations and Trends in ML.

<https://arxiv.org/abs/1811.12560>

1. Abstract representation

The appropriate level of abstraction plays a key role in the bias-overfitting tradeoff and one of the key advantages of using a small but rich abstract representation is to allow for improved generalization.

- ▶ When considering many features on which to base the policy, an RL algorithm may take into consideration spurious correlations, which leads to overfitting.
- ▶ Removing features that discriminate states with a very different role in the dynamics introduces a bias.

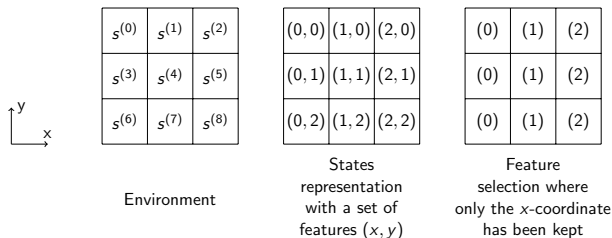


FIGURE – Illustration of the abstract representation.

1. Abstract representation

In POMDPs, policies are based on a state representation built from histories of observations, actions and rewards.

With $\mathcal{H}_t = \Omega \times (\mathcal{A} \times \mathcal{R} \times \Omega)^t$, $\mathcal{H} = \bigcup_{t=0}^{\infty} \mathcal{H}_t$, we consider a mapping $\phi : \mathcal{H} \rightarrow \phi(\mathcal{H})$, where $\phi(\mathcal{H}) = \{\phi(H) | H \in \mathcal{H}\}$.

- ▶ On the one hand, a mapping ϕ with a low cardinality $|\phi(\mathcal{H})|$ reduces the risk of **overfitting** ($|\Pi_{\phi(\mathcal{H})}| \leq |\mathcal{A}|^{|\phi(\mathcal{H})|}$).
- ▶ On the other hand, when ϕ discards information from the history, the state representation $\phi(H)$ might depart from sufficient statistics, which creates an **asymptotic bias**.

2. Modifying the objective function

In order to improve the policy learned by a deep RL algorithm, one can optimize an objective function that diverts from the actual objective. By doing so, a bias is usually introduced but this can in some cases help with generalization. The main approaches to modify the objective function are either

- ▶ to modify the reward of the task to ease learning (reward shaping), or
- ▶ tune the discount factor at training time.

2. Modifying the objective function (online setting)

Motivation from neurosciences :

- ▶ *Empirical studies of cognitive mechanisms in delay of gratification* : The capacity to wait longer for the preferred rewards seems to develop markedly only at about ages 3-4 (“marshmallow experiment”).

In the online setting, the optimal tradeoff evolves as more data is obtained. There is for instance an interest of increasing the discount factor when more data of interest is gathered.

2. Modifying the objective function (online setting)

Increasing the discount factor (e.g., in the DQN algorithm) from 0.95 to 0.99 through learning can significantly improve learning.

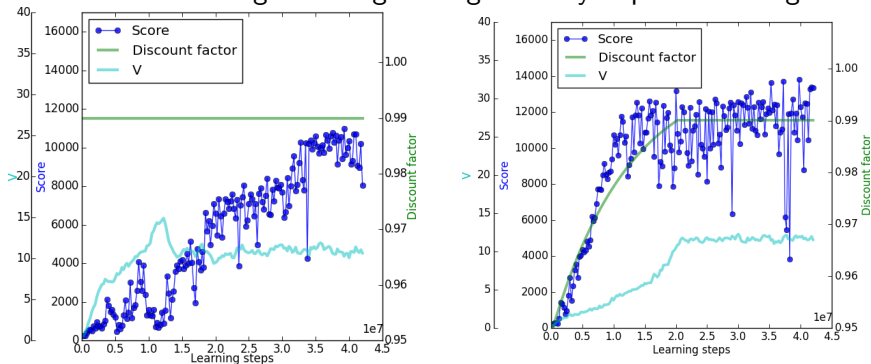
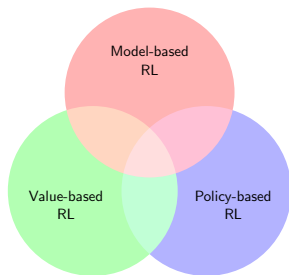
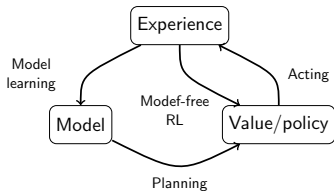


FIGURE – Illustration for the game q-bert of a discount factor γ held fixed on the right and an adaptive discount factor on the right.

