

Value Function Approximation and Deep Q Learning

Clément Rambour

Plan

1 Introduction

2 Approximation de la fonction de valeur

- Stochastic gradient descent
- Monte Carlo with value function approximation
- TD Learning with value function approximation
- Control with value function approximation

3 Deep Reinforcement Learning

- End to end learning of Q
- Double DQN
- Dueling DQN

4 Bibliographie

Cours précédent

- Comment apprendre une bonne politique grâce à l'expérience
- Pour l'instant : représentation tabulaire de la fonction de valeur d'état ou d'état-action
- Souvent trop d'états pour être applicable

Aujourd'hui

Motivation :

- Ne pas stocker ou apprendre explicitement pour chaque état
 - Une dynamique
 - Une récompense
 - Valeur d'action-état
 - Politique
- Recherche d'une représentation compacte qui se généralise mieux

Aujourd'hui

Motivation :

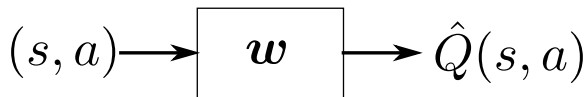
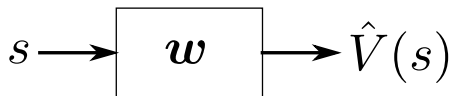
- Ne pas stocker ou apprendre explicitement pour chaque état
 - Une dynamique
 - Une récompense
 - Valeur d'action-état
 - Politique
- Recherche d'une représentation compacte qui se généralise mieux
- Réduction de l'impact mémoire
- Réduction du coût en calcul
- Réduit la quantité d'expérience nécessaire

Plan

- 1 Introduction
- 2 Approximation de la fonction de valeur
 - Stochastic gradient descent
 - Monte Carlo with value function approximation
 - TD Learning with value function approximation
 - Control with value function approximation
- 3 Deep Reinforcement Learning
 - End to end learning of Q
 - Double DQN
 - Dueling DQN
- 4 Bibliographie

Principe

Inférence de la valeur d'un état/ d'un couple action-état grâce à un modèle



Fonctions approximatrices

De nombreuses fonctions n'ont pas de calcul possible pour calculer la fonction de valeur :

- Combinaison linéaire de features
- Réseaux de neurones
- Arbres de décisions
- Plus proches voisins
- Fourier / ondelettes

Fonctions approximatrices

De nombreuses fonctions n'ont pas de calcul possible pour calculer la fonction de valeur :

- Combinaison linéaire de features ← différentiable
- Réseaux de neurones ← différentiable
- Arbres de décisions
- Plus proches voisins
- Fourier / ondelettes

Table of Contents

1 Introduction

2 Approximation de la fonction de valeur

- Stochastic gradient descent
- Monte Carlo with value function approximation
- TD Learning with value function approximation
- Control with value function approximation

3 Deep Reinforcement Learning

- End to end learning of Q
- Double DQN
- Dueling DQN

4 Bibliographie

Apprentissage avec un oracle

Intuition : se ramener à un cas supervisé

- Si on connaît la valeur de $v_\pi(s) \forall s$
- Apprentissage supervisé $(s, v_\pi(s))$
- L'objectif est d'obtenir la meilleure approximation de v_π grâce à la fonction paramétrée $\hat{v}(s, \mathbf{w})$

Descente de gradient stochastique

- **But** : trouver les paramètres \mathbf{w} qui minimisent la loss $\mathcal{L}(v_\pi(s), \hat{v}(s, \mathbf{w}))$
- Généralement, $\mathcal{L} = \text{MSE}$:

$$\mathcal{L}(v_\pi(s), \hat{v}(s, \mathbf{w})) = \mathbb{E}_\pi[(v_\pi(s) - \hat{v}(s, \mathbf{w}))^2]$$

- Minimum trouvé par descente de gradient :

$$\Delta \mathbf{w} = -\frac{1}{2} \alpha \nabla_{\mathbf{w}} \mathcal{L}$$

- Descente de gradient stochastique : gradient approché par une moyenne sur un nombre finit d'échantillons (souvent un seul)

