

Apprentissage par renforcement appliqué

Méthodes approximatives et DRL [revision 3.4]

Brahim Chaib-draa

Brahim.Chaib-Draa@ift.ulaval.ca



Université Laval

2020-07-23

- 1** Concepts clés en RL (la suite)
- 2** Apprentissage par renforcement profond
- 3** DRL et méthode basée sur les valeurs
- 4** Méthode par recherche de politique
- 5** DRL et méthode *Policy Gradient*
- 6** Pour aller plus loin

Concepts clés en RL (la suite)

1 Concepts clés en RL (la suite)

- Méthodes basées sur les valeurs vs méthodes par recherche de politique
- Méthodes OFF-policy vs ON-Policy
- Méthodes tabulaires vs méthodes approximatives

2 Apprentissage par renforcement profond

3 DRL et méthode basée sur les valeurs

4 Méthode par recherche de politique

5 DRL et méthode *Policy Gradient*

6 Pour aller plus loin

Concepts clés **vus** à la partie précédente

Concepts HAUT niveau :

- Apprentissage **sans modèle** vs apprentissage **basé sur un modèle**
[*Model-free vs model-based*]
- **Apprentissage** vs **Planification**
[*Learning vs Planning*]
- Méthodes **basées sur les valeurs** vs méthodes par **recherche de politique**
[*Value based methods vs Policy search methods*]
- Méthodes **tabulaires** vs méthodes **approximatives**
[*Tabular solution methods vs Approximate solution methods*]

Concept BAS niveau :

- **Prédiction** vs **Contrôle**
[*Prediction vs Control*]
- Apprentissage **EN-ligne** vs **HORS-ligne**
[*ON-line learning vs OFF-line learning*]
- **ON-policy** vs **OFF-Policy**

Concepts clés **vus aujourd'hui**

Concepts HAUT niveau :

- Apprentissage **sans modèle** vs apprentissage **basé sur un modèle**
[*Model-free vs model-based*]
- Apprentissage vs Planification
[*Learning vs Planning*]
- Méthodes **basées sur les valeurs** vs méthodes par **recherche de politique**
[*Value based methods vs Policy search methods*]
- Méthodes **tabulaires** vs méthodes **approximatives**
[*Tabular solution methods vs Approximate solution methods*]

Concept BAS niveau :

- Prédiction vs Contrôle
[*Prediction vs Control*]
- Apprentissage **EN-ligne** vs **HORS-ligne**
[*ON-line learning vs OFF-line learning*]
- **ON-policy** vs **OFF-Policy**

Concepts clés en RL (la suite)

Méthodes basées sur les valeurs vs méthodes par recherche de politique

Méthodes basées sur les valeurs vs méthodes par recherche de politique

Concept HAUT niveau

Méthodes basées sur les valeurs [Value based methods]

L'algorithme **apprend à prédire** V^π et/ou Q^π
pour ensuite trouver π^* pas inférence



Méthodes par recherche de politique [Policy search methods]

L'algorithme **apprend directement** π^*



Concepts clés en RL (la suite)

Méthodes OFF-policy vs ON-Policy

Méthodes *ON-Policy* vs méthodes *OFF-policy*

Concept BAS niveau

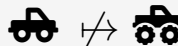
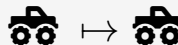
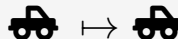
Méthodes *ON-Policy*

La phase d'apprentissage est faite en **suivant l'assomption** que les échantillons proviennent tous de la même distribution.



Si la politique apprise change (même juste un peu), alors on doit utiliser de nouveaux échantillons tirés de cette nouvelle politique.

Analogie On doit réapprendre à conduire de zéro à chaque fois qu'on change de pneu (hivers/été/usure différente/adhérence différente).



Méthodes *ON-Policy* vs méthodes *OFF-policy*

Concept BAS niveau

Méthodes *OFF-policy*

La phase d'apprentissage est faite **sans faire d'assomption** par rapport à la distribution d'où à été tiré chacun des échantillons.



On peut utiliser des échantillons plus anciens, provenant de distribution différente, provenant d'une politique entraînée sur un problème plus général . . .



Analogie : Ma vaste expérience à conduire des autos tamponneuses et des tracteurs à gazon m'a beaucoup appris sur l'art de piloter des *Formules 1*.

« Du point de vue de l'algorithme, les échantillons collectés ne sont pas liés à la politique apprise »

Concepts clés en RL (la suite)

Méthodes tabulaires vs méthodes approximatives

Méthodes tabulaires vs méthodes approximatives

Concept HAUT niveau

Méthodes tabulaires [Tabular solution methods]

L'algorithme **conserve en mémoire** toutes les valeurs prises aux travers **tout l'espace d'état** possible ;

Ex. $Q^\pi(s, a) =$

	a_1	a_2	a_3	a_4
s_1	0.4	0.2	1.0	0.8
s_2	0.6	0.7	0.7	0.6
s_3	0.7	-0.3	0.4	0.2
s_4	0.5	0.1	0.2	0.4



Méthodes approximatives [Approximate solution methods]

L'algorithme **apprend à généraliser** sur la base de son expérience afin de **composer avec l'inconnue** en utilisant approximateur de fonction ;

Ex. $Q^\pi(s, a) \approx \hat{Q}_\theta^\pi(s, a)$



Méthodes approximatives

Cette partie va introduire les algorithmes d'apprentissage par renforcement approximatif **en mettant l'accent** sur les algorithmes utilisant des **réseaux neuronaux profonds** comme approximateur de fonction

... d'où le nom **apprentissage par renforcement profond** [DRL] .

Remarque : méthodes DRL \in méthodes RL approximatives

Les algorithmes d'apprentissage par renforcement profond **ne sont pas les seules méthodes approximatives** en RL. Il est même possible d'avoir de meilleur résultat (dans certaines circonstances) en utilisant d'autres types d'approximateur de fonction :

- ▶ Arbre de décision
- ▶ *K-nearest*
- ▶ Combinaison linéaire de fonction caractéristique (aka. *feature*)
- ▶ ...

Méthodes approximatives

⚠ Problèmes des Méthodes tabulaires :

1. Fonctionnent sous l'assumption qu'il est possible pendant l'apprentissage de **traverser tout l'espace des possibilités**¹
⇒ plus l'espace est large, plus l'apprentissage devient lent ... au point de devenir intractable ;
2. Nécessite de **conserver en mémoire tout l'espace** des possibilités ;

Solution : Apprendre un approximateur de fonction

- ▶ à partir d'un **petit sous-ensemble de l'espace des possibilités**
- ▶ de façon à être capable de **généraliser à des situations jamais rencontré** en entraînement ;

🗨 Con : Les **garanties théoriques** établies pour les méthodes tabulaires **ne se transfèrent pas nécessairement** aux méthodes RL avec approximateurs de fonction.

1. Ici, l'espace des possibilités sous-entend l'espace d'état \mathcal{S} ou d'état-action \mathcal{A}

Méthodes approximatives

⚠ Problèmes des Méthodes tabulaires :

1. Fonctionnent sous l'assumption qu'il est possible pendant l'apprentissage de **traverser tout l'espace des possibilités**¹
⇒ plus l'espace est large, plus l'apprentissage devient lent ... au point de devenir intractable ;
2. Nécessite de **conserver en mémoire tout l'espace** des possibilités ;

Solution : Apprendre un approximateur de fonction

- ▶ à partir d'un **petit sous-ensemble de l'espace des possibilités**
- ▶ de façon à être capable de **généraliser à des situations jamais rencontré** en entraînement ;

🗨 **Con :** Les **garanties théoriques** établies pour les méthodes tabulaires **ne se transfèrent pas nécessairement** aux méthodes RL avec approximateurs de fonction.

1. Ici, l'espace des possibilités sous-entend l'espace d'état \mathcal{S} ou d'état-action \mathcal{A}

Méthodes approximatives

⚠ Problèmes des Méthodes tabulaires :

1. Fonctionnent sous l'assumption qu'il est possible pendant l'apprentissage de **traverser tout l'espace des possibilités**¹
⇒ plus l'espace est large, plus l'apprentissage devient lent ... au point de devenir intractable ;
2. Nécessite de **conserver en mémoire tout l'espace** des possibilités ;

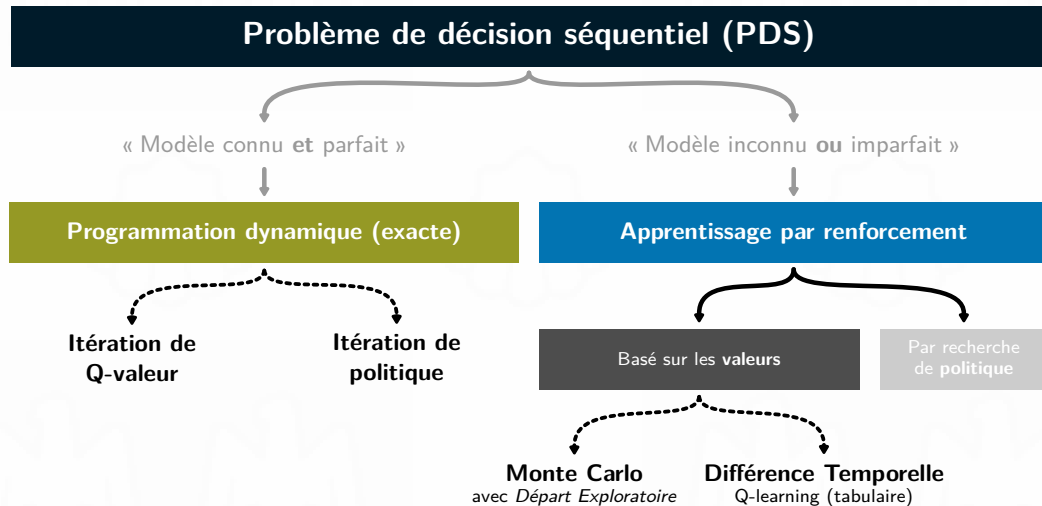
Solution : Apprendre un approximateur de fonction

- ▶ à partir d'un **petit sous-ensemble de l'espace des possibilités**
- ▶ de façon à être capable de **généraliser à des situations jamais rencontré** en entraînement ;

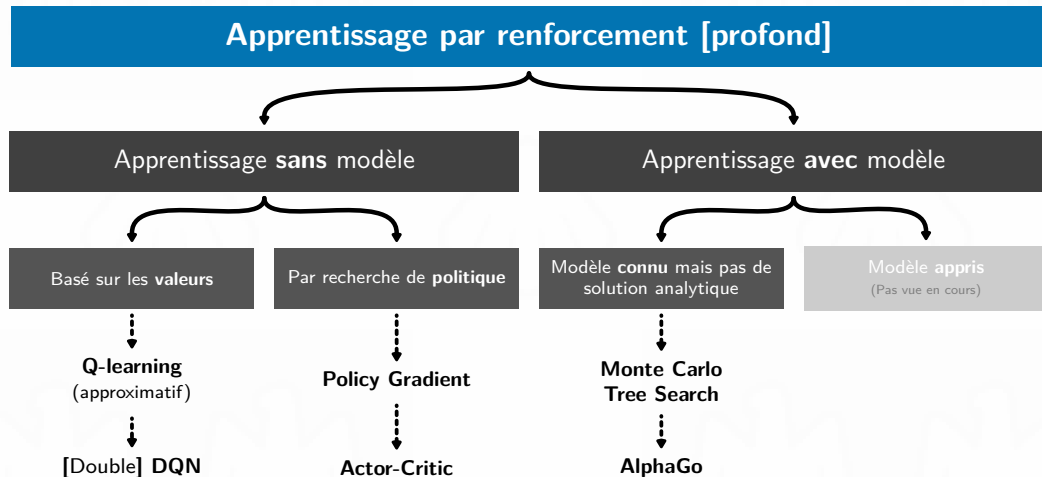
🗨 **Con :** Les **garanties théoriques** établies pour les méthodes tabulaires **ne se transfèrent pas nécessairement** aux méthodes RL avec approximateurs de fonction.

1. Ici, l'espace des possibilités sous-entend l'espace d'état \mathcal{S} ou d'état-action \mathcal{A}

Algorithme vu à la partie précédente



Algorithmes vus dans cette partie



Apprentissage par renforcement profond

1 Concepts clés en RL (la suite)

2 Apprentissage par renforcement profond

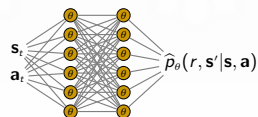
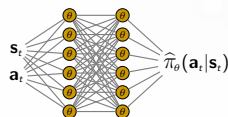
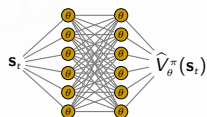
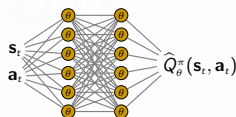
3 DRL et méthode basée sur les valeurs

4 Méthode par recherche de politique

5 DRL et méthode *Policy Gradient*

6 Pour aller plus loin

Apprentissage par renforcement profond [DRL]



Caractéristiques : Apprend à **généraliser** son expérience à **des situations jamais rencontrés** ;

- Idées clés :**
- Utiliser des réseaux neuronaux comme approximateur de fonction pour **représenter** (un ou plusieurs) :
 - fonction de valeur V^π et/ou Q^π
 - politique π
 - modèle $p(r, s' | s, a)$
 - Optimiser l'objectif par **descente de gradient stochastique** ;

Apprentissage par renforcement profond [DRL]

- 👍 **Pro :**
 - Capable de fonctionner sur des espaces de possibilité massifs ;
 - Capable de gérer les situations inconnues (en principe) ;
 - Pas de fonction caractéristique à développer (end-to-end) ;
- 👎 **Con :**
 - Perte de la plupart des garanties théoriques des méthodes tabulaires ;

Enjeux spécifiques au mariage RL/DNN

Processus de décision séquentielle \implies données séquentielles \implies

- les données sont **corrélées** ;
- les données ne sont **pas indépendantes et identiquement distribuées** ;

⚠ **Problème :** les réseaux neuronaux profonds **fonctionnent sous l'assumption** que les données sont iid.

Comment faire pour réconcilier les deux ?

DRL et méthode basée sur les valeurs

1 Concepts clés en RL (la suite)

2 Apprentissage par renforcement profond

3 DRL et méthode basée sur les valeurs

- Q-learning approximatif
- Deep Q-Network
- Double Deep Q-Network

4 Méthode par recherche de politique

5 DRL et méthode *Policy Gradient*

6 Pour aller plus loin

DRL et méthode basée sur les valeurs

Q-learning approximatif

Rappel

★ | « The max trick »

Sauter l'étape d'avoir à représenter explicitement la politique π

$$\max_{\mathbf{a} \in \mathcal{A}(\mathbf{s})} Q^\pi(\mathbf{s}, \mathbf{a}) \iff Q^\pi(\mathbf{s}, \pi_{\text{greedy}}(\mathbf{s})) \quad \text{avec} \quad \pi_{\text{greedy}}(\mathbf{s}) := \arg \max_{\mathbf{a} \in \mathcal{A}(\mathbf{s})} Q(\mathbf{s}, \mathbf{a})$$

↺ Algorithmme | Q-Learning (tabulaire)

Soit $\mu(\cdot|\mathbf{s}_t)$, une politique exploratoire de votre choix, (ex. : aléatoirement, « greedy », « ϵ -greedy », ...)

la politique cible $\pi(\mathbf{s}) \doteq \pi_{\text{greedy}}$, le « learning rate » $\alpha \in]0, 1]$ et le « global step » $k \leftarrow 0$

Répéter pour chaque trajectoire τ jusqu'à convergence vers π^* :

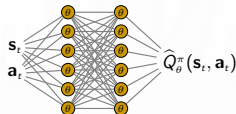
Répéter pour chaque « timestep » $t = 0, 1, 2, \dots, t, t+1, \dots, T$ de la trajectoire τ :

- 1 $r_{t+1}^{(k)}, \mathbf{s}_{t+1}^{(k)} \leftarrow$ exécuter $\mathbf{a}_t \sim \mu(\cdot|\mathbf{s}_t)$ dans l'environnement et observer la réaction
- 2 $Q_{k+1}^\pi(\mathbf{s}_t, \mathbf{a}_t) \leftarrow Q_k^\pi(\mathbf{s}_t, \mathbf{a}_t) + \alpha \left[r_{t+1}^{(k)} + \gamma \max_{\mathbf{a}' \in \mathcal{A}(\mathbf{s}_{t+1})} Q_k^\pi(\mathbf{s}_{t+1}^{(k)}, \mathbf{a}') - Q_k^\pi(\mathbf{s}_t, \mathbf{a}_t) \right]$
- 3 $\mathbf{s}_t \leftarrow \mathbf{s}_{t+1}^{(k)}$ et $k \leftarrow k + 1$

