

# Extensions

Clément Rambour

# Plan

## 1 Introduction

## 2 Policy Gradient

## 3 Introduction aux approches basées modèle

## 4 Bibliographie

# Cours précédents

- Estimation de la valeur par échantillonnage
- Deux grande familles d'approches :
  - MC : approximation de la valeur par moyenne des retour
  - TD : Minimisation de l'erreur de Bellman
- Contrôle : estimation de la fonction de valeur état-action
- Valeurs stockées dans la table  $Q$

# Cours précédents : Approches basées valeurs

- Approximation de la fonction de valeur par une fonction de paramètre  $\theta$  :

$$Q(s, a) \simeq f_{\theta}(s, a)$$

- Optimisation de  $f_{\theta}$  par descente de gradient
- Loss Meas Square Error
  - MC : par rapport au retour
  - TD : par rapport à la TD-target  $\rightarrow$  Q-learning
- Deep Q Networks :  $f_{\theta} \rightarrow$  réseau de neurones

# Cours précédents : Policy gradient

## Motivation :

- Plutôt que de passer par une fonction décrivant le potentiel de chaque état/action de l'agent, pourquoi ne pas apprendre directement la politique ?
- Politique donnée par une distribution paramétrée par  $\theta$  :  $\pi_{\theta}(s, a)$
- Recherche de la politique  $\pi_{\theta}$  paramétrée par  $\theta$ .
- Optimisation d'une fonction objectif :  $J(\theta)$
- **REINFORCE**  $J(\theta) = v_{\pi_{\theta}}(s)$  :

$$\theta = \underset{\theta}{\operatorname{argmax}} v_{\pi_{\theta}}(s)$$

- **Actor Critique** Estimation de la valeur (Critique) et optimisation de la politique (Acteur)

# optimisation de la valeur vs. optimisation de la politique

Avantages de l'optimisation de politique :

- Meilleure convergence
- Fonctionne en grande dimension ou si les actions sont continues
- Peut apprendre des politiques stochastiques

Désavantages :

- Converge vers un minimum local
- Grande variance

# Aujourd'hui

- 1 Récapitulatif sur les approches Actor-Critique
- 2 Introduction aux approches basées modèles

# Plan

- 1 Introduction
- 2 Policy Gradient**
- 3 Introduction aux approches basées modèle
- 4 Bibliographie



# Objectif

- Soit une trajectoire  $\tau$  :

$$\tau = \{s_0, a_0, r_1, \dots, s_{T-1}, a_{T-1}, r_T\}$$

- On note  $R(\tau) = \sum_{t=0}^T r_t$  la somme des récompense.
- $T(\tau) = G_0$  pour  $\gamma = 1$
- La fonction objectif est donc :

$$V(\theta) = \mathbb{E}_{\pi_\theta} \left[ \sum_{t=0}^T r_t \right] = \sum_{\tau} P(\tau; \theta) R(\tau)$$

item  $P(\tau; \theta)$  : probabilité de la trajectoire  $\tau$  en suivant  $\pi_\theta$

- Notre **fonction objectif** est donc :

$$\operatorname{argmax}_{\theta} V(s; \theta) = \operatorname{argmax}_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau)$$

# Gradient de la politique

- Notre **fonction objectif** est donc :

$$\operatorname{argmax}_{\theta} V(s; \theta) = \operatorname{argmax}_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau)$$

- On applique l'astuce vue précédemment pour calculer le gradient :

$$\nabla_{\theta} V(s; \theta) = \sum_{\tau} P(\tau; \theta) \nabla_{\theta} \log P(\tau; \theta) R(\tau)$$

- Espérance approchée par une moyenne sur un batch de  $m$  trajectoires :

$$\nabla_{\theta} V(s; \theta) \simeq \hat{g} = \frac{1}{m} \sum_{i=1}^m R(\tau_i) \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)})$$

- Pas besoin de la dynamique du modèle

# Variance du gradient

- Problème : fonction objectif avec un potentiellement un **très grand nombre de paramètres**
- Conséquence : **grande variance** dans son estimation
- Comment réduire la variance dans l'estimation du gradient  $\nabla_{\theta} V(s; \theta)$  ?