Model Free Approximation

- Ne nécessite pas forcément d'être calculé avec un oracle/les labels
- Model free policy evaluation :
 - On suit une politique π
 - Convergence vers v_π ou q_π
 - Valeurs d'état / d'actions-états stockées en mémoire
 - Mise à jour après un épisode (MC) ou après chaque étape (TD)
- **Modification** de l'approche existente pour inclure la mise à jour des paramètres du modèle

Table of Contents

1 Introduction

2 Approximation de la fonction de valeur

- Stochastic gradient descent
- Monte Carlo with value function approximation
- TD Learning with value function approximation
- Control with value function approximation

3 Deep Reinforcement Learning

- End to end learning of Q
- Double DQN
- Dueling DQN

4 Bibliographie

Approximation linéaire

- Chaque état s est représenté par un ensemble de descripteurs/*features* :
 $\mathbf{x}(s) = [x_1(s), \dots, x_n(s)]$
- L'approximation de la fonction de valeur est obtenue par approximation linéaire : $\hat{v}(s, \mathbf{w}) = \mathbf{w}^T \mathbf{x}(s)$
- La loss est donnée par $\mathcal{L}(\mathbf{w}) = \mathbb{E}_\pi [(v_\pi(s) - \hat{v}(s, \mathbf{w}))^2]$
- Gradient donné par :

$$\nabla_{\mathbf{w}} \mathcal{L} = \mathbb{E}_\pi [2(v_\pi(s) - \hat{v}(s, \mathbf{w}))] \nabla_{\mathbf{w}} \hat{v}(s, \mathbf{w})$$

Monte Carlo Value function approximation

- Monte Carlo : Le retour G_t est un estimateur non biaisé de $v_\pi(s)$
- Apprentissage supervisé : $(s_1, G_1), \dots, (s_T, G_T)$
- Le gradient s'écrit :

$$\nabla_{\mathbf{w}} \mathcal{L} = -2(G_t - \hat{v}(s, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(s, \mathbf{w})$$

- La mise à jour :

$$\begin{aligned} \Delta \mathbf{w} &= -\frac{1}{2} \alpha \nabla_{\mathbf{w}} \mathcal{L} \\ &= \alpha (G_t - \hat{v}(s, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(s, \mathbf{w}) \\ &= \alpha (G_t - \mathbf{x}(s_t)^T \mathbf{w}) \mathbf{x}(s_t) \end{aligned}$$

Monte Carlo Value function approximation

```
1: Initialize  $w = 0$ ,  $k = 1$ 
2: loop
3:   Sample  $k$ -th episode  $(s_{k,1}, a_{k,1}, r_{k,1}, s_{k,2}, \dots, s_{k,L_k})$  given  $\pi$ 
4:   for  $t = 1, \dots, L_k$  do
5:     if First visit to  $(s)$  in episode  $k$  then
6:        $G_t(s) = \sum_{j=t}^{L_k} r_{k,j}$ 
7:       Update weights:

8:     end if
9:   end for
10:   $k = k + 1$ 
11: end loop
```

Table of Contents

1 Introduction

2 Approximation de la fonction de valeur

- Stochastic gradient descent
- Monte Carlo with value function approximation
- TD Learning with value function approximation
- Control with value function approximation

3 Deep Reinforcement Learning

- End to end learning of Q
- Double DQN
- Dueling DQN

4 Bibliographie

TD Learning

- Bootstrapping
- Mise à jour après chaque transition (s, a, r, s') :

$$V(s) = V(s) + \alpha \underbrace{(r + \gamma V(s') - V(s))}_{\text{TD Target}}$$

- TD Target : estimation biaisée de $v_\pi(s)$

TD Learning

- Bootstrapping
- Mise à jour après chaque transition (s, a, r, s') :

$$V(s) = V(s) + \alpha \underbrace{(r + \gamma V(s')) - V(s)}_{\text{TD Target}}$$

- TD Target : estimation biaisée de $v_\pi(s)$
- \Rightarrow Utilisation de la TD Target **approchée** comme supervision :
 $(s_1, r_1 + \gamma \hat{v}(s_2, \mathbf{w})), \dots, (s_T, r_{T+1} + \gamma \hat{v}(s_T, \mathbf{w}))$
- **But** : Trouver les poids qui minimisent la loss :

$$\mathcal{L}(\mathbf{w}) = \mathbb{E}_\pi [(r_{t+1} + \gamma \hat{v}(s_t, \mathbf{w}) - \hat{v}(s_t, \mathbf{w}))^2]$$