Q-learning approximatif [Approximate Q-learning]

$$\hat{Q}_{\theta}^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = \theta_1 f_1(\mathbf{s}_t, \mathbf{a}_t) + \theta_2 f_2(\mathbf{s}_t, \mathbf{a}_t) + \dots + \theta_n f_n(\mathbf{s}_t, \mathbf{a}_t)$$

ou



Caractéristique : Apprentissage **EN-ligne** et **OFF-policy** (comme Q-learning tabulaire)

Idées clés :

- 1 **Représenter** la fonction de valeur par un approximateur de fonction :

$$Q^{\pi}(\mathbf{s}, \mathbf{a}) \approx \hat{Q}_{\theta}^{\pi}(\mathbf{s}, \mathbf{a}) \quad \text{avec } \theta \text{ un vecteur de paramètres}$$

- 2 Définir une **fonction objectif** par *erreur quadratique moyenne* [MSE] et son **gradient** :

$$L_k(\theta_k) = \mathbb{E}_{\mathbf{s}' \sim p(\cdot | \mathbf{s}, \mathbf{a})} \left[\left(\mathbf{y} - \hat{Q}_{\theta}^{\pi}(\mathbf{s}, \mathbf{a}) \right)^2 \middle| \theta = \theta_k \right] \quad \text{avec } \mathbf{y} = r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}' \in \mathcal{A}(\mathbf{s}')} \hat{Q}_{\theta}^{\pi}(\mathbf{s}', \mathbf{a}')$$

$$\frac{1}{2} \nabla_{\theta_k} L(\theta_k) = \mathbb{E}_{\mathbf{s}' \sim p(\cdot | \mathbf{s}, \mathbf{a})} \left[\left(r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}' \in \mathcal{A}(\mathbf{s}')} \hat{Q}_{\theta}^{\pi}(\mathbf{s}', \mathbf{a}') - \hat{Q}_{\theta}^{\pi}(\mathbf{s}, \mathbf{a}) \right) \nabla_{\theta_k} \hat{Q}_{\theta}^{\pi}(\mathbf{s}, \mathbf{a}) \middle| \theta = \theta_k \right]$$

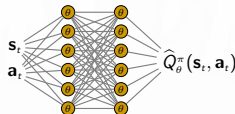
- 3 **Optimiser** l'objectif par *descente de gradient stochastique* [SGD] :

$$\theta_{k+1} \longleftarrow \theta_k - \alpha \nabla_{\theta_k} L(\theta_k)$$

Q-learning approximatif [Approximate Q-learning]

$$\widehat{Q}_{\theta}^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = \theta_1 f_1(\mathbf{s}_t, \mathbf{a}_t) + \theta_2 f_2(\mathbf{s}_t, \mathbf{a}_t) + \dots + \theta_n f_n(\mathbf{s}_t, \mathbf{a}_t)$$

ou



Caractéristique : Apprentissage **EN-ligne** et **OFF-policy** (comme Q-learning tabulaire)

Idées clés :

- 1 **Représenter** la fonction de valeur par un approximateur de fonction :

$$Q^{\pi}(\mathbf{s}, \mathbf{a}) \approx \widehat{Q}_{\theta}^{\pi}(\mathbf{s}, \mathbf{a}) \quad \text{avec } \theta \text{ un vecteur de paramètres}$$

- 2 Définir une **fonction objectif** par *erreur quadratique moyenne [MSE]* et son **gradient** :

$$L_k(\theta_k) = \mathbb{E}_{\mathbf{s}' \sim p(\cdot | \mathbf{s}, \mathbf{a})} \left[\left(\mathbf{y} - \widehat{Q}_{\theta}^{\pi}(\mathbf{s}, \mathbf{a}) \right)^2 \middle| \theta = \theta_k \right] \quad \text{avec } \mathbf{y} = r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}' \in \mathcal{A}(\mathbf{s}')} \widehat{Q}_{\theta}^{\pi}(\mathbf{s}', \mathbf{a}')$$

$$\frac{1}{2} \nabla_{\theta_k} L(\theta_k) = \mathbb{E}_{\mathbf{s}' \sim p(\cdot | \mathbf{s}, \mathbf{a})} \left[\left(r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}' \in \mathcal{A}(\mathbf{s}')} \widehat{Q}_{\theta}^{\pi}(\mathbf{s}', \mathbf{a}') - \widehat{Q}_{\theta}^{\pi}(\mathbf{s}, \mathbf{a}) \right) \nabla_{\theta_k} \widehat{Q}_{\theta}^{\pi}(\mathbf{s}, \mathbf{a}) \middle| \theta = \theta_k \right]$$

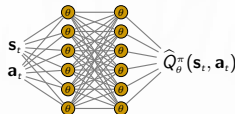
- 3 **Optimiser** l'objectif par *descente de gradient stochastique [SGD]* :

$$\theta_{k+1} \longleftarrow \theta_k - \alpha \nabla_{\theta_k} L(\theta_k)$$

Q-learning approximatif [Approximate Q-learning]

$$\widehat{Q}_{\theta}^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = \theta_1 f_1(\mathbf{s}_t, \mathbf{a}_t) + \theta_2 f_2(\mathbf{s}_t, \mathbf{a}_t) + \dots + \theta_n f_n(\mathbf{s}_t, \mathbf{a}_t)$$

ou



Caractéristique : Apprentissage **EN-ligne** et **OFF-policy** (comme Q-learning tabulaire)

Idées clés :

- 1 **Représenter** la fonction de valeur par un approximateur de fonction :

$$Q^{\pi}(\mathbf{s}, \mathbf{a}) \approx \widehat{Q}_{\theta}^{\pi}(\mathbf{s}, \mathbf{a}) \quad \text{avec } \theta \text{ un vecteur de paramètres}$$

- 2 Définir une **fonction objectif** par *erreur quadratique moyenne* [MSE] et son **gradient** :

$$L_k(\theta_k) = \mathbb{E}_{\mathbf{s}' \sim p(\cdot | \mathbf{s}, \mathbf{a})} \left[\left(\mathbf{y} - \widehat{Q}_{\theta}^{\pi}(\mathbf{s}, \mathbf{a}) \right)^2 \middle| \theta = \theta_k \right] \quad \text{avec } \mathbf{y} = r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}' \in \mathcal{A}(\mathbf{s}')} \widehat{Q}_{\theta}^{\pi}(\mathbf{s}', \mathbf{a}')$$

$$\frac{1}{2} \nabla_{\theta_k} L(\theta_k) = \mathbb{E}_{\mathbf{s}' \sim p(\cdot | \mathbf{s}, \mathbf{a})} \left[\left(r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}' \in \mathcal{A}(\mathbf{s}')} \widehat{Q}_{\theta}^{\pi}(\mathbf{s}', \mathbf{a}') - \widehat{Q}_{\theta}^{\pi}(\mathbf{s}, \mathbf{a}) \right) \nabla_{\theta_k} \widehat{Q}_{\theta}^{\pi}(\mathbf{s}, \mathbf{a}) \middle| \theta = \theta_k \right]$$

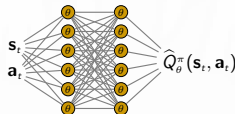
- 3 **Optimiser** l'objectif par *descente de gradient stochastique* [SGD] :

$$\theta_{k+1} \longleftarrow \theta_k - \alpha \nabla_{\theta_k} L(\theta_k)$$

Q-learning approximatif [Approximate Q-learning]

$$\widehat{Q}_{\theta}^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = \theta_1 f_1(\mathbf{s}_t, \mathbf{a}_t) + \theta_2 f_2(\mathbf{s}_t, \mathbf{a}_t) + \dots + \theta_n f_n(\mathbf{s}_t, \mathbf{a}_t)$$

ou



Caractéristique : Apprentissage **EN-ligne** et **OFF-policy** (comme Q-learning tabulaire)

Idées clés :

- 1 **Représenter** la fonction de valeur par un approximateur de fonction :

$$Q^{\pi}(\mathbf{s}, \mathbf{a}) \approx \widehat{Q}_{\theta}^{\pi}(\mathbf{s}, \mathbf{a}) \quad \text{avec } \theta \text{ un vecteur de paramètres}$$

- 2 Définir une **fonction objectif** par *erreur quadratique moyenne* [MSE] et son **gradient** :

$$L_k(\theta_k) = \mathbb{E}_{\mathbf{s}' \sim p(\cdot | \mathbf{s}, \mathbf{a})} \left[\left(\mathbf{y} - \widehat{Q}_{\theta}^{\pi}(\mathbf{s}, \mathbf{a}) \right)^2 \middle| \theta = \theta_k \right] \quad \text{avec } \mathbf{y} = r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}' \in \mathcal{A}(\mathbf{s}')} \widehat{Q}_{\theta}^{\pi}(\mathbf{s}', \mathbf{a}')$$

$$\frac{1}{2} \nabla_{\theta_k} L(\theta_k) = \mathbb{E}_{\mathbf{s}' \sim p(\cdot | \mathbf{s}, \mathbf{a})} \left[\left(r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}' \in \mathcal{A}(\mathbf{s}')} \widehat{Q}_{\theta}^{\pi}(\mathbf{s}', \mathbf{a}') - \widehat{Q}_{\theta}^{\pi}(\mathbf{s}, \mathbf{a}) \right) \nabla_{\theta_k} \widehat{Q}_{\theta}^{\pi}(\mathbf{s}, \mathbf{a}) \middle| \theta = \theta_k \right]$$

- 3 **Optimiser** l'objectif par *descente de gradient stochastique* [SGD] :

$$\theta_{k+1} \longleftarrow \theta_k - \alpha \nabla_{\theta_k} L(\theta_k)$$

↺ Algorithmes | Q-Learning (approximatif)

Soit la politique exploratoire $\mu(\cdot|s)$, la politique cible $\pi(s) \doteq \pi_{greedy}$, le « learning rate » $\alpha \in]0, 1]$,
le « global step » $k \leftarrow 0$, et l'approximateur de fonction \widehat{Q}_θ^π avec le vecteur de paramètre θ

Répéter pour chaque trajectoire τ jusqu'à convergence vers π^* :

Répéter pour chaque « timestep » $t = 0, 1, 2, \dots, t, t+1, \dots, T$ de la trajectoire τ :

```

1       $r_{t+1}^{(k)}, s_{t+1}^{(k)} \leftarrow$  exécuter  $a_t \sim \mu(\cdot|s_t)$  dans l'environnement et observer la réaction
2
3       $y = r_{t+1}^{(k)} + \gamma \max_{a' \in \mathcal{A}(s_t)} \widehat{Q}_\theta^\pi(s_{t+1}^{(k)}, a')$ 
4
5       $\theta_{k+1} \leftarrow \theta_k - \alpha \left( y - \widehat{Q}_\theta^\pi(s_t, a_t) \right) \frac{\partial}{\partial \theta_{s,a}} \widehat{Q}_\theta^\pi(s_t, a_t)$ 
6
7       $s_t \leftarrow s_{t+1}^{(k)}$  et  $k \leftarrow k + 1$ 

```

👍 **Pro :** Fonctionne bien avec un approximateurs linéaire

👎 **Con :** Oscille ou diverge avec un réseau neuronaux ⚠️

↺ Algorithmes | Q-Learning (approximatif)

Soit la politique exploratoire $\mu(\cdot|s)$, la politique cible $\pi(s) \doteq \pi_{greedy}$, le « learning rate » $\alpha \in]0, 1]$, le « global step » $k \leftarrow 0$, et l'approximateur de fonction \widehat{Q}_θ^π avec le vecteur de paramètre θ

Répéter pour chaque trajectoire τ jusqu'à convergence vers π^* :

Répéter pour chaque « timestep » $t = 0, 1, 2, \dots, t, t+1, \dots, T$ de la trajectoire τ :

```

1       $r_{t+1}^{(k)}, s_{t+1}^{(k)} \leftarrow$  exécuter  $a_t \sim \mu(\cdot|s_t)$  dans l'environnement et observer la réaction
2
3       $y = r_{t+1}^{(k)} + \gamma \max_{a' \in \mathcal{A}(s_t)} \widehat{Q}_\theta^\pi(s_{t+1}^{(k)}, a')$ 
4
5       $\theta_{k+1} \leftarrow \theta_k - \alpha \left( y - \widehat{Q}_\theta^\pi(s_t, a_t) \right) \frac{\partial}{\partial \theta_{s,a}} \widehat{Q}_\theta^\pi(s_t, a_t)$ 
6
7       $s_t \leftarrow s_{t+1}^{(k)}$  et  $k \leftarrow k + 1$ 

```

👍 **Pro :** Fonctionne bien avec un approximateurs linéaire

👎 **Con :** **Oscille ou diverge** avec un réseau neuronaux ⚠️

DRL et méthode basée sur les valeurs

Deep Q-Network

⚠ Q-learning avec réseau neuronaux est instable

Causes 1 : Les données sont séquentielles \implies les échantillons sont corrélés et non-iid

$$\tau = \mathbf{s}_0, \mathbf{a}_0, r_1, \mathbf{s}_1, \mathbf{a}_1, r_2, \mathbf{s}_2, \dots, \mathbf{s}_t, \mathbf{a}_t, r_{t+1}, \mathbf{s}_{t+1}, \dots, \mathbf{s}_{T-1}, \mathbf{a}_{T-1}, r_T, \mathbf{s}_T$$

Conséquence : l'algorithme « *overfit* » une multitude de sous-séquences de l'espace

Solution : Emmagaziner les transitions expérimentées dans un « replay buffer »

$$\left\{ \left(\mathbf{s}_0^{(1)}, \mathbf{a}_0^{(1)}, r_1^{(1)}, \mathbf{s}_1^{(1)} \right), \left(\mathbf{s}_1^{(2)}, \mathbf{a}_1^{(2)}, r_2^{(2)}, \mathbf{s}_2^{(2)} \right), \dots, \left(\mathbf{s}_t^{(n)}, \mathbf{a}_t^{(n)}, r_{t+1}^{(n)}, \mathbf{s}_{t+1}^{(n)} \right) \right\}$$

et apprendre à partir de transitions sélectionnées au hasard (dans le « replay buffer »)¹.

L'apprentissage se fait donc sur des données qui ne sont plus séquentielles puisque l'ordre d'apparition des transitions est mélangé² ce qui brise la corrélation des échantillons et nous ramène dans un état iid.

1. Cette technique s'appelle « Experience replay ». Voir [Deep Q-Network | Ressource théorique](#) en fin de section

2. On peut le faire parce que l'algorithme Q-learning est « OFF-policy »

⚠ Q-learning avec réseau neuronaux est instable

Causes 1 : Les données sont séquentielles \implies les échantillons sont corrélés et non-iid

$$\tau = \mathbf{s}_0, \mathbf{a}_0, r_1, \mathbf{s}_1, \mathbf{a}_1, r_2, \mathbf{s}_2, \dots, \mathbf{s}_t, \mathbf{a}_t, r_{t+1}, \mathbf{s}_{t+1}, \dots, \mathbf{s}_{T-1}, \mathbf{a}_{T-1}, r_T, \mathbf{s}_T$$

Conséquence : l'algorithme « *overfit* » une multitude de sous-séquences de l'espace

Solution : Emmagaziner les transitions expérimentées dans un « replay buffer »

$$\left\{ \left(\mathbf{s}_0^{(1)}, \mathbf{a}_0^{(1)}, r_1^{(1)}, \mathbf{s}_1^{(1)} \right), \left(\mathbf{s}_1^{(2)}, \mathbf{a}_1^{(2)}, r_2^{(2)}, \mathbf{s}_2^{(2)} \right), \dots, \left(\mathbf{s}_t^{(n)}, \mathbf{a}_t^{(n)}, r_{t+1}^{(n)}, \mathbf{s}_{t+1}^{(n)} \right) \right\}$$

et apprendre à partir de transitions sélectionnées au hasard (dans le « replay buffer »)¹.

L'apprentissage se fait donc sur des données qui ne sont plus séquentielles puisque l'ordre d'apparition des transitions est mélangé² ce qui brise la corrélation des échantillons et nous ramène dans un état iid.

1. Cette technique s'appelle « Experience replay ». Voir [Deep Q-Network | Ressource théorique](#) en fin de section

2. On peu le faire parce que l'algorithme Q-learning est « OFF-policy »

⚠ Q-learning avec réseau neuronaux est instable

Causes 1 : Les données sont séquentielles \implies les échantillons sont corrélés et non-iid

$$\tau = \mathbf{s}_0, \mathbf{a}_0, r_1, \mathbf{s}_1, \mathbf{a}_1, r_2, \mathbf{s}_2, \dots, \mathbf{s}_t, \mathbf{a}_t, r_{t+1}, \mathbf{s}_{t+1}, \dots, \mathbf{s}_{T-1}, \mathbf{a}_{T-1}, r_T, \mathbf{s}_T$$

Conséquence : l'algorithme « *overfit* » une multitude de sous-séquences de l'espace

Solution : **Emmagasiner les transitions** expérimentées dans un « **replay buffer** »

$$\left\{ \left(\mathbf{s}_0^{(1)}, \mathbf{a}_0^{(1)}, r_1^{(1)}, \mathbf{s}_1^{(1)} \right), \left(\mathbf{s}_1^{(2)}, \mathbf{a}_1^{(2)}, r_2^{(2)}, \mathbf{s}_2^{(2)} \right), \dots, \left(\mathbf{s}_t^{(n)}, \mathbf{a}_t^{(n)}, r_{t+1}^{(n)}, \mathbf{s}_{t+1}^{(n)} \right) \right\}$$

et **apprendre à partir de transitions sélectionnées au hasard** (dans le « replay buffer »)¹.