# Variance du gradient

## 1 Structure temporelle → causalité :

$$\nabla_{\theta} V(\theta) = \mathbb{E} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) G_t \right]$$

## 2 Bootstrapping → TD

- Mise à jour de la politique par montée de gradient → Acteur
- Fonction de valeur  $Q$  apprise via TD (par exemple) → Critique

$$\nabla_{\theta} V(\theta) \simeq \mathbb{E} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) Q(s_t, a_t) \right]$$

## 3 Calibration de la somme des récompenses → Utilisation de l'avantage $Q(s, a) - V(s)$ :

$$\nabla_{\theta} V(\theta) = \mathbb{E} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \underbrace{(Q(s_t, a_t) - V(s))}_{A(s_t, a_t)} \right]$$

## 4 L'erreur TD est un bon estimateur de l'avantage :

$$A(s_t, a_t) \simeq \delta_{\pi_{\theta}} = r + \gamma V_{\pi_{\theta}}(s') - V_{\pi_{\theta}}(s)$$

# Off-line Actor Critic : schéma général

## Offline Batch TD Actor-Critic

Jusqu'à convergence, répéter :

- 1 Échantillonner  $\{s_t^{(i)}, a_t^{(i)}\}$  une trajectoire suivant  $\pi_\theta(a|s)$
- 2 Calculer  $y_t^{(i)} = r_{t+1}^{(i)} + \gamma V_{\pi_\theta}(s_{t+1}^{(i)}; w)$
- 3  $w \leftarrow w - \lambda \nabla_w \sum_i \mathcal{L}(V_{\pi_\theta}(s_t^{(i)}; w), y_t^{(i)})$
- 4 Évaluer  $A_{\pi_\theta}(s_t^{(i)}, a_t^{(i)}) = r_{t+1}^{(i)} + \gamma V_{\pi_\theta}(s_{t+1}^{(i)}; w) - V_{\pi_\theta}(s_t^{(i)}; w)$
- 5  $\nabla_\theta J(\theta) \simeq \sum_i \nabla_\theta \log \pi_\theta(a_t^{(i)} | s_t^{(i)}) A_{\pi_\theta}(s_t^{(i)}, a_t^{(i)})$
- 6  $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

# Online Actor Critic : schéma général

**Online TD Actor-Critic** Jusqu'à convergence, répéter :

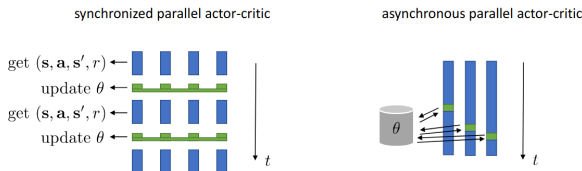
- 1 Prendre l'action  $a_t \sim \pi_\theta(a_t|s_t)$  et observer la transition  $(s_t, a_t, s_{t+1})$
- 2 Calculer  $y_t = r_{t+1} + \gamma V_{\pi_\theta}(s_{t+1}; w)$
- 3  $w \leftarrow w - \lambda \nabla_w \mathcal{L}(V_{\pi_\theta}(s_t; w), y_t)$
- 4 Évaluer  $A_{\pi_\theta}(s_t, a_t) = r_{t+1} + \gamma V_{\pi_\theta}(s_{t+1}; w) - V_{\pi_\theta}(s_t; w)$
- 5  $\nabla_\theta J(\theta) \simeq \nabla_\theta \log \pi_\theta(a_t|s_t) A_{\pi_\theta}(s, a)$
- 6  $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