- Pour TD(0), mise à jour cas linéaire :

$$\begin{aligned} \Delta \mathbf{w} &= \alpha (r + \gamma \hat{v}(s', \mathbf{w}) - \hat{v}(s, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(s, \mathbf{w}) \\ &= \alpha (r + \gamma \mathbf{x}(s')^T \mathbf{w} - \mathbf{x}(s)^T \mathbf{w}) \mathbf{x}(s) \end{aligned} \quad (1)$$

TD(0) Linear value function approximation

-
-
- 1: Initialize $\mathbf{w} = 0$, $k = 1$
 - 2: **loop**
 - 3: Sample tuple (s_k, a_k, r_k, s_{k+1}) given π
 - 4: Update weights:

$$\mathbf{w} = \mathbf{w} + \alpha(r + \gamma \mathbf{x}(s')^T \mathbf{w} - \mathbf{x}(s)^T \mathbf{w}) \mathbf{x}(s)$$

- 5: $k = k + 1$
 - 6: **end loop**
-

Table of Contents

1 Introduction

2 Approximation de la fonction de valeur

- Stochastic gradient descent
- Monte Carlo with value function approximation
- TD Learning with value function approximation
- Control with value function approximation

3 Deep Reinforcement Learning

- End to end learning of Q
- Double DQN
- Dueling DQN

4 Bibliographie

Control et approximation de la valeur

- Approximation de la fonction de valeur état-action :

$$Q_{\pi}(s, a) \simeq \hat{Q}(s, a; \mathbf{w})$$

- $\hat{Q}(s, a; \mathbf{w})$ paramétrée par \mathbf{w}
- Schéma classique possible (légèrement modifié) :
 - Approximation de la fonction de valeur
 - ϵ -greedy amélioration

Control et approximation de la valeur

- Approximation de la fonction de valeur état-action :

$$Q_{\pi}(s, a) \simeq \hat{Q}(s, a; \mathbf{w})$$

- $\hat{Q}(s, a; \mathbf{w})$ paramétrée par \mathbf{w}
- Schéma classique possible (légèrement modifié) :
 - Approximation de la fonction de valeur
 - ϵ -greedy amélioration
- Parfois instable
- Amélioration possible :
 - Approximation bien choisie
 - Bootstrapping
 - Off-policy learning

Contrôle avec un oracle

- $Q_\pi(s, a) \simeq \hat{Q}(s, a; \mathbf{w})$ Mean Square Error :

$$\mathcal{L}(\mathbf{w}) = \mathbb{E}_\pi[(Q_\pi(s, a) - \hat{Q}(s, a; \mathbf{w}))^2]$$

- Minimum obtenu par stochastic gradient descent :

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \frac{1}{2} \mathbb{E}_\pi \left[(Q_\pi(s, a) - \hat{Q}(s, a; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w}) \right]$$

SARSA et approximation de la valeur

- Résumé en deux étapes :

- À partir de s choix d'une action a associée à une politique ϵ -greedy(π) $\rightarrow s'$
- Mise à jour de la fonction de valeur en fonction du couple (s', a') associé à la prochaine action :

$$Q_{\pi}(s, a) \leftarrow Q_{\pi}(s, a) + \alpha(r + \gamma Q_{\pi}(s', a') - Q(s, a))$$

- SARSA avec approximation $Q_{\pi}(s, a) \simeq \hat{Q}(s, a; \mathbf{w})$

- Loss avec objectif SARSA :

$$\mathcal{L}(\mathbf{w}) = \mathbb{E}_{\pi} [(r + \gamma Q_{\pi}(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w}))^2]$$

$$\mathcal{L}(\mathbf{w}) \simeq (r + \gamma Q_{\pi}(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w}))^2$$

- gradient de la loss % \mathbf{w} :

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = 2(r + \gamma \hat{Q}(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w})$$

- Mise à jour :

$$\mathbf{w} \leftarrow \mathbf{w} - \frac{1}{2} \alpha \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w})$$

SARSA et approximation de la valeur : cas linéaire

- Couple états actions représentés par un vecteur de features :
 $\mathbf{x}(s, a) = [x_1(s, a), \dots, x_n(s, a)]$