L'apprentissage se fait donc sur des données qui ne sont plus séquentielles puisque l'ordre d'apparition des transitions est mélangé² ce qui brise la corrélation des échantillons et nous ramène dans un état iid.

1. Cette technique s'appelle « Experience replay ». Voir [Deep Q-Network | Ressource théorique](#) en fin de section

2. On peut le faire parce que l'algorithme Q-learning est « OFF-policy »

⚠ Q-learning avec réseau neuronaux est instable

Causes 2 : Le « *TD-target* » :

- ▷ est non-stationnaire contrairement au contexte de l'apprentissage supervisé ;
- ▷ est corrélé avec l'approximateur \widehat{Q}_θ^π qu'on cherche à apprendre ;

$$L_k(\theta_k) = \mathbb{E}_{\mathbf{s}' \sim p(\cdot | \mathbf{s}, \mathbf{a})} \left[\underbrace{\left(r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}' \in \mathcal{A}(\mathbf{s}')} \widehat{Q}_\theta^\pi(\mathbf{s}', \mathbf{a}') \right)}_{\text{TD-target}} - \widehat{Q}_\theta^\pi(\mathbf{s}, \mathbf{a}) \right]^2$$

Conséquence : \widehat{Q}_θ^π change lentement et π change potentiellement drastiquement

Solution : **Geler périodiquement l'approximateur du « TD target ».** Pratiquement, ça veut dire créer un deuxième réseau neuronaux appelé le « **target Q-network** » et « fetcher » les paramètres du « Q-network » à tout les N « timestep ».

$$\widehat{Q}_{\theta_{\text{target}}}^\pi \longleftarrow \widehat{Q}_\theta^\pi$$

- ▷ Diminue la possibilité d'oscillation ou de divergence de la politique ;
- ▷ Réduit la corrélation entre le « TD target » $r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \widehat{Q}_{\theta_{\text{target}}}^\pi(\mathbf{s}', \mathbf{a}')$ et l'estimateur $\widehat{Q}_\theta^\pi(\mathbf{s}, \mathbf{a})$

⚠ Q-learning avec réseau neuronaux est instable

Causes 2 : Le « *TD-target* » :

- ▷ est non-stationnaire contrairement au contexte de l'apprentissage supervisé ;
- ▷ est corrélé avec l'approximateur \widehat{Q}_θ^π qu'on cherche à apprendre ;

$$L_k(\theta_k) = \mathbb{E}_{\mathbf{s}' \sim p(\cdot | \mathbf{s}, \mathbf{a})} \left[\underbrace{\left(r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}' \in \mathcal{A}(\mathbf{s}')} \widehat{Q}_\theta^\pi(\mathbf{s}', \mathbf{a}') \right)}_{\text{TD-target}} - \widehat{Q}_\theta^\pi(\mathbf{s}, \mathbf{a}) \right]^2$$

Conséquence : \widehat{Q}_θ^π change lentement et π change potentiellement drastiquement

Solution : **Geler périodiquement l'approximateur du « TD target ».** Pratiquement, ça veut dire créer un deuxième réseau neuronaux appelé le « **target Q-network** » et « fetcher » les paramètres du « Q-network » à tout les N « timestep ».

$$\widehat{Q}_{\theta_{\text{target}}}^\pi \longleftarrow \widehat{Q}_\theta^\pi$$

- ▷ Diminue la possibilité d'oscillation ou de divergence de la politique ;
- ▷ Réduit la corrélation entre le « TD target » $r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \widehat{Q}_{\theta_{\text{target}}}^\pi(\mathbf{s}', \mathbf{a}')$ et l'estimateur $\widehat{Q}_\theta^\pi(\mathbf{s}, \mathbf{a})$

⚠ Q-learning avec réseau neuronaux est instable

Causes 2 : Le « *TD-target* » :

- ▷ est non-stationnaire contrairement au contexte de l'apprentissage supervisé ;
- ▷ est corrélé avec l'approximateur \widehat{Q}_θ^π qu'on cherche à apprendre ;

$$L_k(\theta_k) = \mathbb{E}_{\mathbf{s}' \sim p(\cdot | \mathbf{s}, \mathbf{a})} \left[\underbrace{\left(r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}' \in \mathcal{A}(\mathbf{s}')} \widehat{Q}_\theta^\pi(\mathbf{s}', \mathbf{a}') \right)}_{\text{TD-target}} - \widehat{Q}_\theta^\pi(\mathbf{s}, \mathbf{a}) \right]^2$$

Conséquence : \widehat{Q}_θ^π change lentement et π change potentiellement drastiquement

Solution : **Geler périodiquement l'approximateur du « TD target ».** Pratiquement, ça veut dire créer un deuxième réseau neuronaux appelé le « **target Q-network** » et « fetcher » les paramètres du « Q-network » à tout les N « timestep ».

$$\widehat{Q}_{\theta_{\text{target}}}^\pi \longleftarrow \widehat{Q}_\theta^\pi$$

- ▷ Diminue la possibilité d'oscillation ou de divergence de la politique ;
- ▷ Réduit la corrélation entre le « TD target » $r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \widehat{Q}_{\theta_{\text{target}}}^\pi(\mathbf{s}', \mathbf{a}')$ et l'estimateur $\widehat{Q}_\theta^\pi(\mathbf{s}, \mathbf{a})$

Deep Q-Network [DQN]



Caractéristique : Apprentissage **EN-ligne** et **OFF-policy** (comme Q-learning approximatif)

Idées clés :

- 1 **Représenter** $Q^\pi \approx \hat{Q}_\theta^\pi$ un réseau neuronal de vecteur de paramètres θ (le « Q-network »)
- 2 Utiliser un « target Q-network » $\hat{Q}_{\theta_{\text{target}}}^\pi$ cloner périodiquement de \hat{Q}_θ^π pour réduire l'oscillation
- 3 Utiliser un « replay buffer » \mathcal{D} pour briser la corrélation entre les échantillons
- 4 Définir une **fonction objectif** par *erreur quadratique moyenne* [MSE] et son **gradient** :

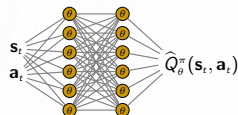
$$L_k(\theta_k) = \mathbb{E}_{(s, a, r, s') \sim \mathcal{U}(\mathcal{D})} \left[\left(y - \hat{Q}_\theta^\pi(s, a) \right)^2 \middle| \theta = \theta_k \right] \quad \text{avec } y = r(s, a, s') + \gamma \max_{a' \in \mathcal{A}(s')} \hat{Q}_{\theta_{\text{target}}}^\pi(s', a')$$

$$\frac{1}{2} \nabla_{\theta_k} L(\theta_k) = \mathbb{E}_{(s, a, r, s') \sim \mathcal{U}(\mathcal{D})} \left[\left(r(s, a, s') + \gamma \max_{a' \in \mathcal{A}(s')} \hat{Q}_{\theta_{\text{target}}}^\pi(s', a') - \hat{Q}_\theta^\pi(s, a) \right) \nabla_{\theta_k} \hat{Q}_\theta^\pi(s, a) \middle| \theta = \theta_k \right]$$

- 5 **Optimiser** l'objectif par *descente de gradient stochastique* [SGD] :

$$\theta_{k+1} \longleftarrow \theta_k - \alpha \nabla_{\theta_k} L(\theta_k)$$

Deep Q-Network [DQN]



Caractéristique : Apprentissage **EN-ligne** et **OFF-policy** (comme Q-learning approximatif)

Idées clés :

- 1 **Représenter** $Q^\pi \approx \hat{Q}_\theta^\pi$ un réseau neuronal de vecteur de paramètres θ (le « **Q-network** »)
- 2 Utiliser un « **target Q-network** » $\hat{Q}_{\theta_{\text{target}}}^\pi$ cloner périodiquement de \hat{Q}_θ^π pour réduire l'oscillation
- 3 Utiliser un « **replay buffer** » \mathcal{D} pour briser la corrélation entre les échantillons
- 4 Définir une **fonction objectif** par *erreur quadratique moyenne* [MSE] et son **gradient** :

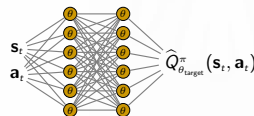
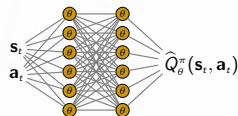
$$L_k(\theta_k) = \mathbb{E}_{(s, a, r, s') \sim \mathcal{U}(\mathcal{D})} \left[\left(y - \hat{Q}_\theta^\pi(s, a) \right)^2 \middle| \theta = \theta_k \right] \quad \text{avec } y = r(s, a, s') + \gamma \max_{a' \in \mathcal{A}(s')} \hat{Q}_{\theta_{\text{target}}}^\pi(s', a')$$

$$\frac{1}{2} \nabla_{\theta_k} L(\theta_k) = \mathbb{E}_{(s, a, r, s') \sim \mathcal{U}(\mathcal{D})} \left[\left(r(s, a, s') + \gamma \max_{a' \in \mathcal{A}(s')} \hat{Q}_{\theta_{\text{target}}}^\pi(s', a') - \hat{Q}_\theta^\pi(s, a) \right) \nabla_{\theta_k} \hat{Q}_\theta^\pi(s, a) \middle| \theta = \theta_k \right]$$

- 5 **Optimiser** l'objectif par *descente de gradient stochastique* [SGD] :

$$\theta_{k+1} \leftarrow \theta_k - \alpha \nabla_{\theta_k} L(\theta_k)$$

Deep Q-Network [DQN]



Caractéristique : Apprentissage **EN-ligne** et **OFF-policy** (comme Q-learning approximatif)

Idées clés :

- 1 **Représenter** $Q^\pi \approx \hat{Q}_\theta^\pi$ un réseau neuronal de vecteur de paramètres θ (le « **Q-network** »)
- 2 Utiliser un « **target Q-network** » $\hat{Q}_{\theta_{\text{target}}}^\pi$ cloner périodiquement de \hat{Q}_θ^π pour réduire l'oscillation
- 3 Utiliser un « **replay buffer** » \mathcal{D} pour briser la corrélation entre les échantillons
- 4 Définir une **fonction objectif** par *erreur quadratique moyenne* [MSE] et son **gradient** :

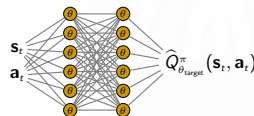
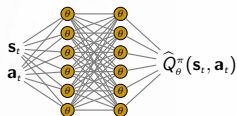
$$L_k(\theta_k) = \mathbb{E}_{(s, a, r, s') \sim \mathcal{U}(\mathcal{D})} \left[\left(y - \hat{Q}_\theta^\pi(s, a) \right)^2 \middle| \theta = \theta_k \right] \quad \text{avec } y = r(s, a, s') + \gamma \max_{a' \in \mathcal{A}(s')} \hat{Q}_{\theta_{\text{target}}}^\pi(s', a')$$

$$\frac{1}{2} \nabla_{\theta_k} L(\theta_k) = \mathbb{E}_{(s, a, r, s') \sim \mathcal{U}(\mathcal{D})} \left[\left(r(s, a, s') + \gamma \max_{a' \in \mathcal{A}(s')} \hat{Q}_{\theta_{\text{target}}}^\pi(s', a') - \hat{Q}_\theta^\pi(s, a) \right) \nabla_{\theta_k} \hat{Q}_\theta^\pi(s, a) \middle| \theta = \theta_k \right]$$

- 5 **Optimiser** l'objectif par *descente de gradient stochastique* [SGD] :

$$\theta_{k+1} \longleftarrow \theta_k - \alpha \nabla_{\theta_k} L(\theta_k)$$

Deep Q-Network [DQN]



Caractéristique : Apprentissage **EN-ligne** et **OFF-policy** (comme Q-learning approximatif)

Idées clés :

- 1 **Représenter** $Q^\pi \approx \hat{Q}_\theta^\pi$ un réseau neuronal de vecteur de paramètres θ (le « **Q-network** »)
- 2 Utiliser un « **target Q-network** » $\hat{Q}_{\theta_{\text{target}}}^\pi$ cloner périodiquement de \hat{Q}_θ^π pour réduire l'oscillation
- 3 Utiliser un « **replay buffer** » \mathcal{D} pour briser la corrélation entre les échantillons
- 4 Définir une **fonction objectif** par *erreur quadratique moyenne* [MSE] et son **gradient** :

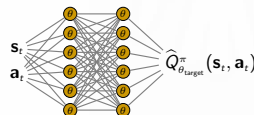
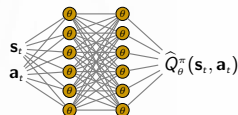
$$L_k(\theta_k) = \mathbb{E}_{(s, a, r, s') \sim \mathcal{U}(\mathcal{D})} \left[\left(y - \hat{Q}_\theta^\pi(s, a) \right)^2 \middle| \theta = \theta_k \right] \quad \text{avec } y = r(s, a, s') + \gamma \max_{a' \in \mathcal{A}(s')} \hat{Q}_{\theta_{\text{target}}}^\pi(s', a')$$

$$\frac{1}{2} \nabla_{\theta_k} L(\theta_k) = \mathbb{E}_{(s, a, r, s') \sim \mathcal{U}(\mathcal{D})} \left[\left(r(s, a, s') + \gamma \max_{a' \in \mathcal{A}(s')} \hat{Q}_{\theta_{\text{target}}}^\pi(s', a') - \hat{Q}_\theta^\pi(s, a) \right) \nabla_{\theta_k} \hat{Q}_\theta^\pi(s, a) \middle| \theta = \theta_k \right]$$

- 5 **Optimiser** l'objectif par *descente de gradient stochastique* [SGD] :

$$\theta_{k+1} \leftarrow \theta_k - \alpha \nabla_{\theta_k} L(\theta_k)$$

Deep Q-Network [DQN]



Caractéristique : Apprentissage **EN-ligne** et **OFF-policy** (comme Q-learning approximatif)

Idées clés :

- 1 **Représenter** $Q^\pi \approx \hat{Q}_\theta^\pi$ un réseau neuronal de vecteur de paramètres θ (le « **Q-network** »)
- 2 Utiliser un « **target Q-network** » $\hat{Q}_{\theta_{\text{target}}}^\pi$ cloner périodiquement de \hat{Q}_θ^π pour réduire l'oscillation
- 3 Utiliser un « **replay buffer** » \mathcal{D} pour briser la corrélation entre les échantillons
- 4 Définir une **fonction objectif** par *erreur quadratique moyenne* [MSE] et son **gradient** :

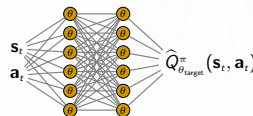
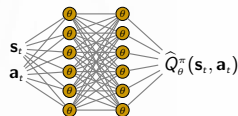
$$L_k(\theta_k) = \mathbb{E}_{(\mathbf{s}, \mathbf{a}, r, \mathbf{s}') \sim \mathcal{U}(\mathcal{D})} \left[\left(y - \hat{Q}_\theta^\pi(\mathbf{s}, \mathbf{a}) \right)^2 \middle| \theta = \theta_k \right] \quad \text{avec } y = r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}' \in \mathcal{A}(\mathbf{s}')} \hat{Q}_{\theta_{\text{target}}}^\pi(\mathbf{s}', \mathbf{a}')$$

$$\frac{1}{2} \nabla_{\theta_k} L(\theta_k) = \mathbb{E}_{(\mathbf{s}, \mathbf{a}, r, \mathbf{s}') \sim \mathcal{U}(\mathcal{D})} \left[\left(r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}' \in \mathcal{A}(\mathbf{s}')} \hat{Q}_{\theta_{\text{target}}}^\pi(\mathbf{s}', \mathbf{a}') - \hat{Q}_\theta^\pi(\mathbf{s}, \mathbf{a}) \right) \nabla_{\theta_k} \hat{Q}_\theta^\pi(\mathbf{s}, \mathbf{a}) \middle| \theta = \theta_k \right]$$

- 5 **Optimiser** l'objectif par *descente de gradient stochastique* [SGD] :

$$\theta_{k+1} \longleftarrow \theta_k - \alpha \nabla_{\theta_k} L(\theta_k)$$

Deep Q-Network [DQN]