# Online Actor Critic : schéma général

## Online AC

- 1 Prendre l'action  $a \sim \pi_\theta(a|s)$
- 2 Fit  $V_{\pi_\theta}(s_t^{(i)}; w)$  vers la TD target
- 3 Évaluer  $A_{\pi_\theta}(s_t^{(i)}, a_t^{(i)}) = r_{t+1}^{(i)} + \gamma V_{\pi_\theta}(s_{t+1}^{(i)}; w) - V_{\pi_\theta}(s_t^{(i)}; w)$
- 4  $\nabla_\theta J(\theta) \simeq \sum_i \nabla_\theta \log \pi_\theta(a_t^{(i)}|s_t^{(i)}) A_{\pi_\theta}(s_t^{(i)}, a_t^{(i)})$
- 5  $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

**Problème :** batch de taille 1  $\Rightarrow$  grande corrélation entre les transitions



# Soft Actor Critic

Constat :

- Les algorithmes vus jusqu'à présents convergent souvent vers une politique déterministe
- Particulièrement vrai pour SARSA ou Q-learning
- Également avec policy gradient... donc aussi avec AC
- Lorsqu'un chemin fait monter la somme des récompenses celui-ci est trop vite privilégié
- Problème : politique déterministe pas forcément optimale

**Intuition** : Encourager le modèle à continuer à explorer



# Principe

Comment favoriser l'exploration ?

- Value based algorithms :  $\epsilon - greedy$  → au final converge vers déterministe
- Policy gradient : variance du gradient → non souhaitable

Ajout d'un terme opposée à la confiance du réseau : **l'entropie**

$$\mathcal{H}(\pi(.|s)) = -\mathbb{E}[\log \pi(.|s)]$$

# Soft RL

On cherche la politique qui maximise la somme des récompense et d'entropie maximale :

$$J(\theta) = \mathbb{E}\left[\sum_{t=0}^{T-1} r_t\right] + \lambda \mathcal{H}(\pi(\cdot|s_t))$$

# Soft RL

On cherche la politique qui maximise la somme des récompense et d'entropie maximale :

$$J(\theta) = \mathbb{E}\left[\sum_{t=0}^{T-1} r_t\right] + \lambda \mathcal{H}(\pi(\cdot|s_t))$$

- Modification de la récurrence de Bellman par ajout d'un terme d'entropie  $Q$  :

$$Q(s, a) \leftarrow \mathbb{E}[R_{t+1} + Q(S_{t+1}, A_{t+1}|S_t = s, A_t = a) - \lambda \log \pi(A_{t+1}|S_{t+1})]$$

- Mise à jour de la politique par minimisation de la KL-divergence :

$$\pi \leftarrow \operatorname{argmin} D_{KL}(\pi || \exp Q(s, \cdot)/Z)$$

- cf. TP

# Soft Actor Critic

- De nombreuses versions dans la littératures
- Dernière en date : pas d'utilisation de la fonction de valeur
- Peut être utilisée comme A2C vu précédemment : avec  $Q$  et  $V$ , avec  $A...$

# Plan

- 1 Introduction
- 2 Policy Gradient
- 3 Introduction aux approches basées modèle**
- 4 Bibliographie

# Model based

Et si on **connaît** le modèle *ie.* les transitions entre états ?

- Connu dans de nombreux cas : jeux, systèmes physiques simples, simulations
- Utilisation d'un modèle : plus d'information sur les états suivants
- Possibilité de simuler l'environnement
- Combinaison avec les approches non basées modèles possible

→ Connaître le modèle rend la vie plus simple

# Écriture du problème basé modèle

Notation issues de communautés proches mais différentes : contrôle optimal, contrôle de trajectoire, planification

- $c$  : Coût / rôle inverse de  $r$  la récompense
- $f$  : fonction de transition

## Formulation :

$$a_1, \dots, a_T = \operatorname{argmin}_{a_1, \dots, a_T} \sum_{t=1}^T c(s_t, a_t) \text{ tel que } s_t = f(s_{t-1}, a_{t-1})$$



$$a_1, \dots, a_T = \operatorname{argmax}_{a_1, \dots, a_T} \sum_{t=1}^T r(s_t, a_t) \text{ tel que } s_t = f(s_{t-1}, a_{t-1})$$

