

5. Reinforcement learning

Apprentissage :

- supervisé :

données
labels
 (X_i, Y_i)

(SVM...)

{ on cherche à prédire le label
(la classe...) d'une nouvelle donnée

→ classification
supervisée

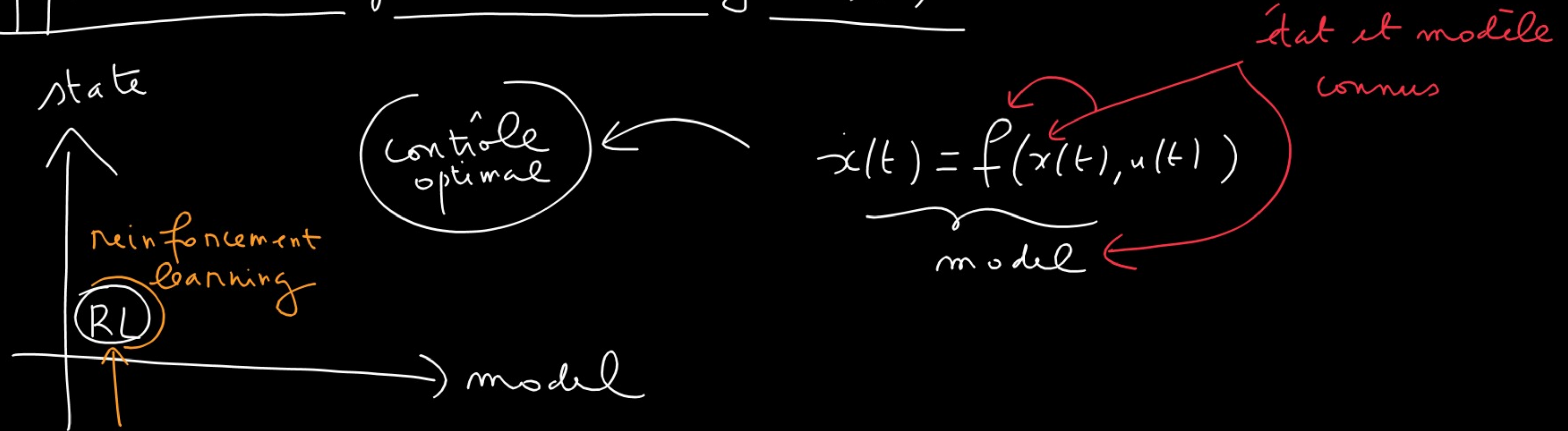
- non-supervisé :



→ clustering : "faire des paquets"
(en postulant un nombre de
classes) (K-means, spectral
clustering...)

- par renforcement

Approximate dynamic learning (ADP):



ex: maps avec

- carte partiellement connue (modèle imparfait)
- position (état) imparfaitement mesuré

→ Contrôle optimal stochastique discret:

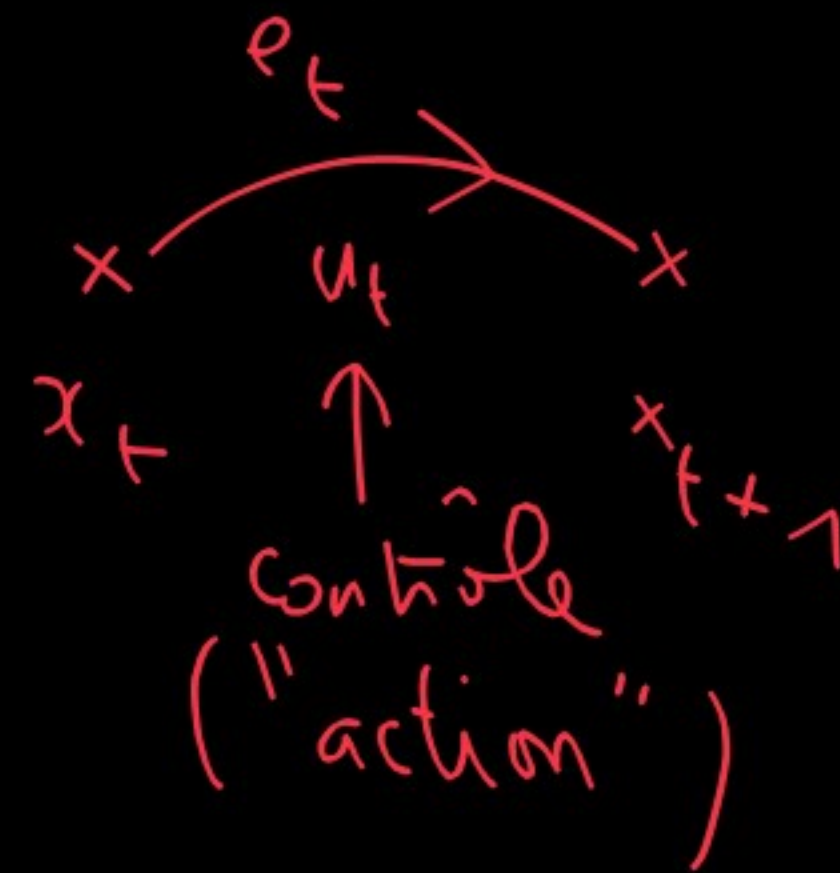
$$E \left(\sum_{t=0}^{N-1} R(x_t, u_t) \right) \rightarrow \max$$

reward

stochastique (e_t : s.a.)

$$x_{t+1} = f(x_t, u_t, e_t), \quad t=0, 1, \dots \quad (\text{temps discret})$$

$$x_0 \text{ connu}, \quad u_t \in U$$



(MDP = Markov Decision Process)

R, f : modèle inconnu

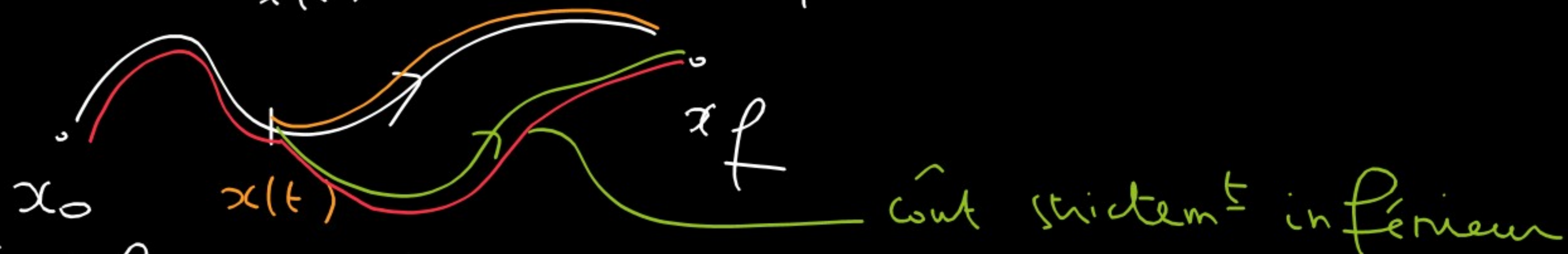
— Sur la base de plusieurs épisodes ("réalisations"),
on fait des moindres carrés pour apprendre R et f :

$$\left. \begin{array}{l} (\tilde{x}_0, \tilde{x}_1, \dots, \tilde{x}_L) \\ (\tilde{u}_0, \dots, \tilde{u}_{L-1}) \\ (\tilde{r}_0, \dots, \tilde{r}_{L-1}) \end{array} \right\} 1 \text{ épisode}$$

— autre approche : Q-learning \Rightarrow on ne cherche pas à apprendre R ou f (= le modèle) mais on utilise le fait que le modèle sous-jacent est du contrôle \hat{u} (optimal stochastique discret)