Caractéristique : Apprentissage **EN-ligne** et **OFF-policy** (comme Q-learning approximatif)

Idées clés :

- 1 **Représenter** $Q^\pi \approx \hat{Q}_\theta^\pi$ un réseau neuronal de vecteur de paramètres θ (le « **Q-network** »)
- 2 Utiliser un « **target Q-network** » $\hat{Q}_{\theta_{\text{target}}}^\pi$ cloner périodiquement de \hat{Q}_θ^π pour réduire l'oscillation
- 3 Utiliser un « **replay buffer** » \mathcal{D} pour briser la corrélation entre les échantillons
- 4 Définir une **fonction objectif** par *erreur quadratique moyenne* [MSE] et son **gradient** :

$$L_k(\theta_k) = \mathbb{E}_{(\mathbf{s}, \mathbf{a}, r, \mathbf{s}') \sim \mathcal{U}(\mathcal{D})} \left[\left(\mathbf{y} - \hat{Q}_\theta^\pi(\mathbf{s}, \mathbf{a}) \right)^2 \middle| \theta = \theta_k \right] \quad \text{avec } \mathbf{y} = r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}' \in \mathcal{A}(\mathbf{s}')} \hat{Q}_{\theta_{\text{target}}}^\pi(\mathbf{s}', \mathbf{a}')$$

$$\frac{1}{2} \nabla_{\theta_k} L(\theta_k) = \mathbb{E}_{(\mathbf{s}, \mathbf{a}, r, \mathbf{s}') \sim \mathcal{U}(\mathcal{D})} \left[\left(r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}' \in \mathcal{A}(\mathbf{s}')} \hat{Q}_{\theta_{\text{target}}}^\pi(\mathbf{s}', \mathbf{a}') - \hat{Q}_\theta^\pi(\mathbf{s}, \mathbf{a}) \right) \nabla_{\theta_k} \hat{Q}_\theta^\pi(\mathbf{s}, \mathbf{a}) \middle| \theta = \theta_k \right]$$

- 5 **Optimiser** l'objectif par *descente de gradient stochastique* [SGD] :

$$\theta_{k+1} \longleftarrow \theta_k - \alpha \nabla_{\theta_k} L(\theta_k)$$

↺ Algorithme | Deep Q-Network

Soit la politique exploratoire $\mu(\cdot|\mathbf{s})$, la politique cible $\pi(\mathbf{s}) \doteq \pi_{\text{greedy}}$, le « learning rate » $\alpha \in]0, 1]$, le « global step » $k \leftarrow 0$, le « replay buffer » \mathcal{D} , le « Q-network » $\widehat{Q}_{\theta}^{\pi}$ et le « Target-network » $\widehat{Q}_{\theta_{\text{target}}}^{\pi}$

Répéter pour chaque trajectoire τ jusqu'à convergence vers π^* :

Répéter pour chaque « timestep » $t = 0, 1, 2, \dots, t, t+1, \dots, T$ de la trajectoire τ :

```

1       $r_{t+1}^{(k)}, \mathbf{s}_{t+1}^{(k)} \leftarrow$  exécuter  $\mathbf{a}_t \sim \mu(\cdot|\mathbf{s}_t)$  dans l'environnement et observer la réaction
2
3       $\mathcal{D} \leftarrow \mathcal{D} \cup \left\{ \left( \mathbf{s}_t, \mathbf{a}_t, r_{t+1}^{(k)}, \mathbf{s}_{t+1}^{(k)} \right) \right\}$ 
4
5      Prélever une « mini-batch » de transitions  $\left\{ \left( \mathbf{s}_t^{(b)}, \mathbf{a}_t^{(b)}, r_{t+1}^{(b)}, \mathbf{s}_{t+1}^{(b)} \right) \in \mathcal{U}(\mathcal{D}) \right\}$ 
6
7       $\mathbf{y} = r_{t+1}^{(b)} + \gamma \max_{\mathbf{a}' \in \mathcal{A}(\mathbf{s}_{t+1})} \widehat{Q}_{\theta_{\text{target}}}^{\pi}(\mathbf{s}_{t+1}^{(b)}, \mathbf{a}')$ 
8
9       $\theta_{k+1} \leftarrow \theta_k - \alpha \sum_{b \in \mathcal{U}(\mathcal{D})} \left( \mathbf{y} - \widehat{Q}_{\theta_k}^{\pi}(\mathbf{s}_t^{(b)}, \mathbf{a}_t^{(b)}) \right) \frac{\partial}{\partial \theta_k} \widehat{Q}_{\theta_k}^{\pi}(\mathbf{s}_t^{(b)}, \mathbf{a}_t^{(b)})$ 
10
11      $\mathbf{s}_t \leftarrow \mathbf{s}_{t+1}^{(k)}$  et  $k \leftarrow k + 1$ 
12
13     Lorsque « timesteps »  $t = N$  :  $\theta_{\text{target}} \leftarrow \theta_k$ 
```

↺ Algorithme | Deep Q-Network

Soit la politique exploratoire $\mu(\cdot|\mathbf{s})$, la politique cible $\pi(\mathbf{s}) \doteq \pi_{\text{greedy}}$, le « learning rate » $\alpha \in]0, 1]$,
 le « global step » $k \leftarrow 0$, le « replay buffer » \mathcal{D} , le « Q-network » $\widehat{Q}_{\theta}^{\pi}$ et le « Target-network » $\widehat{Q}_{\theta_{\text{target}}}^{\pi}$

Répéter pour chaque trajectoire τ jusqu'à convergence vers π^* :

Répéter pour chaque « timestep » $t = 0, 1, 2, \dots, t, t+1, \dots, T$ de la trajectoire τ :

1 $r_{t+1}^{(k)}, \mathbf{s}_{t+1}^{(k)} \leftarrow$ exécuter $\mathbf{a}_t \sim \mu(\cdot|\mathbf{s}_t)$ dans l'environnement et observer la réaction

2 $\mathcal{D} \leftarrow \mathcal{D} \cup \left\{ \left(\mathbf{s}_t, \mathbf{a}_t, r_{t+1}^{(k)}, \mathbf{s}_{t+1}^{(k)} \right) \right\}$

3 Prélever une « mini-batch » de transitions $\left\{ \left(\mathbf{s}_t^{(b)}, \mathbf{a}_t^{(b)}, r_{t+1}^{(b)}, \mathbf{s}_{t+1}^{(b)} \right) \in \mathcal{U}(\mathcal{D}) \right\}$

4 $\mathbf{y} = r_{t+1}^{(b)} + \gamma \max_{\mathbf{a}' \in \mathcal{A}(\mathbf{s}_t)} \widehat{Q}_{\theta_{\text{target}}}^{\pi}(\mathbf{s}_{t+1}^{(b)}, \mathbf{a}')$

5 $\theta_{k+1} \leftarrow \theta_k - \alpha \sum_{b \in \mathcal{U}(\mathcal{D})} \left(\mathbf{y} - \widehat{Q}_{\theta_k}^{\pi}(\mathbf{s}_t^{(b)}, \mathbf{a}_t^{(b)}) \right) \frac{\partial}{\partial \theta_k} \widehat{Q}_{\theta_k}^{\pi}(\mathbf{s}_t^{(b)}, \mathbf{a}_t^{(b)})$

6 $\mathbf{s}_t \leftarrow \mathbf{s}_{t+1}^{(k)}$ et $k \leftarrow k + 1$

7 Lorsque « timesteps » $t = N$: $\theta_{\text{target}} \leftarrow \theta_k$

- 👍 **Pro :**
- « ... DQN **excels at reactive games**. More recent versions are now better at game involving planning ... »¹ ;
 - Est efficaces au niveau de l'utilisation des échantillons ;
- 👎 **Con :**
- Peu être **très long avant de voir signe de vie** pendant l'entraînement ;
 - Peut être **difficile à stabiliser** ;
 - DQN ne fonctionne pas sur les espaces d'actions continues (dans sa forme classique) ;

Autre architecture et amélioration reliée à DQN :

- *Double DQN*[1]
- *Dueling DQN*[2]
- *Q-learning avec N-step returns*[3]
- *Prioritized experience replay*[4]
- Pour espace d'action continue : **Deep Deterministic Policy Gradient**² [5]

1. Vlad Mnih (Deepmind), [Deep RL Bootcamp Lecture 3: Deep Q-Networks, Berkeley, August 2017](#)

2. DDPG est un mélange de méthode *basée valeur* et *recherche de politique* qui utilise un autre réseau neuronaux pour représenter une politique $\mu_{\phi}(s) \approx \arg \max_a Q_{\theta}^{\pi}(s, a)$

- 👍 **Pro :**
- « ... DQN **excels at reactive games**. More recent versions are now better at game involving planning ... »¹;
 - Est efficaces au niveau de l'utilisation des échantillons ;
- 👎 **Con :**
- Peu être **très long avant de voir signe de vie** pendant l'entraînement ;
 - Peut être **difficile à stabiliser** ;
 - DQN ne fonctionne pas sur les espaces d'actions continues (dans sa forme classique) ;

Autre architecture et amélioration reliée à DQN :

- **Double DQN[1]** ← vue à la prochaine section
- *Dueling DQN[2]*
- *Q-learning avec N-step returns[3]*
- *Prioritized experience replay[4]*
- Pour espace d'action continue : **Deep Deterministic Policy Gradient²** [5]

1. Vlad Mnih (Deepmind), [Deep RL Bootcamp Lecture 3: Deep Q-Networks, Berkeley, August 2017](#)

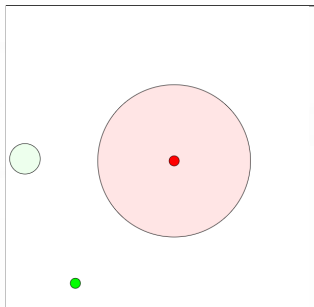
2. DDPG est un mélange de méthode *basée valeur* et *recherche de politique* qui utilise un autre réseau neuronaux pour représenter une politique $\mu_{\phi}(s) \approx \arg \max_a Q_{\theta}^{\pi}(s, a)$

Deep Q-Network [DQN]

Exemple interactif

Deep Q-Network dans l'environnement *PuckWorld*/*WaterWorld*

► Exemple interactif: *PuckWorld*



► Exemple interactif: *WaterWorld*

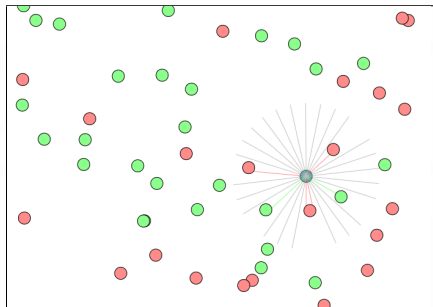


Image : REINFORCEjs, Stanford

■ REINFORCEjs

► [Projet](#)

► [GitHub](#)

DRL et méthode basée sur les valeurs

Double Deep Q-Network

Biais de maximisation

⚠ Problème : Le « TD target » du *DQN* à tendance à être **trop optimiste** dans ses estimations de la Q-valeur du prochain état.

Cause : *DQN* calcule un **maximum sur un approximateur qui est imparfait et bruité** (particulièrement au début de l'entraînement)

$$y = r(s, a, s') + \gamma \max_{a' \in \mathcal{A}(s')} \hat{Q}_{\theta}^{\pi}(s', a')$$

Supposons la vraie Q-fonction

$$Q(s, a) = 0 \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}$$

et son estimateur

$$\hat{Q}^{\pi}(s, a) \in [-1, 1]$$

alors le maximum de la vraie valeur sera 0 et le maximum de son estimateur sera **biaisé de façon positive**.¹

Solution : Utiliser un réseau neuronal **différent pour l'évaluation** des Q-valeur **de celui pour la sélection** d'action.

1. Exemple tiré de « Reinforcement Learning : An introduction », section 6.7 « Maximization Bias and Double Learning » par Sutton & Barto [6]

Biais de maximisation

⚠ Problème : Le « TD target » du *DQN* à tendance à être **trop optimiste** dans ses estimations de la Q-valeur du prochain état.

Cause : *DQN* calcule un **maximum sur un approximateur qui est imparfait et bruité** (particulièrement au début de l'entraînement)

$$y = r(s, a, s') + \gamma \max_{a' \in \mathcal{A}(s')} \hat{Q}_{\theta}^{\pi}(s', a')$$

Supposons la vraie Q-fonction

$$Q(s, a) = 0 \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}$$

et son estimateur

$$\hat{Q}^{\pi}(s, a) \in [-1, 1]$$

alors le maximum de la vraie valeur sera 0 et le maximum de son estimateur sera **biaisé de façon positive**.¹

Solution : Utiliser un réseau neuronal **différent pour l'évaluation** des Q-valeur **de celui pour la sélection** d'action.

1. Exemple tiré de « Reinforcement Learning : An introduction », section 6.7 « Maximization Bias and Double Learning » par Sutton & Barto [6]

Biais de maximisation

⚠ Problème : Le « TD target » du *DQN* à tendance à être **trop optimiste** dans ses estimations de la Q-valeur du prochain état.

Cause : *DQN* calcule un **maximum sur un approximateur qui est imparfait et bruité** (particulièrement au début de l'entraînement)

$$y = r(s, a, s') + \gamma \max_{a' \in \mathcal{A}(s')} \hat{Q}_{\theta}^{\pi}(s', a')$$

Supposons la vraie Q-fonction

$$Q(s, a) = 0 \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}$$

et son estimateur

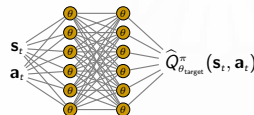
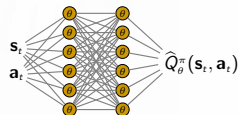
$$\hat{Q}^{\pi}(s, a) \in [-1, 1]$$

alors le maximum de la vraie valeur sera 0 et le maximum de son estimateur sera **biaisé de façon positive**.¹

Solution : Utiliser un réseau neuronal **différent pour l'évaluation** des Q-valeur **de celui pour la sélection** d'action.

1. Exemple tiré de « Reinforcement Learning : An introduction », section 6.7 « Maximization Bias and Double Learning » par Sutton & Barto [6]

Double Deep Q-Network [DDQN]



Idées clés :

- 1 Extension de DQN ;
- 2 Utiliser un **réseau différent pour la sélection d'action** et l'évaluation de Q-valeur dans le « TD target »

$$L_k(\theta_k) = \mathbb{E}_{(s, a, r, s') \sim \mathcal{U}(\mathcal{D})} \left[\left(y - \widehat{Q}_{\theta}^{\pi}(s, a) \right)^2 \middle| \theta = \theta_k \right]$$

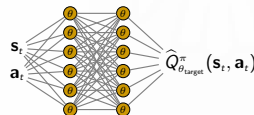
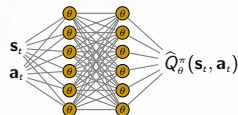
avec le « TD target » $y = r(s, a, s') + \gamma \widehat{Q}_{\theta_{\text{target}}}^{\pi} \left(s', \arg \max_{a' \in \mathcal{A}(s')} \widehat{Q}_{\theta}^{\pi}(s', a') \right)$



Pro :

- Amélioration efficace et simple à implémenter ;
- Aucun effet négatif ;

Double Deep Q-Network [DDQN]



Idées clés :

- 1 Extension de DQN ;
- 2 Utiliser un **réseau différent pour** la sélection d'action et **l'évaluation de Q-valeur** dans le « TD target »

$$L_k(\theta_k) = \mathbb{E}_{(s, a, r, s') \sim \mathcal{U}(\mathcal{D})} \left[\left(y - \hat{Q}_{\theta}^{\pi}(s, a) \right)^2 \middle| \theta = \theta_k \right]$$

avec le « TD target » $y = r(s, a, s') + \gamma \hat{Q}_{\theta_{\text{target}}}^{\pi} \left(s', \arg \max_{a' \in \mathcal{A}(s')} \hat{Q}_{\theta}^{\pi}(s', a') \right)$



Pro :

- Amélioration efficace et simple à implémenter ;
- Aucun effet négatif ;

Ressource théorique

Experience replay

- « Reinforcement learning for robots using neural networks » [► Publication](#)
par Lin, L.-J. [7]

Deep Q-network

- « Playing Atari with Deep Reinforcement Learning » [► Publication](#)
par Mnih et al, 2013 [8]
- « Human-level control through deep reinforcement learning », Nature Vol 518 [► Publication](#)
par Mnih et al, 2015 [9]

Double DQN

- « Deep Reinforcement Learning with Double Q-Learning » [► Publication](#)
par Hasselt et al, 2015 [1]
- « Reinforcement Learning : An introduction », section 6.7 « Maximization Bias and Double Learning » [6] [► Publication](#)
par Sutton & Barto