# Random shooting methods

Optimisation de la fonction objectif  $J$  par tirage aléatoire

$$(a_1, \dots, a_T) = \mathbf{A} = \operatorname{argmax}_{\mathbf{A}} J(\mathbf{A})$$

Deux étapes :

- 1 Tirer  $\mathbf{A}_1, \dots, \mathbf{A}_N \sim \mathcal{U}$
- 2 Choisir  $\mathbf{A}_i = \operatorname{argmax}_j J(\mathbf{A}_j)$



# Cross-Entropy Methods

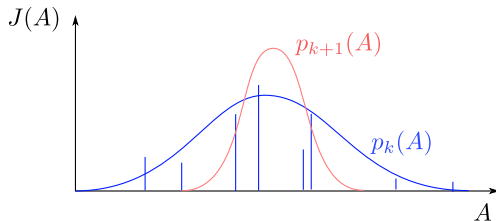
Restriction des échantillons aux actions les plus prometteuses

**Initialisation**  $p(A) = \mathcal{Uou}\mathcal{N}$

**repeat**

- 1 Tirer  $\mathbf{A}_1, \dots, \mathbf{A}_N \sim p(A)$
- 2 Évaluer  $J(\mathbf{A}_1), \dots, J(\mathbf{A}_N)$
- 3 Prendre les  $M$  meilleures trajectoires
- 4 Raffiner  $p(A)$  tel que  $\mathbf{A}_1, \dots, \mathbf{A}_M \sim p(A)$

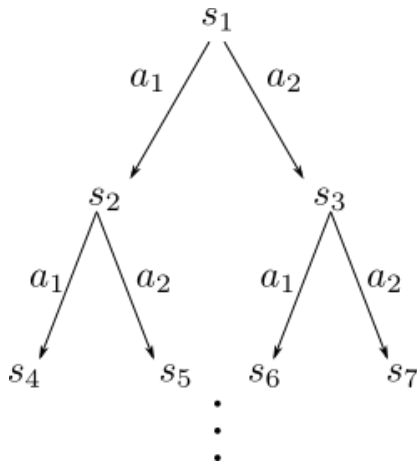
**until** **convergence**;



# Monte Carlo Tree Search

Dans le cas d'un environnement discret, on peut modéliser l'arbre de décision

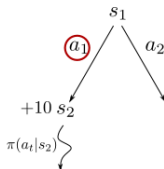
- L'arbre est souvent trop grand pour être parcouru en entier...
- ... en particulier si l'environnement est stochastique
- On peut approcher la valeur des états par échantillonnage



# Monte Carlo Tree Search

Comment parcourir l'arbre ?

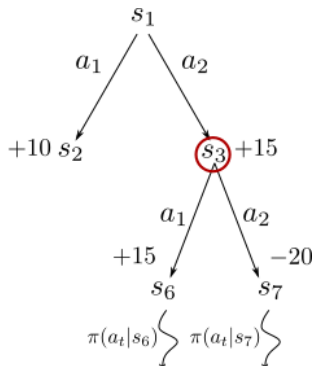
On préfère les noeuds non visités puis ceux de plus grande valeur



# Monte Carlo Tree Search

Comment parcourir l'arbre ?

On préfère les noeuds non visités puis ceux de plus grande valeur



# Monte Carlo Tree Search

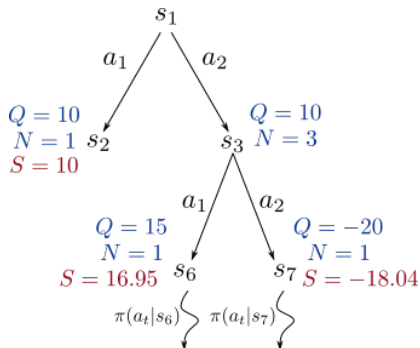
generic MCTS sketch

1. find a leaf  $s_t$  using  $\text{TreePolicy}(s_1)$
  2. evaluate the leaf using  $\text{DefaultPolicy}(s_t)$
  3. update all values in tree between  $s_1$  and  $s_t$
- take best action from  $s_1$