3. Choice of the learning algorithm and function approximator selection

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

- ▶ 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
- ▶ a model of the environment in conjunction with a planning algorithm.



Deep learning has brought its generalization capabilities to RL. ▶

3. Choice of the learning algorithm and function approximator selection

- ▶ The function approximator in deep learning characterizes how the features will be treated into higher levels of abstraction. A fortiori, it is related to feature selections (e.g., an attention mechanism), etc.
- ▶ Depending on the task, finding a performant function approximator is easier in either a model-free or a model-based approach. The choice of relying more on one or the other approach is thus also a crucial element to improve generalization.

3. Choice of the learning algorithm

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

- ▶ If the agent does not have access to a generative model of the environment, the learned model will have some inaccuracies.
- ▶ Second, a model-based approach requires working in conjunction with a planning algorithm, which is often computationally demanding.
- ▶ Third, for some tasks, the model of the environment may be learned more efficiently due to the particular structure of the task.

3. Choice of the learning algorithm

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

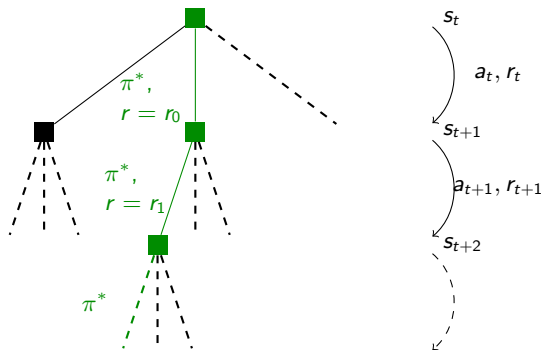


FIGURE – Illustration of model-based.

3. Choice of the learning algorithm : a parallel with neurosciences

In cognitive science, there is a dichotomy between two modes of thoughts (*D. Kahneman. (2011). Thinking, Fast and Slow*) :

- ▶ a "System 1" that is fast and instinctive and
- ▶ a "System 2" that is slower and more logical.



FIGURE – System 1

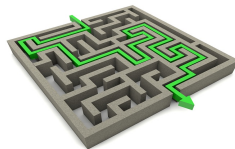


FIGURE – System 2

In deep reinforcement, a similar dichotomy can be observed when we consider the model-free and the model-based approaches.

Best practices in deep RL

How to benchmark deep RL ?

- ▶ Test an algorithm's effectiveness with an average across a few learning trials. If possible study the results with techniques derived from significance testing.

Stochasticity plays a large role in deep RL, both from randomness within initializations of neural networks and stochasticity in environments. Results may vary significantly simply by changing the random seed.

How to benchmark deep RL ?

- ▶ Do not to over-interpret the results.

It is possible that a hypothesis can be shown to hold for one or several given environments and under one or several given set of hyperparameters, but fail in other settings.

How to benchmark deep RL ?

- ▶ Ensure a fair comparison between learning algorithms.

Ensuring that a novel algorithm is indeed performing much better requires proper scientific procedure when choosing such hyperparameters

How to benchmark deep RL ?

- ▶ When choosing metrics to report, it is important to select those that provide a fair comparison.

Using the top-K trials is usually inadequate for fair comparisons.

Combining model-based and model-free via abstract representations

Combining model-based and model-free via abstract representations

Why are we interested in learning everything through one abstract representation ?

- ▶ it can help enforce a good generalization,
- ▶ planning is computationally efficient,
- ▶ it facilitates interpretation of the decisions taken by the agent,
- ▶ it enables strategies for transfer learning, and
- ▶ it can be used to improve exploration.