Méthode par recherche de politique

- 1 Concepts clés en RL (la suite)
- 2 Apprentissage par renforcement profond
- 3 DRL et méthode basée sur les valeurs
- 4 Méthode par recherche de politique**
- 5 DRL et méthode *Policy Gradient*
- 6 Pour aller plus loin

Méthode par recherche de politique (en général)

Caractéristiques :

- **Apprends explicitement la politique π^*** sans passer par V^π ou Q^π ;
- C'est un **problème d'optimisation** \implies on doit définir un objectif et choisir une méthode d'optimisation ;

Idées clés :

- 1 **Représenter** $\pi(a|s) \approx \hat{\pi}_\theta(a|s)$ un approximateur de fonction de paramètres θ ;
- 2 **Définir** la fonction objectif de façon à **mesurer la qualité de la politique $\hat{\pi}_\theta$**

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [G(\tau)]$$

(Ex. du cas épisodique)

- 3 **Trouver le paramètre θ^*** qui maximise $J(\theta)$;

★ | L'objectif des méthodes par recherche de politique en RL

$$\begin{aligned} \theta^* &= \arg \max_{\theta} J(\theta) \\ &= \arg \max_{\theta} \mathbb{E}_{\tau \sim \pi_\theta} [G(\tau)] \\ &= \arg \max_{\theta} \sum_{t=1}^T \mathbb{E}_{(s_t, a_t, s_{t+1}) \sim \pi_\theta(\tau)} [r(s_t, a_t, s_{t+1})] \end{aligned}$$

Méthode par recherche de politique (en général)

Caractéristiques :

- **Apprends explicitement la politique π^*** sans passer par V^π ou Q^π ;
- C'est un **problème d'optimisation** \implies on doit définir un objectif et choisir une méthode d'optimisation ;

Idées clés :

- 1 **Représenter** $\pi(\mathbf{a}|\mathbf{s}) \approx \hat{\pi}_\theta(\mathbf{a}|\mathbf{s})$ un approximateur de fonction de paramètres θ ;
- 2 Définir la fonction objectif de façon à **mesurer la qualité de la politique $\hat{\pi}_\theta$**

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [G(\tau)]$$

(Ex. du cas épisodique)

- 3 **Trouver le paramètre θ^*** qui maximise $J(\theta)$;

★ | L'objectif des méthodes par recherche de politique en RL

$$\begin{aligned} \theta^* &= \arg \max_{\theta} J(\theta) \\ &= \arg \max_{\theta} \mathbb{E}_{\tau \sim \pi_\theta} [G(\tau)] \\ &= \arg \max_{\theta} \sum_{t=1}^T \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) \sim \pi_\theta(\tau)} [r(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1})] \end{aligned}$$

Méthode par recherche de politique (en général)

Caractéristiques :

- **Apprends explicitement la politique π^*** sans passer par V^π ou Q^π ;
- C'est un **problème d'optimisation** \implies on doit définir un objectif et choisir une méthode d'optimisation ;

Idées clés :

- 1 **Représenter** $\pi(\mathbf{a}|\mathbf{s}) \approx \hat{\pi}_\theta(\mathbf{a}|\mathbf{s})$ un approximateur de fonction de paramètres θ ;
- 2 **Définir** la fonction objectif de façon à **mesurer la qualité de la politique $\hat{\pi}_\theta$**

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [G(\tau)] \quad \langle \text{ Ex. du cas épisodique } \rangle$$

- 3 **Trouver le paramètre θ^*** qui maximise $J(\theta)$;

★ | L'objectif des méthodes par recherche de politique en RL

$$\begin{aligned} \theta^* &= \arg \max_{\theta} J(\theta) \\ &= \arg \max_{\theta} \mathbb{E}_{\tau \sim \pi_\theta} [G(\tau)] \\ &= \arg \max_{\theta} \sum_{t=1}^T \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) \sim \pi_\theta(\tau)} [r(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1})] \end{aligned}$$

Méthode par recherche de politique (en général)

Caractéristiques :

- Apprends explicitement la politique π^* sans passer par V^π ou Q^π ;
- C'est un **problème d'optimisation** \implies on doit définir un objectif et choisir une méthode d'optimisation ;

Idées clés :

- 1 Représenter $\pi(\mathbf{a}|\mathbf{s}) \approx \hat{\pi}_\theta(\mathbf{a}|\mathbf{s})$ un approximateur de fonction de paramètres θ ;
- 2 Définir la fonction objectif de façon à **mesurer la qualité de la politique** $\hat{\pi}_\theta$

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [G(\tau)] \quad \langle \text{ Ex. du cas épisodique } \rangle$$

- 3 Trouver le paramètre θ^* qui maximise $J(\theta)$;

★ | L'objectif des méthodes par recherche de politique en RL

$$\begin{aligned} \theta^* &= \arg \max_{\theta} J(\theta) \\ &= \arg \max_{\theta} \mathbb{E}_{\tau \sim \pi_\theta} [G(\tau)] \\ &= \arg \max_{\theta} \sum_{t=1}^T \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) \sim \pi_\theta(\tau)} [r(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1})] \end{aligned}$$

Méthode par recherche de politique (en général)

- 👍 **Pro :**
 - Efficace sur les **espaces en haute dimension** ;
 - Capable de manipuler des **espaces d'actions continues** ;
 - Capable d'apprendre une politique stochastique ;

- 👎 **Con :** **Converge vers un optimum local**¹ (en général) ;

1. Voir [UCL Course on RL](#), Lecture 7: Policy Gradient Methods par David Silver

Méthode par recherche de politique (en général)

Remarque sur les méthodes d'optimisation de politique

Le problème d'optimisation de politique peut être résolu selon deux types d'approches :

- ▶ les approches **utilisant le gradient** de la politique [*Policy gradient*]
- ▶ les approches **sans gradient** [*Gradient-free*]
 - « Hill climbing », « Simulated Annealing »
 - Algorithme génétique, stratégie évolutive¹ [*ES*]
 - Entropy croisé [*CEM*]
 - ...

Prenez note que l'approche « *Gradient-free* » peut être plus approprié dans certaine circonstance par exemple : quand la politique est non différentiable

1. Voir l'excellent post « [Evolution Strategies](#) » par Lilian Weng ainsi que « [A Visual Guide to Evolution Strategies](#) » par David Ha

Méthode par recherche de politique (en général)

Remarque sur les méthodes d'optimisation de politique

Le problème d'optimisation de politique peut être résolu selon deux types d'approches :

- ▶ les approches **utilisant le gradient** de la politique [*Policy gradient*] \Leftarrow **ce que nous allons voir**
- ▶ les approches **sans gradient** [*Gradient-free*]
 - « Hill climbing », « Simulated Annealing »
 - Algorithme génétique, stratégie évolutive¹ [*ES*]
 - Entropy croisé [*CEM*]
 - ...

Prenez note que l'approche « *Gradient-free* » peut être plus approprié dans certaine circonstance par exemple : quand la politique est non différentiable

1. Voir l'excellent post « [Evolution Strategies](#) » par Lilian Weng ainsi que « [A Visual Guide to Evolution Strategies](#) » par David Ha

DRL et méthode *Policy Gradient*

1 Concepts clés en RL (la suite)

2 Apprentissage par renforcement profond

3 DRL et méthode basée sur les valeurs

4 Méthode par recherche de politique

5 DRL et méthode *Policy Gradient*

- REINFORCE
- REINFORCE avec *reward-to-go*
- *Advantage Actor-Critic*

6 Pour aller plus loin

Méthode *Policy Gradient* [Gradient based policy search]

Assomption : 1 La politique $\hat{\pi}_{\theta}$ est **différentiable** par rapport à son paramètre θ ; ¹

2 On sait comment **résoudre le gradient** $\nabla_{\theta} \hat{\pi}_{\theta}(\mathbf{a}|\mathbf{s})$;

★ **Remarque :** On peut utiliser n'importe quel politique dans la mesure où elle respecte l'assomption de différentiabilité. En pratique on utilise une politique stochastique à cause de ses propriétés.

1. Voir « Reinforcement Learning : An introduction », section 13.1 « Policy Approximation and its Advantages » par Sutton & Barto [6]

Méthode *Policy Gradient* [Gradient based policy search]

Politique stochastique paramétrée $\pi_{\theta}(\mathbf{a}|\mathbf{s})$

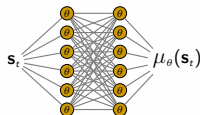
$$\pi_{\theta}(\mathbf{a}|\mathbf{s}) = \mathbb{P}[A_t = \mathbf{a}, S_t = \mathbf{s} | \theta] \quad \text{avec le vecteur de paramètres } \theta$$

Propriétés :

- ▶ **fonction lisse** + continue \implies différentiable ;
- ▶ capable de manipuler des espaces d'**actions continues** ;
- ▶ capable de gérer les cas où deux actions sont de valeur égale ;

Exemple de politique stochastique paramétré compatible avec les espaces d'actions continues :

$$\pi_{\theta}(\mathbf{a}|\mathbf{s}) = \mathcal{N}(\mu_{\theta}(\mathbf{s}), \Sigma^2) = \frac{1}{\sqrt{2\pi\Sigma^2}} \exp\left(-\frac{(\mathbf{a} - \mu_{\theta}(\mathbf{s}))^2}{\Sigma^2}\right)$$



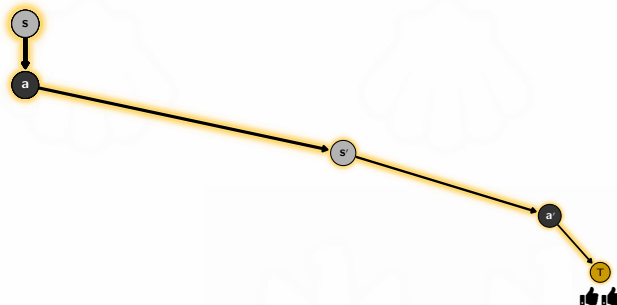
DRL et méthode *Policy Gradient*

REINFORCE

REINFORCE [Monte-Carlo Policy Gradient]

« It's a formalization of the notion of trial & error » – Sergey Levine

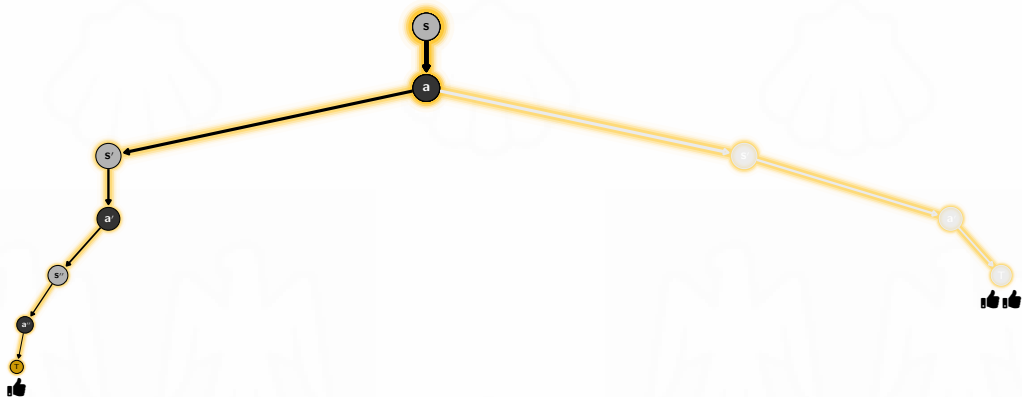
Idée générale : Attribuer un poids à chaque échantillon de trajectoire en fonction de sa qualité.
(Plus le retour est élevé, plus le poids est gros).



REINFORCE [Monte-Carlo Policy Gradient]

« It's a formalization of the notion of trial & error » – Sergey Levine

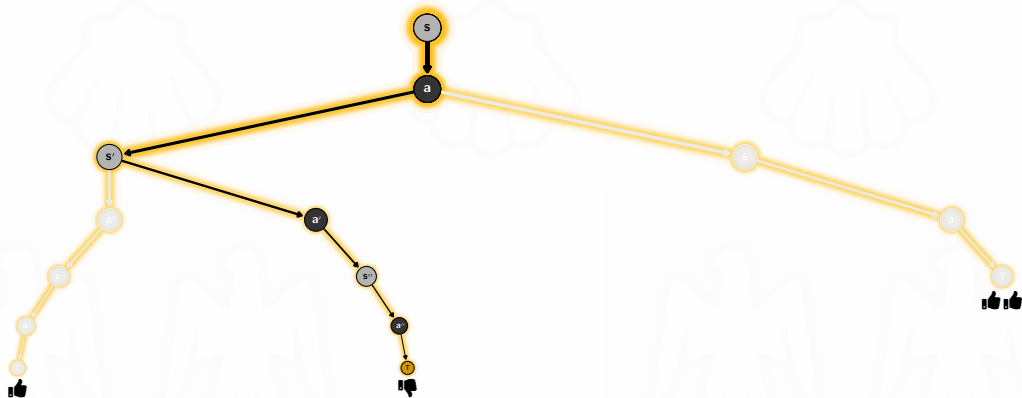
Idée générale : Attribuer un poids à chaque échantillon de trajectoire en fonction de sa qualité.
(Plus le retour est élevé, plus le poids est gros).



REINFORCE [Monte-Carlo Policy Gradient]

« It's a formalization of the notion of trial & error » – Sergey Levine

Idee générale : Attribuer un poids à chaque échantillon de trajectoire en fonction de sa qualité.
(Plus le retour est élevé, plus le poids est gros).

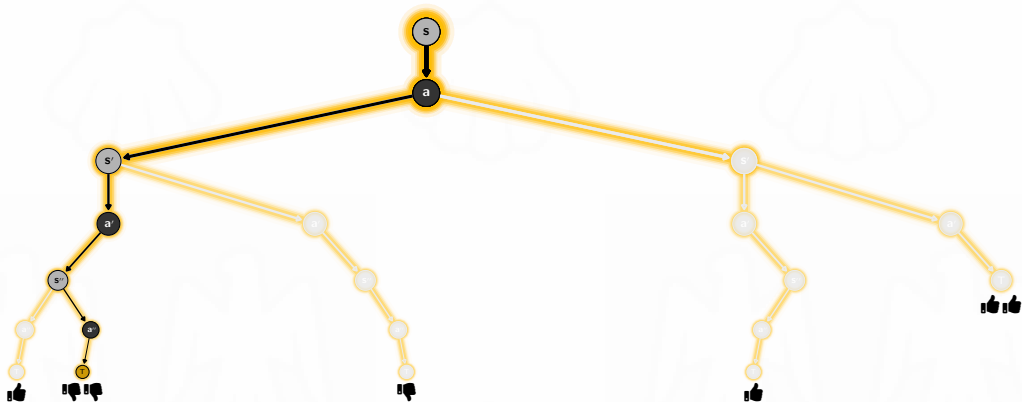


« It's a formalization of the notion of trial & error » – Sergey Levine

REINFORCE [Monte-Carlo Policy Gradient]

« It's a formalization of the notion of trial & error » – Sergey Levine

Idee générale : Attribuer un poids à chaque échantillon de trajectoire en fonction de sa qualité.
(Plus le retour est élevé, plus le poids est gros).



REINFORCE [Monte-Carlo Policy Gradient]

Enjeux 1 : On veut résoudre $\nabla_{\theta} J(\theta)$ pour savoir **dans quelle direction faire l'ascension du gradient** dans le but d'améliorer π_{θ} ;

Étant donné l'assomption que la politique est différentiable, on procède par différentiation directe de π_{θ} ;

Différentiation directe de la politique

$$\begin{aligned}
 J(\theta) &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} [G(\tau)] \\
 \nabla_{\theta} J(\theta) &= \nabla_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} [G(\tau)] \\
 &= \int \nabla_{\theta} \pi_{\theta}(\tau) G(\tau) d\tau \\
 &= \dots \quad \left(\text{ Voir l'annexe Complément théorique pour la } \underline{\text{dérivation complète}} \right) \\
 &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) G(\tau) \right]
 \end{aligned}$$

Observation : Ce résultat montre qu'on n'a pas besoin du modèle pour calculer $\nabla_{\theta} J(\theta)$. On a seulement besoin de la politique π_{θ} et de la fonction de récompense $r(\mathbf{s}, \mathbf{a})$

REINFORCE [Monte-Carlo Policy Gradient]

Enjeux 1 : On veut résoudre $\nabla_{\theta} J(\theta)$ pour savoir **dans quelle direction faire l'ascension du gradient** dans le but d'améliorer π_{θ} ;

Étant donné l'assomption que la politique est différentiable, on procède par différentiation directe de π_{θ} ;

Différentiation directe de la politique

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} [G(\tau)]$$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} [G(\tau)]$$

$$= \int \nabla_{\theta} \pi_{\theta}(\tau) G(\tau) d\tau$$

$$= \dots$$

(Voir l'annexe *Complément théorique* pour la dérivation complète)

$$= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) G(\tau) \right]$$

Observation : Ce résultat montre qu'on n'a **pas besoin du modèle** pour calculer $\nabla_{\theta} J(\theta)$. On a seulement besoin de la politique π_{θ} et de la fonction de récompense $r(\mathbf{s}, \mathbf{a})$

REINFORCE [*Monte-Carlo Policy Gradient*]

Enjeux 2 : Avant de pouvoir optimiser θ **on doit évaluer** $J(\theta)$;

⚠ **Problème** : Les trajectoires τ possibles sont très nombreuses et en haute dimension, donc **on ne peut pas calculer τ pour toute les éventualités** et encore moins de façon exacte ;

Solution : On peut approximer τ en échantillonnant à partir de sa distribution, c.-à-d. en exécutant des trajectoires dans un environnement par rapport à la politique courante ; (C'est un estimé de Monte-Carlo).