- Approximation linéaire :

$$\hat{Q}(s, a; \mathbf{w}) = \mathbf{w}^t \mathbf{x}(s, a)$$

- Loss avec objectif SARSA :

$$\mathcal{L}(\mathbf{w}) \simeq (r + \gamma Q_\pi(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w}))^2$$

- gradient de la loss % \mathbf{w} cas linéaire :

$$\begin{aligned} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) &= 2(r + \gamma \hat{Q}(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w}) \\ &= (r + \gamma \mathbf{x}(s', a')^T \mathbf{w} - \mathbf{x}(s, a)^T \mathbf{w}) \mathbf{x}(s, a) \end{aligned}$$

- Mise à jour :

$$\mathbf{w} \leftarrow \mathbf{w} - \frac{1}{2} \alpha \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w})$$

Q-learning et approximation de la valeur

- Résumé en deux étapes :

- À partir de s choix d'une action a associée à une politique ϵ -greedy(π) $\leftarrow s'$
- Mise à jour de la fonction de valeur en fonction du couple (s', a') associé à la **meilleure** action possible :

$$Q_{\pi}(s, a) \leftarrow Q_{\pi}(s, a) + \alpha(r + \gamma \max_{a'} Q_{\pi}(s', a') - Q(s, a))$$

- Q-learning avec approximation $Q_{\pi}(s, a) \simeq \hat{Q}(s, a; \mathbf{w})$
- Loss avec objectif Q-learning :

$$\mathcal{L}(\mathbf{w}) = \mathbb{E}_{\pi} [(r + \gamma \max_{a'} Q_{\pi}(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w}))^2]$$

$$\mathcal{L}(\mathbf{w}) \simeq (r + \max_{a'} \gamma Q_{\pi}(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w}))^2$$

- gradient de la loss $\% \mathbf{w}$:

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = 2(r + \gamma \max_{a'} \hat{Q}(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w})$$

- Mise à jour :

$$\mathbf{w} \leftarrow \mathbf{w} - \frac{1}{2} \alpha \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w})$$

Q-learning et approximation de la valeur : cas linéaire

- Couple états actions représentés par un vecteur de features :

$$\mathbf{x}(s, a) = [x_1(s, a), \dots, x_n(s, a)]$$

- Approximation linéaire :

$$\hat{Q}(s, a; \mathbf{w}) = \mathbf{w}^t \mathbf{x}(s, a)$$

- Loss avec objectif Q-learning :

$$\mathcal{L}(\mathbf{w}) \simeq (r + \max_{a'} \gamma Q_{\pi}(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w}))^2$$

- gradient de la loss % \mathbf{w} cas linéaire :

$$\begin{aligned} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) &= 2(r + \gamma \max_{a'} \hat{Q}(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w}) \\ &= (r + \gamma \max_{a'} \mathbf{x}(s', a')^T \mathbf{w} - \mathbf{x}(s, a)^T \mathbf{w}) \mathbf{x}(s, a) \end{aligned}$$

- Mise à jour :

$$\mathbf{w} \leftarrow \mathbf{w} - \frac{1}{2} \alpha \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w})$$

Plan

- 1 Introduction
- 2 Approximation de la fonction de valeur
 - Stochastic gradient descent
 - Monte Carlo with value function approximation
 - TD Learning with value function approximation
 - Control with value function approximation
- 3 Deep Reinforcement Learning
 - End to end learning of Q
 - Double DQN
 - Dueling DQN
- 4 Bibliographie

Deep RL

- Utiliser des réseaux de neurones pour approcher
 - La fonction de valeur (v ou Q)
 - La politique
 - Le modèle (de l'environnement)
- Optimisation par descente de gradient stochastique

Table of Contents

1 Introduction

2 Approximation de la fonction de valeur

- Stochastic gradient descent
- Monte Carlo with value function approximation
- TD Learning with value function approximation
- Control with value function approximation

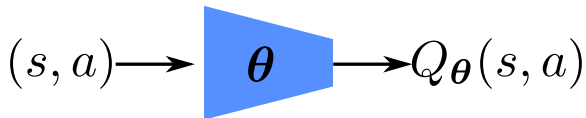
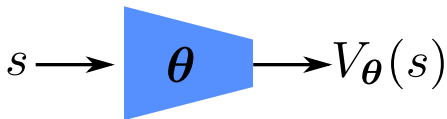
3 Deep Reinforcement Learning