More details : *Combined Reinforcement Learning via Abstract Representations*, V. Francois-Lavet, Y. Bengio, D. Precup, J. Pineau, 2018 (AAAI).

Combined Reinforcement Learning via Abstract Representations (CRAR)

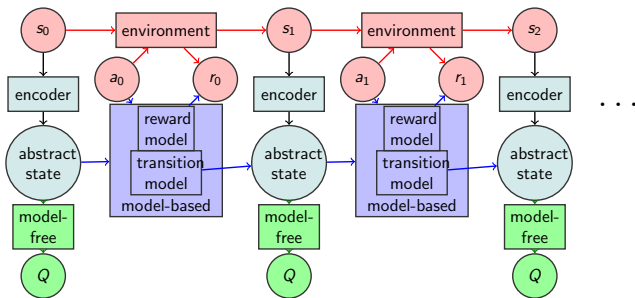


FIGURE – Illustration of the integration of model-based and model-free RL in the CRAR architecture.

The value function and the model are trained using off-policy data, via the abstract representation and without auto-encoder.

Learning the value function

Training of the value function is done with DDQN :

$$Y_k^{DDQN} = r + \gamma Q \left(e(s'; \theta_e^-), \operatorname{argmax}_{a \in \mathcal{A}} Q(e(s'; \theta_e), a; \theta_Q); \theta_Q^- \right),$$

The training is done by minimizing the loss

$$\mathcal{L}_{mf}(\theta_e, \theta_Q) = \left(Q(e(s; \theta_e), a; \theta_Q) - Y_k^{DDQN} \right)^2.$$

This loss trains the weights of both the encoder and the model-free component.

Learning the model

We have one loss for learning the reward, one for the discount factor and one for learning the transition :

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

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

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

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

In practice, there is a pressure to decrease the amount of information being represented.

In our model, we introduce :

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

where s_1 and s_2 are random states stored in the replay memory and C_d is a constant.

The risk of obtaining very large values for the features of the state representation is avoided by the following loss that penalizes abstract states that are out of an L_∞ ball of radius 1 :

$$\mathcal{L}_{d2}(\theta_e) = \max(\|e(s_1; \theta_e)\|_\infty^2 - 1, 0).$$

The loss \mathcal{L}_d , called the representation loss, is a combination of both losses.

Simple labyrinth

Representation of one state for a labyrinth task
(without any reward).

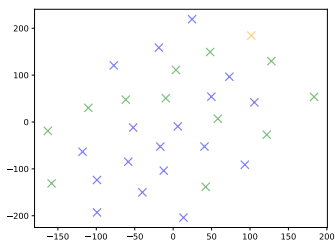
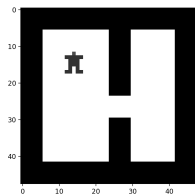


FIGURE – 2D representation using t-SNE (blue represents states where the agent is on the left part, green on the right part and orange in the junction).

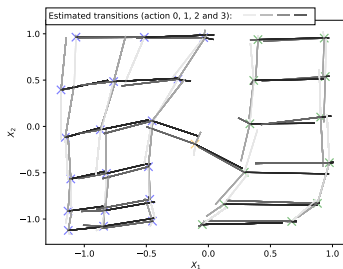


FIGURE – The CRAR agent is able to reconstruct a sensible representation of its environment in 2 dimensions.

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.

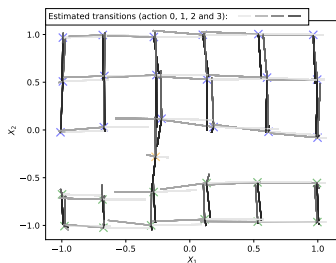


FIGURE – With enforcing $\mathcal{L}_{\text{interpr}}$ and $v(a_0) = [1, 0]$

Catcher

This environment has only a few important features

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

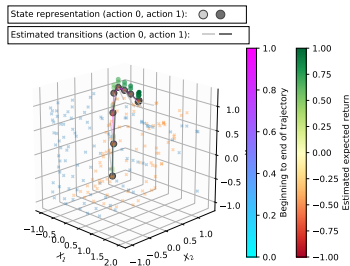
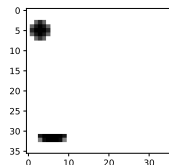


FIGURE –
Without interpretability loss.