Estimé de Monte-Carlo

$$\begin{aligned} J(\theta) &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} \left[\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right] \\ &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T r_t^{(i)} \end{aligned}$$

avec N le nombre d'échantillons de trajectoires τ

REINFORCE [Monte-Carlo Policy Gradient]

Enjeux 2 : Avant de pouvoir optimiser θ **on doit évaluer** $J(\theta)$;

⚠ Problème : Les trajectoires τ possibles sont très nombreuses et en haute dimension, donc **on ne peut pas calculer τ pour toute les éventualités** et encore moins de façon exacte ;

Solution : On peut approximer τ en échantillonnant à partir de sa distribution, c.-à-d. en exécutant des trajectoires dans un environnement par rapport à la politique courante ; (C'est un estimé de Monte-Carlo).

Estimé de Monte-Carlo

$$\begin{aligned} J(\theta) &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} \left[\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right] \\ &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T r_t^{(i)} \end{aligned}$$

avec N le nombre d'échantillons de trajectoires \mathcal{T}

REINFORCE [*Monte-Carlo Policy Gradient*]

Enjeux 2 : Avant de pouvoir optimiser θ **on doit évaluer** $J(\theta)$;

⚠ Problème : Les trajectoires τ possibles sont très nombreuses et en haute dimension, donc **on ne peut pas calculer τ pour toute les éventualités** et encore moins de façon exacte ;

Solution : On peut approximer τ en échantillonnant à partir de sa distribution, c.-à-d. en exécutant des trajectoires dans un environnement par rapport à la politique courante ; (C'est un estimé de Monte-Carlo).

Estimé de Monte-Carlo

$$\begin{aligned} J(\theta) &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} \left[\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right] \\ &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T r_t^{(i)} \end{aligned}$$

avec N le nombre d'échantillons de trajectoires τ

REINFORCE [Monte-Carlo Policy Gradient]

En combinant le résultat de la **différentiation directe de la politique** et celui de l'estimé de **Monte-Carlo** on obtient :

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) G(\tau) \right] \\
 &= \dots \quad \langle \text{ Voir l'annexe Complément théorique pour la } \underline{\text{dérivation complète}} \rangle \\
 &\approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau^{(i)}) G(\tau^{(i)})
 \end{aligned}$$

REINFORCE [Monte-Carlo Policy Gradient]

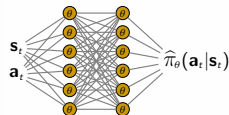
En combinant le résultat de la **différentiation directe de la politique** et celui de l'**estimé de Monte-Carlo** on obtient :

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) G(\tau) \right] \\
 &= \dots \quad \langle \text{ Voir l'annexe Complément théorique pour la dérivation complète } \rangle \\
 &\approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau^{(i)}) G(\tau^{(i)})
 \end{aligned}$$

- ▶ G est utilisé comme **mesure de qualité** de τ
- ▶ G est un **poids** attribué à chaque échantillon



REINFORCE [Monte-Carlo Policy Gradient]



Caractéristiques :

- Apprentissage **HORS-ligne** et **ON-policy** ;
- Optimise une politique stochastique $\pi_{\theta}(a_t | s_t) = \mathbb{P}[A_t = s, S_t = a | \theta]$

Idées clés :

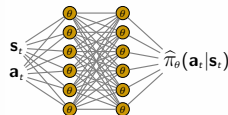
- 1 **Représenter** $\pi(a|s) \approx \hat{\pi}_{\theta}(a|s)$ un réseau neuronal de vecteur de paramètres θ (le « Policy-network ») ;
- 2 **Définir** la fonction objectif : $J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} [G(\tau)]$
- 3 Résoudre l'objectif par **différentiation direct de la politique**
- 4 et **approximer τ par un estimé de Monte-Carlo** (c.a.d. exécuter la politique dans l'environnement)

$$\nabla_{\theta_k} J(\theta_k) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta_k} \log \hat{\pi}_{\theta_k}(\tau^{(i)}) G(\tau^{(i)})$$

- 5 **Optimiser** l'objectif par *ascension du gradient* :

$$\theta_{k+1} \leftarrow \theta_k + \alpha \nabla_{\theta} J(\theta_k)$$

REINFORCE [Monte-Carlo Policy Gradient]



Caractéristiques :

- Apprentissage **HORS-ligne** et **ON-policy** ;
- Optimise une politique stochastique $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) = \mathbb{P}[A_t = \mathbf{s}, S_t = \mathbf{a} | \theta]$

Idées clés :

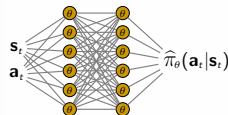
- 1 **Représenter** $\pi(\mathbf{a} | \mathbf{s}) \approx \hat{\pi}_{\theta}(\mathbf{a} | \mathbf{s})$ un réseau neuronal de vecteur de paramètres θ (le « Policy-network ») ;
- 2 **Définir** la fonction objectif : $J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} [G(\tau)]$
- 3 Résoudre l'objectif par **différentiation direct de la politique**
- 4 et **approximer τ par un estimé de Monte-Carlo** (c.a.d. exécuter la politique dans l'environnement)

$$\nabla_{\theta_k} J(\theta_k) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta_k} \log \hat{\pi}_{\theta_k}(\tau^{(i)}) G(\tau^{(i)})$$

- 5 **Optimiser** l'objectif par *ascension du gradient* :

$$\theta_{k+1} \leftarrow \theta_k + \alpha \nabla_{\theta} J(\theta_k)$$

REINFORCE [Monte-Carlo Policy Gradient]



Caractéristiques :

- Apprentissage **HORS-ligne** et **ON-policy** ;
- Optimise une politique stochastique $\pi_{\theta}(a_t | s_t) = \mathbb{P}[A_t = s, S_t = a | \theta]$

Idées clés :

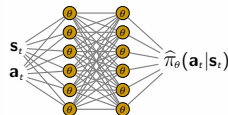
- 1 **Représenter** $\pi(a|s) \approx \hat{\pi}_{\theta}(a|s)$ un réseau neuronal de vecteur de paramètres θ (le « Policy-network ») ;
- 2 **Définir** la fonction objectif : $J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} [G(\tau)]$
- 3 Résoudre l'objectif par **différentiation direct de la politique**
- 4 et **approximer τ par un estimé de Monte-Carlo** (c.a.d. exécuter la politique dans l'environnement)

$$\nabla_{\theta_k} J(\theta_k) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta_k} \log \hat{\pi}_{\theta_k}(\tau^{(i)}) G(\tau^{(i)})$$

- 5 **Optimiser** l'objectif par *ascension du gradient* :

$$\theta_{k+1} \leftarrow \theta_k + \alpha \nabla_{\theta} J(\theta_k)$$

REINFORCE [Monte-Carlo Policy Gradient]



Caractéristiques :

- Apprentissage **HORS-ligne** et **ON-policy** ;
- Optimise une politique stochastique $\pi_{\theta}(a_t | s_t) = \mathbb{P}[A_t = s, S_t = a | \theta]$

Idées clés :

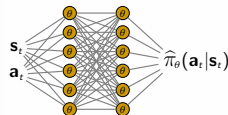
- 1 **Représenter** $\pi(a|s) \approx \hat{\pi}_{\theta}(a|s)$ un réseau neuronal de vecteur de paramètres θ (le « Policy-network ») ;
- 2 **Définir** la fonction objectif : $J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} [G(\tau)]$
- 3 Résoudre l'objectif par **différentiation direct de la politique**
- 4 et **approximer τ par un estimé de Monte-Carlo** (c.a.d. exécuter la politique dans l'environnement)

$$\nabla_{\theta_k} J(\theta_k) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta_k} \log \hat{\pi}_{\theta_k}(\tau^{(i)}) G(\tau^{(i)})$$

- 5 **Optimiser** l'objectif par *ascension du gradient* :

$$\theta_{k+1} \leftarrow \theta_k + \alpha \nabla_{\theta} J(\theta_k)$$

REINFORCE [Monte-Carlo Policy Gradient]



Caractéristiques :

- Apprentissage **HORS-ligne** et **ON-policy** ;
- Optimise une politique stochastique $\pi_{\theta}(a_t | s_t) = \mathbb{P}[A_t = s, S_t = a | \theta]$

Idées clés :

- 1 **Représenter** $\pi(a|s) \approx \hat{\pi}_{\theta}(a|s)$ un réseau neuronal de vecteur de paramètres θ (le « Policy-network ») ;
- 2 **Définir** la fonction objectif : $J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} [G(\tau)]$
- 3 Résoudre l'objectif par **différentiation direct de la politique**
- 4 et **approximer τ par un estimé de Monte-Carlo** (c.a.d. exécuter la politique dans l'environnement)

$$\nabla_{\theta_k} J(\theta_k) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta_k} \log \hat{\pi}_{\theta_k}(\tau^{(i)}) G(\tau^{(i)})$$

- 5 **Optimiser** l'objectif par *ascension du gradient* :

$$\theta_{k+1} \leftarrow \theta_k + \alpha \nabla_{\theta} J(\theta_k)$$

↺ Algorithmme | REINFORCE

Soit le « policy network » $\hat{\pi}_{\theta}(\mathbf{a}|\mathbf{s})$ (une politique différentiable), une « batch » de trajectoire \mathcal{D} , le « learning rate » $\alpha \in]0, 1]$, l'échantillon $^{(i)}$ et le « global step » $k \leftarrow 0$

Répéter pour chaque « global step » $k = 1, 2, 3, \dots$ jusqu'à convergence vers $\approx \pi^*$:

1 Répéter pour chaque échantillon $i = 1, 2, 3, \dots, |\mathcal{D}|$:

 ► Générer une trajectoire $\tau^{(i)}$ en suivant π_{θ_k} jusqu'à terminaison

$$\tau^{(i)} = s_0, a_0, r_1, s_1, a_1, r_2, \dots, s_t, a_t, r_{t+1}, \dots, s_{T-1}, a_{T-1}, r_T$$

 ► $\mathcal{D} \leftarrow \mathcal{D} \cup \{\tau^{(i)}\}$

2 $\nabla_{\theta_k} J(\theta_k) \approx \frac{1}{|\mathcal{D}|} \sum_{\tau^{(i)} \in \mathcal{D}} \nabla_{\theta_k} \log \hat{\pi}_{\theta_k}(\tau^{(i)}) G(\tau^{(i)})$

3 $\theta_{k+1} \leftarrow \theta_k + \alpha \nabla_{\theta_k} J(\theta_k)$

REINFORCE [Monte-Carlo Policy Gradient]

- 👍 **Pro :**
 - Fonctionne sur les **espaces d'action discret et continue** ;
 - **N'utilise pas la propriété de Markov**
 - ⇒ Peu être utilisé sur des MDP partiellement observé sans modification [POMDP] ;
 - L'algorithme est agnostique à la dynamique du système
 - ⇒ Apprentissage **sans modèle** ;
 - Non biaisé ;
- 👎 **Con :**
 - Algorithme **HORS-ligne + ON-policy** ⇒ l'efficience de l'algorithme **dépend de la capacité de l'environnement à générer rapidement des trajectoires complètes**.
(Plus la génération des trajectoires est longue, plus l'entraînement devient inefficace)
 - **Beaucoup de variances** ⇒ mauvaise performance ⚠
 - ★ Impératif d'utiliser une **technique de réduction de variance** pour mitiger le problème :
 - Utiliser le « reward-to-go »
 - Ajouter une « baseline » : retour moyen, récompense espérée pondérée,
 - ...

REINFORCE [Monte-Carlo Policy Gradient]

- 👍 **Pro :**
 - Fonctionne sur les **espaces d'action discret et continue** ;
 - **N'utilise pas la propriété de Markov**
 - ⇒ Peu être utilisé sur des MDP partiellement observé sans modification [POMDP] ;
 - L'algorithme est agnostique à la dynamique du système
 - ⇒ Apprentissage **sans modèle** ;
 - Non biaisé ;
- 👎 **Con :**
 - Algorithme **HORS-ligne + ON-policy** ⇒ l'efficience de l'algorithme **dépend de la capacité de l'environnement à générer rapidement des trajectoires complètes**.
(Plus la génération des trajectoires est longue, plus l'entraînement devient inefficace)
 - **Beaucoup de variances** ⇒ mauvaise performance ⚠
 - ★ Impératif d'utiliser une **technique de réduction de variance** pour mitiger le problème :
 - Utiliser le « reward-to-go »
 - Ajouter une « baseline » : retour moyen, récompense espérée pondérée, ...
 - ...

REINFORCE [*Monte-Carlo Policy Gradient*]

- 👍 **Pro :**
 - Fonctionne sur les **espaces d'action discret et continu** ;
 - **N'utilise pas la propriété de Markov**
⇒ Peu être utilisé sur des MDP partiellement observé sans modification [POMDP] ;
 - L'algorithme est agnostique à la dynamique du système
⇒ Apprentissage **sans modèle** ;
 - Non biaisé ;
- 👎 **Con :**
 - Algorithme **HORS-ligne + ON-policy** ⇒ l'efficacité de l'algorithme **dépend de la capacité de l'environnement à générer rapidement des trajectoires complètes**.
(Plus la génération des trajectoires est longue, plus l'entraînement devient inefficace)
 - **Beaucoup de variances** ⇒ mauvaise performance ⚠
 - ★ Impératif d'utiliser une **technique de réduction de variance** pour mitiger le problème :
 - Utiliser le « reward-to-go »
 - Ajouter une « baseline » : retour moyen, récompense espérée pondérée, ...
 - ...

REINFORCE [Monte-Carlo Policy Gradient]

- 👍 **Pro :**
 - Fonctionne sur les **espaces d'action discret et continue** ;
 - **N'utilise pas la propriété de Markov**
 - ⇒ Peu être utilisé sur des MDP partiellement observé sans modification [POMDP] ;
 - L'algorithme est agnostique à la dynamique du système
 - ⇒ Apprentissage **sans modèle** ;
 - Non biaisé ;
- 👎 **Con :**
 - Algorithme **HORS-ligne + ON-policy** ⇒ l'efficience de l'algorithme **dépend de la capacité de l'environnement à générer rapidement des trajectoires complètes**.
(Plus la génération des trajectoires est longue, plus l'entraînement devient inefficace)
 - **Beaucoup de variances** ⇒ mauvaise performance ⚠
 - ★ Impératif d'utiliser une **technique de réduction de variance** pour mitiger le problème :
 - Utiliser le « reward-to-go »
 - Ajouter une « baseline » : retour moyen, récompense espérée pondérée, ...
 - ...

REINFORCE [Monte-Carlo Policy Gradient]

- 👍 **Pro :**
 - Fonctionne sur les **espaces d'action discret et continue** ;
 - **N'utilise pas la propriété de Markov**
⇒ Peu être utilisé sur des MDP partiellement observé sans modification [POMDP] ;
 - L'algorithme est agnostique à la dynamique du système
⇒ Apprentissage **sans modèle** ;
 - Non biaisé ;
- 👎 **Con :**
 - Algorithme **HORS-ligne + ON-policy** ⇒ l'efficacité de l'algorithme **dépend de la capacité de l'environnement à générer rapidement des trajectoires complètes**.
(Plus la génération des trajectoires est longue, plus l'entraînement devient inefficace)
 - **Beaucoup de variances** ⇒ mauvaise performance ⚠
 - ★ Impératif d'utiliser une **technique de réduction de variance** pour mitiger le problème :
 - Utiliser le « **reward-to-go** » ← **c'est ce qu'on va voir maintenant**
 - Ajouter une « **baseline** » : retour moyen, récompense espérée pondérée, ...
 - ...