- End to end learning of Q
- Double DQN
- Dueling DQN

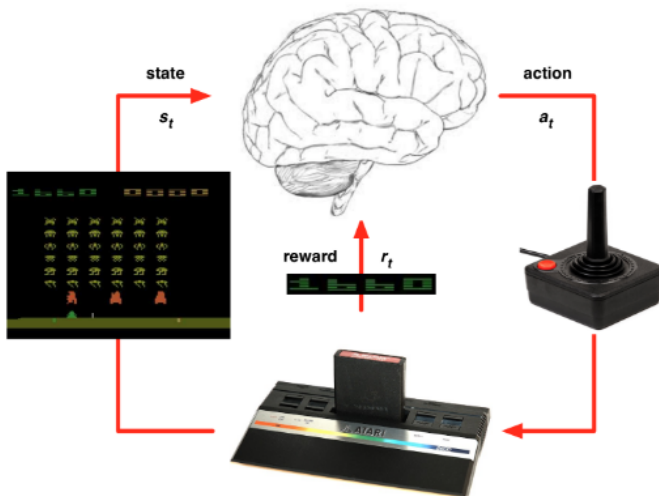
4 Bibliographie

Deep Q-Networks (DQNs)

Représentation de la table de valeur état-action par un Q -network de paramètre θ

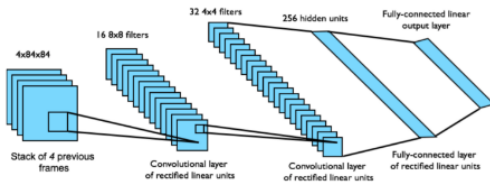


DQNs dans Atari



DQNs dans Atari

- Apprentissage *end to end* des valeurs de $Q(s, a)$ d'après les pixels $\mathbf{s} \leftarrow$ vecteur de features décrivant l'état
- Entrée du réseau \mathbf{s} concaténation des dernières frames observées (ex : 4 dernières)
- Sortie : $Q(s, a)$ pour les actions réalisables (~ 18 boutons)
- Architectures et hyperparamètres fixés pour toute la partie

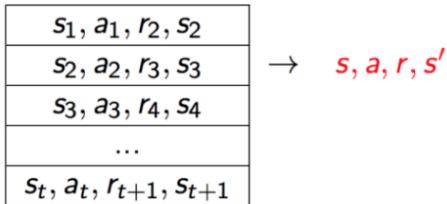


DQ-learning

- Loss : MSE optimisée par descente de gradient stochastique
- Deux problèmes (principaux) pouvant amener à diverger :
 - Correlation entre les échantillons
 - *Targets*/labels non stationnaires
- Deux solutions proposées pour le Deep Q-Learning (Mnih et Al. 2015)
 - Experience replay
 - Fixed Q-targets

DQNs : Experience replay

- Pour éviter la forte corrélation entre échantillons successifs : les données sont stockées dans un *buffer* $\mathcal{D} \leftarrow \text{replay buffer}$



- *Experience replay* consiste à répéter les étapes suivantes :

- Échantillonnage d'une expérience $(s, a, r, s' \sim \mathcal{D})$
- Calcul de la TD *target* pour l'échantillon : $r + \gamma \max_{a'} Q_{\theta}(s', a')$
- Calcul de la loss :

$$\mathcal{L}(\theta) = (r + \gamma \max_{a'} Q_{\theta}(s', a') - Q_{\theta}(s, a))^2$$

- Mise à jour des paramètres :

$$\theta \leftarrow \theta - \alpha \frac{1}{2} \nabla_{\theta} \mathcal{L}(\theta)$$

$$\theta \leftarrow \theta + \alpha (r + \gamma \max_{a'} Q_{\theta}(s', a') - Q_{\theta}(s, a)) \nabla_{\theta} Q_{\theta}(s, a)$$