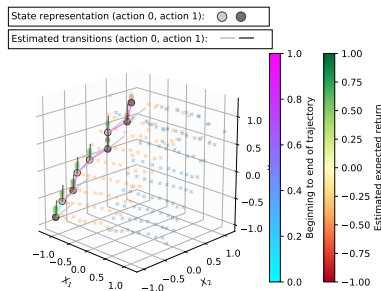


FIGURE – With interpretability loss :
 $v(a^{(1)}) = (1, 1)$ and $v(a^{(2)}) = (-1, 1)$.

Planning

The dynamics for some sequence of actions is 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 $Q(\hat{x}_t, a; \theta_Q)$ ($\mathcal{A}^* \subseteq \mathcal{A}$).

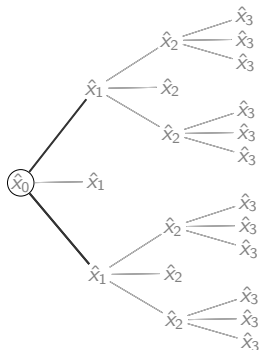


FIGURE – Expansion from current state representation x_0 .

Planning

The dynamics for some sequence of actions is 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 $Q(\hat{x}_t, a; \theta_Q)$ ($\mathcal{A}^* \subseteq \mathcal{A}$).

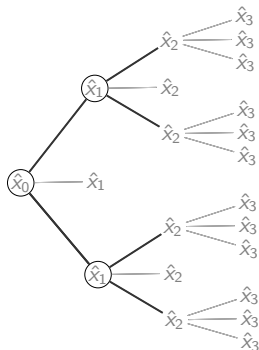


FIGURE – Expansion from current state representation x_0 .

Planning

The dynamics for some sequence of actions is 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 $Q(\hat{x}_t, a; \theta_Q)$ ($\mathcal{A}^* \subseteq \mathcal{A}$).

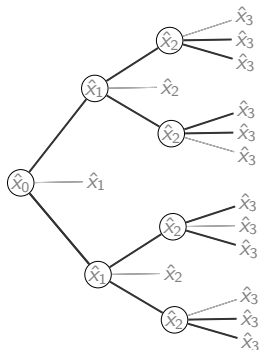


FIGURE – Expansion from current state representation x_0 .

We define recursively the depth- d estimated expected return as

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

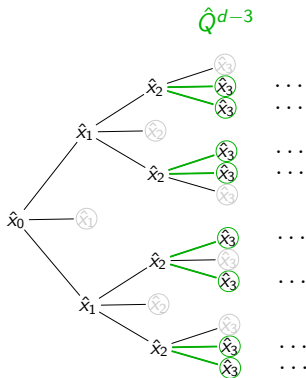


FIGURE – Backup.

We define recursively the depth- d estimated expected return as

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

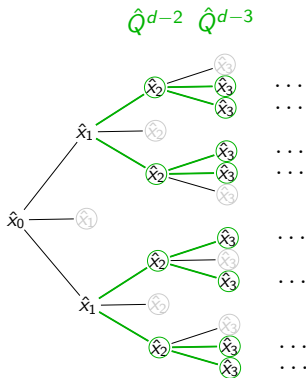


FIGURE – Backup.

We define recursively the depth- d estimated expected return as

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

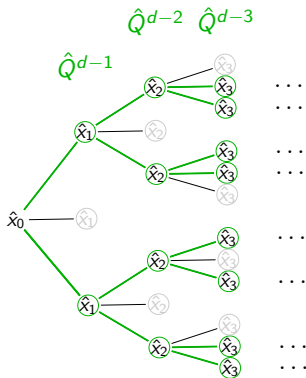


FIGURE – Backup.

We define recursively the depth- d estimated expected return as

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

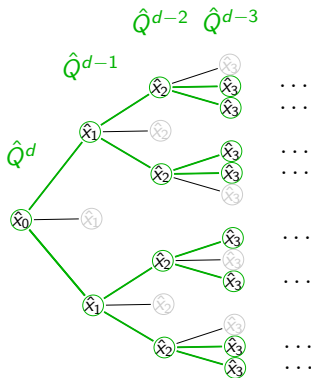


FIGURE – Backup.