DRL et méthode *Policy Gradient*

REINFORCE avec *reward-to-go*

REINFORCE avec *reward-to-go*

⚠ **Problème** : REINFORCE souffre d'avoir **trop de variance** ;

Solution : Utiliser le **principe de causalité** \implies « le futur n'affecte pas le passé »¹ ;

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau^{(i)}) G(\tau^{(i)}) \\
 &= \dots \quad \langle \text{ Voir l'annexe Complément théorique pour les [détails du calcul](#) } \rangle \\
 &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \widehat{Q}_t(\tau^{(i)})
 \end{aligned}$$

Pourquoi ça marche ?

Enlever les récompenses passées **rend la somme des récompenses plus petite**

\implies

la variance $V = \mathbb{E}[(X - \mathbb{E}[X])^2]$ devient également plus petite

Le « reward-to-go »

L'**estimé de Monte-Carlo d'un retour** $G_t(\tau)$ est appelé intuitivement le « **reward-to-go** » $\widehat{Q}_t(\tau^{(i)})$

$$G_t(\tau) = \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \approx \sum_{t'=t}^T r_{t'}^{(i)} = \widehat{Q}_t(\tau^{(i)})$$

1. Voir CS 294–112 Deep Reinforcement Learning : [lecture 5 - Policy Gradient \(à partir de 35:09\)](#) donné par professeur Sergey Levine à l'université Berkeley en Californie ;

REINFORCE avec *reward-to-go*

⚠ **Problème** : REINFORCE souffre d'avoir **trop de variance** ;

Solution : Utiliser le **principe de causalité** \implies « le futur n'affecte pas le passé »¹ ;

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau^{(i)}) G(\tau^{(i)}) \\
 &= \dots \quad \langle \text{ Voir l'annexe Complément théorique pour les [détails du calcul](#) } \rangle \\
 &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \hat{Q}_t(\tau^{(i)})
 \end{aligned}$$

Pourquoi ça marche ?

Enlever les récompenses passées **rend la somme des récompenses plus petite**



la variance $\mathbb{V} = \mathbb{E}[(X - \mathbb{E}[X])^2]$ devient également plus petite

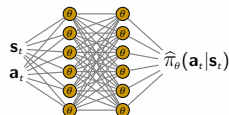
Le « reward-to-go »

L'**estimé de Monte-Carlo d'un retour** $G_t(\tau)$ est appelé intuitivement le « **reward-to-go** » $\hat{Q}_t(\tau^{(i)})$

$$G_t(\tau) = \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \approx \sum_{t'=t}^T r_{t'}^{(i)} = \hat{Q}_t(\tau^{(i)})$$

1. Voir CS 294-112 Deep Reinforcement Learning : [lecture 5 - Policy Gradient](#) (à partir de 35:09) donné par professeur Sergey Levine à l'université Berkeley en Californie ;

REINFORCE avec reward-to-go [Monte-Carlo Policy Gradient]



Idées clés :

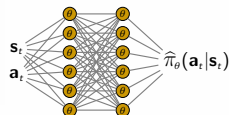
- 1 Extension de « REINFORCE »
- 2 **Technique de réduction de la variance**
- 3 Remplacer le retour $G(\tau^{(i)})$ par le « reward-to-go » $\hat{Q}_t^\pi(\tau^{(i)})$:

$$\hat{Q}_t^\pi(\tau^{(i)}) = \sum_{t'=t}^T r_{t'}^{(i)}$$

$$\nabla_{\theta_k} J(\theta_k) \approx \frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} \sum_{t=1}^T \nabla_{\theta_k} \log \pi_{\theta_k}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \hat{Q}_t(\tau^{(i)})$$

👍 Pro : Simple à implémenter ;

REINFORCE avec reward-to-go [Monte-Carlo Policy Gradient]



Idées clés :

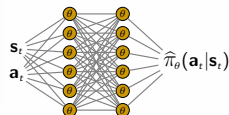
- 1 Extension de « REINFORCE »
- 2 **Technique de réduction de la variance**
- 3 **Remplacer** le retour $G(\tau^{(i)})$ par le « **reward-to-go** » $\hat{Q}_t^{\pi}(\tau^{(i)})$:

$$\hat{Q}_t^{\pi}(\tau^{(i)}) = \sum_{t'=t}^T r_{t'}^{(i)}$$

$$\nabla_{\theta_k} J(\theta_k) \approx \frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} \sum_{t=1}^T \nabla_{\theta_k} \log \pi_{\theta_k}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \hat{Q}_t(\tau^{(i)})$$

👍 Pro : Simple à implémenter ;

REINFORCE avec reward-to-go [Monte-Carlo Policy Gradient]



Idées clés :

- 1 Extension de « REINFORCE »
- 2 **Technique de réduction de la variance**
- 3 **Remplacer** le retour $G(\tau^{(i)})$ par le « **reward-to-go** » $\hat{Q}_t^{\pi}(\tau^{(i)})$:

$$\hat{Q}_t^{\pi}(\tau^{(i)}) = \sum_{t'=t}^T r_{t'}^{(i)}$$

$$\nabla_{\theta_k} J(\theta_k) \approx \frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} \sum_{t=1}^T \nabla_{\theta_k} \log \pi_{\theta_k}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \hat{Q}_t^{\pi}(\tau^{(i)})$$

👍 **Pro : Simple à implémenter ;**

↺ Algorithme | REINFORCE avec reward-to-go

Soit le « policy network » $\hat{\pi}_{\theta}(\mathbf{a}|\mathbf{s})$ (une politique différentiable), une « batch » de trajectoire \mathcal{D} , le « learning rate » $\alpha \in]0, 1]$, l'échantillon $^{(i)}$ et le « global step » $k \leftarrow 0$

Répéter pour chaque « global step » $k = 1, 2, 3, \dots$ jusqu'à convergence vers $\approx \pi^*$:

```

1  Répéter pour chaque échantillon  $i = 1, 2, 3, \dots, |\mathcal{D}|$  :
    ▶ Générer une trajectoire  $\tau^{(i)}$  en suivant  $\pi_{\theta_k}$  jusqu'à terminaison
      
$$\tau^{(i)} = \mathbf{s}_0^{(i)}, \mathbf{a}_0^{(i)}, r_1^{(i)}, \mathbf{s}_1^{(i)}, \mathbf{a}_1^{(i)}, r_2^{(i)}, \dots, \mathbf{s}_t^{(i)}, \mathbf{a}_t^{(i)}, r_{t+1}^{(i)}, \dots, \mathbf{s}_{T-1}^{(i)}, \mathbf{a}_{T-1}^{(i)}, r_T^{(i)}$$

    ▶  $\mathcal{D} \leftarrow \mathcal{D} \cup \{ \tau^{(i)} \}$ 

2   $\nabla_{\theta_k} J(\theta_k) \approx \frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} \sum_{t=1}^T \nabla_{\theta_k} \log \pi_{\theta_k}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \left( \sum_{t'=t}^T r_{t'}^{(i)} \right)$ 

3   $\theta_{k+1} \leftarrow \theta_k + \alpha \nabla_{\theta_k} J(\theta_k)$ 

```

DRL et méthode *Policy Gradient*

Advantage Actor-Critic

Élargir la portée de *REINFORCE* avec « reward-to-go »

$$\nabla_{\theta_k} J(\theta_k) \approx \frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} \sum_{t=1}^T \nabla_{\theta_k} \log \pi_{\theta_k}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \hat{Q}_t(\tau^{(i)})$$

⚠ Problèmes :

- 1 $\hat{Q}_t(\tau^{(i)})$ est **l'estimé d'une espérance basée sur un seul échantillon**
 \Rightarrow haute variance ;
- 2 $\hat{Q}_t(\tau^{(i)})$ **tien compte d'un seul futur possible**, celui qui a été échantillonné dans la trajectoire $\tau^{(i)}$.
 \Rightarrow c'est un estimateur avec une vision très étroite de la réalité

Solution : Remplacer l'estimé à un seul échantillon de trajectoire $\hat{Q}_t(\tau^{(i)})$ par l'**Avantage** :

$$\begin{aligned} A^\pi(\mathbf{s}_t, \mathbf{a}_t) &= Q^\pi(\mathbf{s}_t, \mathbf{a}_t) - V^\pi(\mathbf{s}_t) \\ &\quad \langle \text{Compromis : estimée basée sur un seul échantillon de récompense} \rangle \\ &\approx r_t^{(i)} + V^\pi(\mathbf{s}_{t+1}) - V^\pi(\mathbf{s}_t) \end{aligned}$$

L'**Avantage** $A^\pi(\mathbf{s}_t, \mathbf{a}_t)$ répond à la question suivante :

« **à quel point meilleure** est l'action \mathbf{a}_t dans l'état \mathbf{s}_t **par rapport à la moyenne** de toutes les actions \mathbf{a}'_t possibles dans cet état \mathbf{s}_t ? »

Élargir la portée de *REINFORCE* avec « *reward-to-go* »

$$\nabla_{\theta_k} J(\theta_k) \approx \frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} \sum_{t=1}^T \nabla_{\theta_k} \log \pi_{\theta_k}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \hat{Q}_t(\tau^{(i)})$$

⚠ Problèmes :

- 1 $\hat{Q}_t(\tau^{(i)})$ est **l'estimé d'une espérance basée sur un seul échantillon**
 \Rightarrow haute variance ;
- 2 $\hat{Q}_t(\tau^{(i)})$ **tient compte d'un seul futur possible**, celui qui a été échantillonné dans la trajectoire $\tau^{(i)}$.
 \Rightarrow c'est un estimateur avec une vision très étroite de la réalité

Solution : Remplacer l'estimé à un seul échantillon de trajectoire $\hat{Q}_t(\tau^{(i)})$ par l'**Avantage** :

$$\begin{aligned} A^\pi(\mathbf{s}_t, \mathbf{a}_t) &= Q^\pi(\mathbf{s}_t, \mathbf{a}_t) - V^\pi(\mathbf{s}_t) \\ &\quad \langle \text{Compromis : estimé basée sur un seul échantillon de récompense} \rangle \\ &\approx r_t^{(i)} + V^\pi(\mathbf{s}_{t+1}) - V^\pi(\mathbf{s}_t) \end{aligned}$$

L'**Avantage** $A^\pi(\mathbf{s}_t, \mathbf{a}_t)$ répond à la question suivante :

« **à quel point meilleure** est l'action \mathbf{a}_t dans l'état \mathbf{s}_t **par rapport à la moyenne** de tout les actions \mathbf{a}'_t possibles dans cet état state \mathbf{s}_t ? »

Élargir la portée de *REINFORCE* avec « *reward-to-go* »

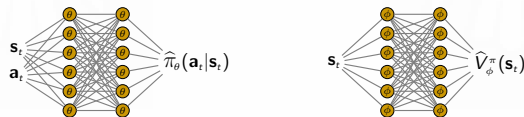
Remarque sur l'*Avantage*

La démonstration théorique expliquant le fondement de la formulation de l'*Avantage* et ses bénéfices est **beaucoup plus consistante** et demande entre autre l'**introduction du concept de « baseline »**.

Pour aller plus loin :

- *Reinforcement Learning : An introduction*, chapitre 13.4 *REINFORCE with Baseline* [6] [► Publication](#)
- *Optimizing expectations : from deep reinforcement learning to stochastic computation graphs*, [10]
[► Thèse de doctorat](#) par John Schulman
- *CS 294-112 Deep Reinforcement Learning*, Université Berkeley
donné par professeur Sergey Levine
 - *lecture 5 - Policy Gradient* (à partir de 35 :00) [► Présentation](#)
 - *lecture 6 - Actor-Critic Algorithms* [► Présentation](#)

ON-line Advantage Actor-Critic



Caractéristiques : Apprentissage **ON-policy** et **EN-ligne** (Il existe aussi une version HORS-ligne) ;

Idées clés :

- 1 Apprendre deux réseaux neuronaux :

$\hat{\pi}_{\theta}(a|s) \rightarrow$ **l'acteur** : celui responsable de prendre les **décisions** dans l'environnement ;

$\hat{V}_{\phi}^{\pi}(s) \rightarrow$ **le critique** : celui responsable d'**évaluer** si $\hat{\pi}_{\theta}$ fait du bon travail ;

- 2 Apprendre $\hat{V}_{\phi}^{\pi}(s_t)$ par MSE avec le « bootstrap target »¹ : $y_k = r_{t+1}^{(i)} + \hat{V}_{\phi}^{\pi}(s_{t+1}^{(i)})$

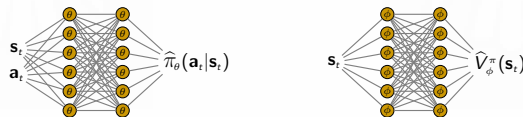
- 3 Calculer l'Avantage² : $\hat{A}^{\pi}(s_t, r_t^{(i)}, s_{t+1}^{(i)}) = r_{t+1}^{(i)} + \hat{V}_{\phi}^{\pi}(s_{t+1}^{(i)}) - \hat{V}_{\phi}^{\pi}(s_t)$

- 4 Résoudre l'objectif : $\nabla_{\theta_k} J(\theta_k) \approx \sum_{t=1}^T \nabla_{\theta_k} \log \pi_{\theta_k}(a_t^{(i)} | s_t) \hat{A}^{\pi}(s_t, r_t^{(i)}, s_{t+1}^{(i)})$

1. On peut également utiliser un « Monte-Carlo target » : $y_k = \hat{Q}_t^{(i)}$

2. Notation : s_t est le point de départ de l'échantillon courant et a été échantillonné à l'itération précédente.

ON-line Advantage Actor-Critic



Caractéristiques : Apprentissage **ON-policy** et **EN-ligne** (Il existe aussi une version HORS-ligne) ;

Idées clés :

1 **Apprendre deux réseaux neuronaux :**

$\hat{\pi}_{\theta}(\mathbf{a}|\mathbf{s}) \rightarrow$ **l'acteur** : celui responsable de prendre les **décisions** dans l'environnement ;

$\hat{V}_{\phi}^{\pi}(\mathbf{s}) \rightarrow$ **le critique** : celui responsable d'**évaluer** si $\hat{\pi}_{\theta}$ fait du bon travail ;

2 **Apprendre $\hat{V}_{\phi}^{\pi}(\mathbf{s}_t)$** par MSE avec le « bootstrap target »¹ : $\mathbf{y}_k = r_{t+1}^{(i)} + \hat{V}_{\phi}^{\pi}(\mathbf{s}_{t+1}^{(i)})$

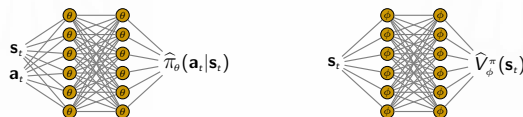
3 **Calculer l'Avantage**² : $\hat{A}^{\pi}(\mathbf{s}_t, r_t^{(i)}, \mathbf{s}_{t+1}^{(i)}) = r_{t+1}^{(i)} + \hat{V}_{\phi}^{\pi}(\mathbf{s}_{t+1}^{(i)}) - \hat{V}_{\phi}^{\pi}(\mathbf{s}_t)$

4 **Résoudre l'objectif** : $\nabla_{\theta_k} J(\theta_k) \approx \sum_{t=1}^T \nabla_{\theta_k} \log \pi_{\theta_k}(\mathbf{a}_t^{(i)} | \mathbf{s}_t) \hat{A}^{\pi}(\mathbf{s}_t, r_t^{(i)}, \mathbf{s}_{t+1}^{(i)})$

1. On peut également utiliser un « Monte-Carlo target » : $\mathbf{y}_k = \hat{Q}_t^{(i)}$

2. Notation : \mathbf{s}_t est le point de départ de l'échantillon courant et a été échantillonné à l'itération précédente.

ON-line Advantage Actor-Critic



Caractéristiques : Apprentissage **ON-policy** et **EN-ligne** (Il existe aussi une version HORS-ligne) ;

Idées clés :

1 **Apprendre deux réseaux neuronaux :**

$\hat{\pi}_{\theta}(\mathbf{a}|\mathbf{s}) \rightarrow$ **l'acteur** : celui responsable de prendre les **décisions** dans l'environnement ;

$\hat{V}_{\phi}^{\pi}(\mathbf{s}) \rightarrow$ **le critique** : celui responsable d'**évaluer** si $\hat{\pi}_{\theta}$ fait du bon travail ;

2 **Apprendre $\hat{V}_{\phi}^{\pi}(\mathbf{s}_t)$** par MSE avec le « bootstrap target »¹ : $\mathbf{y}_k = r_{t+1}^{(i)} + \hat{V}_{\phi}^{\pi}(\mathbf{s}_{t+1}^{(i)})$

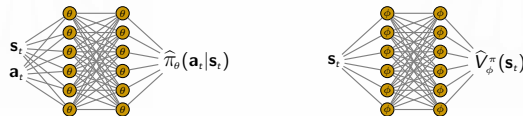
3 **Calculer l'Avantage**² : $\hat{A}^{\pi}(\mathbf{s}_t, r_t^{(i)}, \mathbf{s}_{t+1}^{(i)}) = r_{t+1}^{(i)} + \hat{V}_{\phi}^{\pi}(\mathbf{s}_{t+1}^{(i)}) - \hat{V}_{\phi}^{\pi}(\mathbf{s}_t)$

4 **Résoudre l'objectif** : $\nabla_{\theta_k} J(\theta_k) \approx \sum_{t=1}^T \nabla_{\theta_k} \log \pi_{\theta_k}(\mathbf{a}_t^{(i)} | \mathbf{s}_t) \hat{A}^{\pi}(\mathbf{s}_t, r_t^{(i)}, \mathbf{s}_{t+1}^{(i)})$

1. On peut également utiliser un « Monte-Carlo target » : $\mathbf{y}_k = \hat{Q}_t^{(i)}$

2. Notation : \mathbf{s}_t est le point de départ de l'échantillon courant et a été échantillonné à l'itération précédente.