UCT  $\text{TreePolicy}(s_t)$

- if  $s_t$  not fully expanded, choose new  $a_t$   
 else choose child with best  $\text{Score}(s_{t+1})$

$$\text{Score}(s_t) = \frac{Q(s_t)}{N(s_t)} + 2C \sqrt{\frac{2 \ln N(s_{t-1})}{N(s_t)}}$$



# Optimisation de trajectoire

On cherche à résoudre :

$$a_1, \dots, a_T = \underset{a_1, \dots, a_T}{\operatorname{argmin}} \sum_{t=1}^T c(s_t, a_t) \text{ tel que } s_t = f(s_{t-1}, a_{t-1})$$

**Cas linéaire :**

- $f(s_t, a_t) = \mathbf{F}_t \begin{bmatrix} s_t \\ a_t \end{bmatrix} + \mathbf{f}_t$
- $c(s_t, a_t) = \frac{1}{2} \begin{bmatrix} s_t \\ a_t \end{bmatrix}^T \mathbf{C}_t \begin{bmatrix} s_t \\ a_t \end{bmatrix} + \begin{bmatrix} s_t \\ a_t \end{bmatrix}^T \mathbf{c}_t$

# Optimisation de trajectoire

## Linear Quadratic Regulator

### Backward Recursion

$$(V_T, v_T) = (0, 0)$$

for  $t = T$  to 1 do

$$\mathbf{1} \quad (Q_t, q_t) = g(C_t, F_t, V_t, c_t, f_t, v_t)$$

$$\mathbf{2} \quad J(s_t, a_t) = \frac{1}{2} \begin{bmatrix} s_t \\ a_t \end{bmatrix}^T Q_t \begin{bmatrix} s_t \\ a_t \end{bmatrix} + \begin{bmatrix} s_t \\ a_t \end{bmatrix}^T q_t$$

$$\mathbf{3} \quad a_t^* = \operatorname{argmax}_{a_t} J(s_t, a_t) = K_t s_t + k_t$$

$$\mathbf{4} \quad (V_t, v_t) = h(C_t, F_t, V_t, c_t, f_t, v_t)$$

end

Avec  $g$  et  $h$  deux fonctions linéaires

### Forward Recursion

$$(V_T, v_T) = (0, 0)$$

for  $t = 1$  to  $T$  do

$$\mathbf{1} \quad a_t = K_t s_t + k_t$$

$$\mathbf{2} \quad s_{t+1} = f(s_t, a_t)$$

end

# Optimisation de trajectoire

- Dans le cas de transitions stochastiques, DQR est toujours optimal si la distribution est symétrique *eg* gaussienne.
- Cas non linéaire → On se ramène au cas linéaire en partant d'une séquence de référence  $(s'_1, a'_1), \dots (s'_T, a'_T)$  :

$$\mathbf{C}_t = \nabla_{s_t, a_t}^2 c(s'_t, a'_t)$$

$$\mathbf{c}_t = \nabla_{s_t, a_t} c(s'_t, a'_t)$$

$$\mathbf{F}_t = \nabla_{s_t, a_t} f(s'_t, a'_t)$$

- Dans le cas non-linéaire, on fait généralement quelques itérations seulement en suivant la récursion vers l'avant avant de replanifier.



# Plan

- 1 Introduction
- 2 Policy Gradient
- 3 Introduction aux approches basées modèle
- 4 Bibliographie**

# Bibliographie

- Reinforcement Learning : an introduction, second edition, Richard S. Sutton and Adrew G. Barto
- Reinforcement Learning courses, David Silver, DeepMind (<https://www.davidsilver.uk/>)
- A3C (Mnih et al. ICML 2016)
- T. P. Lillicrap et al., Continuous control with deep reinforcement learning
- Browne, Powley, Whitehouse, Lucas, Cowling, Rohlfshagen, Tavener, Perez, Samothrakis, Colton. (2012). A Survey of Monte Carlo Tree Search Methods