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 $\operatorname{argmax}_{a \in \mathcal{A}} Q_{plan}^D(\hat{x}_t, a)$.

Meta-learning with limited off-policy data

The CRAR agent is successfully able to learn from a small set of off-policy data (2×10^5 tuples) in a complex distribution of tasks, while using planning in an abstract state space.

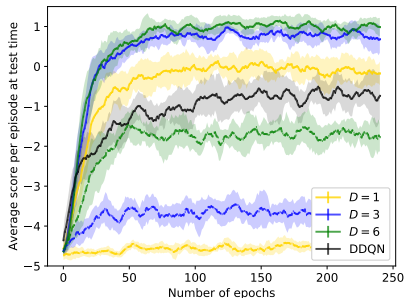
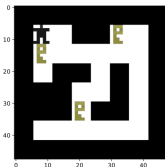


FIGURE – Meta-learning score on a distribution of labyrinths where the training is done with a limited number of transitions obtained off-line by a random policy. 2×10^5 tuples, ~ 500 labyrinths.

Another important challenge : transfer learning

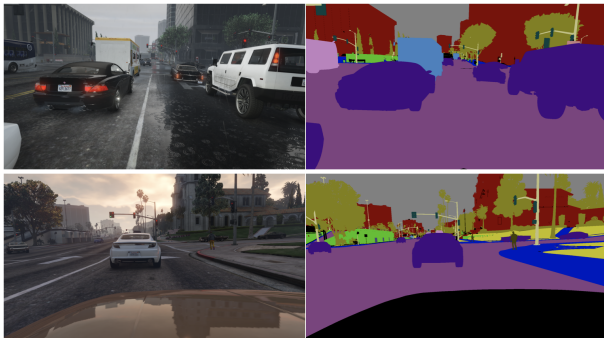


FIGURE – Transfer learning between different renderings. Picture from *"Playing for Data : Ground Truth from Computer Games"*, Richter, S. and Vineet, V., et al

Transfer learning with the CRAR agent

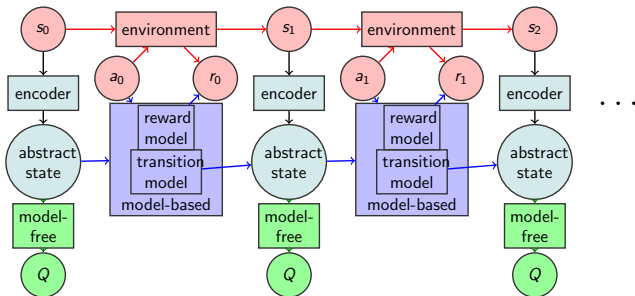


FIGURE – The abstract state can be enforced to be the same for semantically identical observations.

Transfer learning with the CRAR agent



FIGURE – After the first 250 epochs, training and test are done on the same distribution of tasks but the pixels have the opposite values.

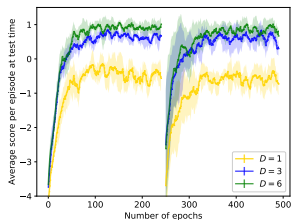


FIGURE – Without transfer

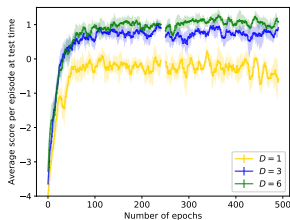


FIGURE – With transfer

Conclusions

Further resources

- ▶ Richard Sutton and Andrew G. Barto. Reinforcement learning : An introduction. Vol. 1. No. 1. Cambridge : MIT press, 1998.
- ▶ RL Course by David Silver on Youtube :
<https://www.youtube.com/watch?v=2pWv7GOvuf0>
- ▶ V François-Lavet, et al. "*An introduction to deep reinforcement learning*". Foundations and Trends in ML.
<https://arxiv.org/abs/1811.12560>

Questions ?