ON-line Advantage Actor-Critic



Caractéristiques : Apprentissage **ON-policy** et **EN-ligne** (Il existe aussi une version HORS-ligne) ;

Idées clés :

1 **Apprendre deux réseaux neuronaux :**

$\hat{\pi}_{\theta}(\mathbf{a}|\mathbf{s}) \rightarrow$ **l'acteur** : celui responsable de prendre les **décisions** dans l'environnement ;

$\hat{V}_{\phi}^{\pi}(\mathbf{s}) \rightarrow$ **le critique** : celui responsable d'**évaluer** si $\hat{\pi}_{\theta}$ fait du bon travail ;

2 **Apprendre $\hat{V}_{\phi}^{\pi}(\mathbf{s}_t)$** par MSE avec le « bootstrap target »¹ : $\mathbf{y}_k = r_{t+1}^{(i)} + \hat{V}_{\phi}^{\pi}(\mathbf{s}_{t+1}^{(i)})$

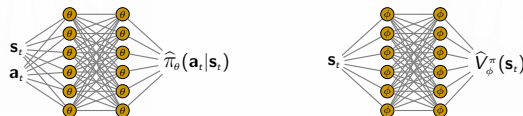
3 **Calculer l'Avantage**² : $\hat{A}^{\pi}(\mathbf{s}_t, r_t^{(i)}, \mathbf{s}_{t+1}^{(i)}) = r_{t+1}^{(i)} + \hat{V}_{\phi}^{\pi}(\mathbf{s}_{t+1}^{(i)}) - \hat{V}_{\phi}^{\pi}(\mathbf{s}_t)$

4 **Résoudre l'objectif :** $\nabla_{\theta_k} J(\theta_k) \approx \sum_{t=1}^T \nabla_{\theta_k} \log \pi_{\theta_k}(\mathbf{a}_t^{(i)} | \mathbf{s}_t) \hat{A}^{\pi}(\mathbf{s}_t, r_t^{(i)}, \mathbf{s}_{t+1}^{(i)})$

1. On peut également utiliser un « Monte-Carlo target » : $\mathbf{y}_k = \hat{Q}_t^{(i)}$

2. Notation : \mathbf{s}_t est le point de départ de l'échantillon courant et a été échantillonné à l'itération précédente.

ON-line Advantage Actor-Critic



Caractéristiques : Apprentissage **ON-policy** et **EN-ligne** (Il existe aussi une version HORS-ligne) ;

Idées clés :

1 Apprendre deux réseaux neuronaux :

$\hat{\pi}_{\theta}(\mathbf{a}|\mathbf{s}) \rightarrow$ **l'acteur** : celui responsable de prendre les **décisions** dans l'environnement ;

$\hat{V}_{\phi}^{\pi}(\mathbf{s}) \rightarrow$ **le critique** : celui responsable d'**évaluer** si $\hat{\pi}_{\theta}$ fait du bon travail ;

2 Apprendre $\hat{V}_{\phi}^{\pi}(\mathbf{s}_t)$ par MSE avec le « bootstrap target »¹ : $\mathbf{y}_k = r_{t+1}^{(i)} + \hat{V}_{\phi}^{\pi}(\mathbf{s}_{t+1}^{(i)})$

3 Calculer l'Avantage² : $\hat{A}^{\pi}(\mathbf{s}_t, r_t^{(i)}, \mathbf{s}_{t+1}^{(i)}) = r_{t+1}^{(i)} + \hat{V}_{\phi}^{\pi}(\mathbf{s}_{t+1}^{(i)}) - \hat{V}_{\phi}^{\pi}(\mathbf{s}_t)$

4 Résoudre l'objectif : $\nabla_{\theta_k} J(\theta_k) \approx \sum_{t=1}^T \nabla_{\theta_k} \log \pi_{\theta_k}(\mathbf{a}_t^{(i)} | \mathbf{s}_t) \hat{A}^{\pi}(\mathbf{s}_t, r_t^{(i)}, \mathbf{s}_{t+1}^{(i)})$

1. On peut également utiliser un « Monte-Carlo target » : $\mathbf{y}_k = \hat{Q}_t^{(i)}$

2. Notation : \mathbf{s}_t est le point de départ de l'échantillon courant et a été échantillonné à l'itération précédente.

↺ Algorithme | Advantage Actor-Critic

Soit « l'acteur » $\hat{\pi}_{\theta}(\mathbf{a}|\mathbf{s})$ (une politique différentiable), le « critique » $\hat{V}_{\phi}^{\pi}(\mathbf{s})$, une « batch » de trajectoire \mathcal{D} , le « learning rate » $\alpha \in]0, 1]$, et le « global step » $k \leftarrow 0$

Répéter jusqu'à convergence vers $\approx \pi^*$:

Répéter pour chaque « timestep » $t = 0, 1, 2, \dots, t, t+1, \dots, T$ de la trajectoire τ :

```

1    $r_{t+1}^{(k)}, \mathbf{s}_{t+1}^{(k)} \leftarrow$  exécuter  $\mathbf{a}_t \sim \mu(\cdot|\mathbf{s}_t)$  dans l'environnement et observer la réaction
2   mettre à jour  $\hat{V}_{\phi}^{\pi}(\mathbf{s}_t)$  par MSE avec le « bootstrap target »  $\mathbf{y}_k = r_{t+1}^{(k)} + \hat{V}_{\phi}^{\pi}(\mathbf{s}_{t+1}^{(k)})$ 
3    $\nabla_{\theta_k} J(\theta_k) \approx \nabla_{\theta_k} \log \pi_{\theta_k}(\mathbf{a}_t|\mathbf{s}_t) \left( r_{t+1}^{(k)} + \hat{V}_{\phi}^{\pi}(\mathbf{s}_{t+1}^{(k)}) - \hat{V}_{\phi}^{\pi}(\mathbf{s}_t) \right)$ 
4    $\theta_{k+1} \leftarrow \theta_k + \alpha \nabla_{\theta_k} J(\theta_k)$ 
5    $\mathbf{s}_t \leftarrow \mathbf{s}_{t+1}^{(k)}$  et  $k \leftarrow k + 1$ 
```

👍 **Pro** : Moins de variance que REINFORCE avec « reward-to-go » ;

👎 **Con** : N'est plus non-biaisé par rapport à REINFORCE ;

Ressource théorique

REINFORCE

- *Simple statistical gradient-following algorithms for connectionist reinforcement learning* ▶ Publication
par Ronald J. Williams (1992) [11]
- *Policy Gradient Methods for Reinforcement Learning with Function Approximation* ▶ Publication
par Richard S. Sutton, David McAllester, Satinder Singh, Yishay Mansour (2017) [12]
- *Reinforcement Learning : An introduction*, chapitre 13 [6] ▶ Publication
par Sutton & Barto
- *CS 294-112 Deep Reinforcement Learning : lecture 5 - Policy Gradient* ▶ Présentation
donné par professeur Sergey Levine à l'université Berkeley en Californie

Actor-Critic

- *Reinforcement Learning : An introduction*, chapitre 13, section 5 [6] ▶ Publication
par Sutton & Barto
- *CS 294-112 Deep Reinforcement Learning : lecture 6 - Actor-Critic Algorithms* ▶ Présentation
donné par professeur Sergey Levine à l'université Berkeley en Californie

Pour aller plus loin

1 Concepts clés en RL (la suite)

2 Apprentissage par renforcement profond

3 DRL et méthode basée sur les valeurs

4 Méthode par recherche de politique

5 DRL et méthode *Policy Gradient*

6 Pour aller plus loin

- Complément théorique
- Références

Pour aller plus loin

Complément théorique

REINFORCE [Monte-Carlo Policy Gradient]

Différentiation directe de la politique (dérivation complète)

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &= \nabla_{\theta} \left(\mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} [G(\tau)] \right) && \langle \text{L'objectif de l'apprentissage par renforcement} \rangle \\
 &= \int \nabla_{\theta} \pi_{\theta}(\tau) G(\tau) d\tau && \left\langle \text{Identité : } \nabla_{\theta} \pi_{\theta}(\tau) = \frac{\pi_{\theta}(\tau)}{\pi_{\theta}(\tau)} \nabla_{\theta} \pi_{\theta}(\tau) = \pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau) \right\rangle \\
 &= \int \pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau) G(\tau) d\tau \\
 &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) G(\tau)] && \left\langle \pi_{\theta}(\tau) = p(\mathbf{s}_1) \prod_{t=1}^T \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \right\rangle \\
 &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} \left[\nabla_{\theta} \log \left(p(\mathbf{s}_1) \prod_{t=1}^T \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \right) G(\tau) \right] && \left\langle \text{Identité : } \log(A \cdot B) = \log A + \log B \right\rangle \\
 &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} \left[\nabla_{\theta} \left(\log p(\mathbf{s}_1) + \sum_{t=1}^T \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) + \log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \right) G(\tau) \right] && \langle \text{les termes qui ne dépendent pas de } \theta \text{ s'éliminent} \rangle \\
 &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} \left[\nabla_{\theta} \left(\sum_{t=1}^T \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) G(\tau) \right] && \left\langle \text{Propriété de distributivité : } \nabla(A+B) = \nabla A + \nabla B \right\rangle \\
 &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) G(\tau) \right]
 \end{aligned}$$

REINFORCE [Monte-Carlo Policy Gradient]

Combinaison du résultat de la **différentiation directe de la politique** et celui de l'**estimé de Monte-Carlo (dérivation complète)** :

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) G(\tau) \right] \\
 &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} \left[\left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \left(\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right) \right] \\
 &\quad \langle \text{estimez le gradient par rapport à } N \text{ échantillons} \rangle \\
 &\approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \right) \left(\sum_{t=1}^T r_t^{(i)} \right) \\
 &= \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau^{(i)}) G(\tau^{(i)})
 \end{aligned}$$

REINFORCE avec *reward-to-go*

Application du **principe de causalité** à la structure temporelle (**détails du calcul**) :

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau^{(i)}) G(\tau^{(i)}) \\
 &= \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \right) \left(\sum_{t=1}^T r_t^{(i)} \right) \\
 &\quad \langle \text{Distribuer les récompenses à l'intérieur de la somme des politiques} \rangle \\
 &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \left(\sum_{t'=1}^T r_{t'}^{(i)} \right) \\
 &\quad \langle \text{Appliquer le principe de causalité en changeant l'index de départ du retour } t'=1 \rightarrow t'=t \rangle \\
 &\geq \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \left(\sum_{t'=t}^T r_{t'}^{(i)} \right) \\
 &\quad \langle \text{« reward-to-go »} \rangle \\
 &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \hat{Q}_t(\tau^{(i)})
 \end{aligned}$$

REINFORCE avec *reward-to-go*

Application du **principe de causalité** à la structure temporelle (**détails du calcul**) :

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau^{(i)}) G(\tau^{(i)}) \\
 &= \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \right) \left(\sum_{t=1}^T r_t^{(i)} \right) \\
 &\quad \langle \text{Distribuer les récompenses à l'intérieur de la somme des politiques} \rangle \\
 &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \left(\sum_{t'=1}^T r_{t'}^{(i)} \right) \\
 &\quad \langle \text{Appliquer le principe de causalité en changeant l'index de départ du retour } t'=1 \rightarrow t'=t \rangle \\
 &\geq \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \left(\sum_{t'=t}^T r_{t'}^{(i)} \right) \\
 &\quad \langle \text{« reward-to-go »} \rangle \\
 &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \hat{Q}_t(\tau^{(i)})
 \end{aligned}$$

REINFORCE avec *reward-to-go*

Application du **principe de causalité** à la structure temporelle (**détails du calcul**) :

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau^{(i)}) G(\tau^{(i)}) \\
 &= \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \right) \left(\sum_{t=1}^T r_t^{(i)} \right) \\
 &\quad \langle \text{Distribuer les récompenses à l'intérieur de la somme des politiques} \rangle \\
 &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \left(\sum_{t'=1}^T r_{t'}^{(i)} \right) \\
 &\quad \langle \text{Appliquer le principe de causalité en changeant l'index de départ du retour } t'=1 \rightarrow t'=t \rangle \\
 &\geq \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \left(\sum_{t'=t}^T r_{t'}^{(i)} \right) \\
 &\quad \langle \text{« reward-to-go »} \rangle \\
 &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \hat{Q}_t(\tau^{(i)})
 \end{aligned}$$

REINFORCE avec *reward-to-go*

Application du **principe de causalité** à la structure temporelle (**détails du calcul**) :

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau^{(i)}) G(\tau^{(i)}) \\
 &= \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \right) \left(\sum_{t=1}^T r_t^{(i)} \right) \\
 &\quad \langle \text{Distribuer les récompenses à l'intérieur de la somme des politiques} \rangle \\
 &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \left(\sum_{t'=1}^T r_{t'}^{(i)} \right) \\
 &\quad \langle \text{Appliquer le principe de causalité en changeant l'index de départ du retour } t'=1 \rightarrow t'=t \rangle \\
 &\geq \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \left(\sum_{t'=t}^T r_{t'}^{(i)} \right) \\
 &\quad \langle \text{« reward-to-go »} \rangle \\
 &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)}) \hat{Q}_t(\tau^{(i)})
 \end{aligned}$$

Pour aller plus loin

Références

References I

1. Van HASSELT, H., GUEZ, A. & SILVER, D. Deep Reinforcement Learning with Double Q-learning. *Artificial Intelligence* **230**, 173-191. arXiv : 1509.06461. <http://arxiv.org/abs/1509.06461> (sept. 2015).
2. WANG, Z. *et al.* Dueling Network Architectures for Deep Reinforcement Learning. arXiv : 1511.06581. <http://arxiv.org/abs/1511.06581> (2015).
3. MUNOS, R., STEPLETON, T., HARUTYUNYAN, A. & BELLEMARE, M. G. *Safe and efficient off-policy reinforcement learning*. in *Advances in Neural Information Processing Systems* (2016). arXiv : 1606.02647.
4. SCHAU, T., QUAN, J., ANTONOGLOU, I. & SILVER, D. *Prioritized experience replay*. in *4th International Conference on Learning Representations, ICLR 2016 - Conference Track Proceedings* (2016). arXiv : 1511.05952.
5. LILICRAP, T. P. *et al.* Continuous control with deep reinforcement learning. arXiv : 1509.02971. <http://arxiv.org/abs/1509.02971> (2015).
6. SUTTON, R. S. & BARTO, A. G. *Reinforcement learning: An introduction*. 2^e éd. (éd. MIT PRESS) ISBN : 978-0262039246. <http://incompleteideas.net/book/RLbook2018.pdf> (Cambridge, MA, 2018).

References II

7. LIN, L. J. Programming Robots Using Reinforcement Learning and Teaching. *Aaai-91 the Ninth National Conference on Artificial Intelligence*, 781-786.
<http://www.aaai.org/Library/AAAI/1991/aaai91-122.php> (1991).
8. MNIH, V. *et al.* Playing Atari with Deep Reinforcement Learning. 1-9. arXiv : 1312.5602.
<http://arxiv.org/abs/1312.5602> (2013).
9. MNIH, V. *et al.* Human-level control through deep reinforcement learning. *Nature* **518**, 529-533. ISSN : 14764687. <https://web.stanford.edu/class/psych209/Readings/MnihEtAlHassibis15NatureControlDeepRL.pdf> (2015).
10. SCIENCE, C. John Schulman Thesis: Optimizing Expectations: From Deep Reinforcement Learning to Stochastic Computation Graphs. *Ph.D thesis* (2016).
11. WILLIAMS, R. J. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning* **8**, 229-256. ISSN : 0885-6125 (1992).
12. SUTTON, R. S., MCALLESTER, D., SINGH, S. & MANSOUR, Y. Policy Gradient Methods for Reinforcement Learning with Function Approximation Richard. *Advances in Neural Information Processing Systems* **12**, 1057-1063. arXiv : 1706.06643.
<http://arxiv.org/abs/1706.06643> (2017).