Th. (Bellman / programmation dynamique):

$x(t)$ = trajectoire optimale



$x: [0, t_f] \rightarrow \mathbb{R}^n$ optimale

$\Rightarrow (\forall t \in [0, t_f]) : x|_{[t, t_f]} : [t, t_f] \rightarrow \mathbb{R}^n$ est encore optimale
pour relier $x(t)$ à x_f .

dém.: f. à-dessous. \square

Déf.: Q-function : on considère le pb précédent mais à horizon infini :

$$\underline{E} \left(\sum_{t=0}^{\infty} \gamma^t R(x_t, u_t) \right) \rightarrow \max$$

$\gamma \in]0, 1[$ (discount factor)

$$x_{t+1} = f(x_t, u_t, e_t)$$

$$x_0 \text{ connu}, u_t \in \bigcup_{n \in \mathbb{N}} \mathbb{R}^m, t = 0, 1, 2, \dots$$

$$Q(x, u) := \max_{\underbrace{u_1, u_2, \dots}_{\text{suite}}} \left\{ E \left(\sum_{t=0}^{\infty} \gamma^t R(x_t, u_t) \right) \mid x_{t+1} = f(x_t, u_t), x_0 = \underbrace{x}, u_0 = \underbrace{u} \right\}$$

Rg: i) $v(x) = \max_{u \in U} Q(x, u)$: fonction valeur

ii) si on connaît Q , on a accès à l'action (la "politique") optimale à jouer dans l'état x :

$$u \in \arg \max_{s \in U} Q(x, s)$$

Bellman

On:

$$Q(x, u) = R(x, u) + E \left(\max_{u_1, u_2, \dots} \left\{ E \left(\sum_{t=1}^{\infty} \gamma^t R(x_t, u_t) \right) \mid \begin{array}{l} x_{t+1} = f(x_t, u_t, e_t), t=1, 2, \dots \\ x_1 = f(x, u, e_0), u_1 = u' \end{array} \right\} \right)$$

Si on apprend la fonction Q directement (sans apprendre le modèle R, f, \dots), on sait quelle action optimale prendre dans l'état x .

$\Rightarrow Q(x, u) = R(x, u) + E \left(\max_{u' \in U} \gamma \cdot Q(f(x, u, \underline{e_0}), u') \right)$: la fonction est solution

de cette équation de pt fixe.

\rightarrow Apprentissage : $\left(\begin{array}{l} \tilde{x}_0, \dots, \tilde{x}_L \\ \tilde{u}_0, \dots, \tilde{u}_{L-1} \\ \tilde{\pi}_0, \dots, \tilde{\pi}_{L-1} \end{array} \right)$ 1 épisode $\rightarrow K$ épisodes

Si $\hat{Q}(x, u, \gamma_k)$ est un estimateur $\textcircled{*}$ (paramétré par γ_k), on l'améliore à chaque nouvel épisode en cherchant γ_{k+1} tq :

$\textcircled{*}$ en pratique $\hat{Q}(x, u, \gamma) := NN_{\gamma}(x, u)$ (γ = poids synaptiques)

$$\hat{Q}(\tilde{x}_t, \tilde{u}_t, y_{k+1}) = (1 - \eta_k) \cdot \left(\hat{\tilde{\pi}}_t + \max_{u' \in U} \gamma \cdot \hat{Q}(\tilde{x}_{t+1}, u', y_{k+1}) \right) + \eta_k \cdot \hat{Q}(\tilde{x}_t, \tilde{u}_t, y_k), \quad t=0, \dots, L-1$$

Si $(y_k)_k \rightarrow \bar{y}$, on a CV vers $\hat{Q}(\cdot, \cdot, \bar{y})$ qui vérifie le pt fixe souhaité sur les réalisations.

$$\rightarrow u \in \arg \max_{v \in U} \hat{Q}(x, v, \bar{y})$$