Fixed Q-Targets

- Pour améliorer la stabilité, on peut fixer les poids du réseaux pour évaluer la TD *target*
- Deux réseaux : un réseaux cible et un réseaux apprenant la fonction de valeur
- θ^- : les paramètres du réseau cible
- θ : les paramètres du réseaux mis à jour
- Léger changement dans le schéma précédent :
 - Échantillonnage d'une expérience $(s, a, r, s') \sim \mathcal{D}$
 - Calcul de la TD *target* pour l'échantillon : $r + \gamma \max_{a'} Q_{\theta^-}(s', a')$
 - Calcul de la loss :

$$\mathcal{L}(\theta) = (r + \gamma \max_{a'} Q_{\theta^-}(s', a') - Q_{\theta}(s, a))^2$$

- Mise à jour des paramètres :

$$\theta \leftarrow \theta - \alpha \frac{1}{2} \nabla_{\theta} \mathcal{L}(\theta)$$

$$\theta \leftarrow \theta + \alpha (r + \gamma \max_{a'} Q_{\theta^-}(s', a') - Q_{\theta}(s, a)) \nabla_{\theta} Q_{\theta}(s, a)$$

DQN Pseudocode

```

1: Input  $C, \alpha, D = \{\}$ , Initialize  $\mathbf{w}, \mathbf{w}^- = \mathbf{w}, t = 0$ 
2: Get initial state  $s_0$ 
3: loop
4:   Sample action  $a_t$  given  $\epsilon$ -greedy policy for current  $\hat{Q}(s_t, a; \mathbf{w})$ 
5:   Observe reward  $r_t$  and next state  $s_{t+1}$ 
6:   Store transition  $(s_t, a_t, r_t, s_{t+1})$  in replay buffer  $D$ 
7:   Sample random minibatch of tuples  $(s_j, a_j, r_j, s_{j+1})$  from  $D$ 
8:   for  $j$  in minibatch do
9:     if episode terminated at step  $i + 1$  then
10:       $y_j = r_j$ 
11:     else
12:       $y_j = r_j + \gamma \max_{a'} \hat{Q}(s_{j+1}, a'; \mathbf{w}^-)$ 
13:     end if
14:     Do gradient descent step on  $(y_j - \hat{Q}(s_j, a_j; \mathbf{w}))^2$  for parameters  $\mathbf{w}$ :  $\Delta \mathbf{w} = \alpha(y_j - \hat{Q}(s_j, a_j; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s_j, a_j; \mathbf{w})$ 
15:   end for
16:    $t = t + 1$ 
17:   if  $\text{mod}(t, C) == 0$  then
18:      $\mathbf{w}^- \leftarrow \mathbf{w}$ 
19:   end if
20: end loop

```

DQN results dans Atari

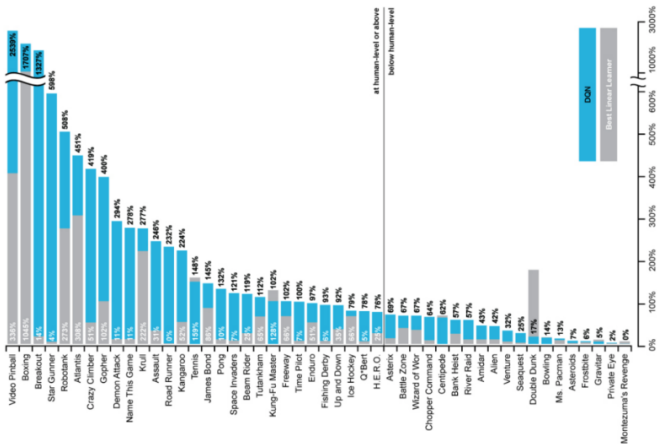


FIGURE – Human-level control through deep reinforcement learning, Mnih et al, 2015

Vanilla DQN vs. Fixed Target and experience replay

Game	Linear	Deep Network	DQN w/ fixed Q	DQN w/ replay	DQN w/replay and fixed Q
Breakout	3	3	10	241	317
Enduro	62	29	141	831	1006
River Raid	2345	1453	2868	4102	7447
Seaquest	656	275	1003	823	2894
Space Invaders	301	302	373	826	1089

Table of Contents

1 Introduction

2 Approximation de la fonction de valeur

- Stochastic gradient descent
- Monte Carlo with value function approximation
- TD Learning with value function approximation
- Control with value function approximation

3 Deep Reinforcement Learning

- End to end learning of Q
- Double DQN
- Dueling DQN

4 Bibliographie

Double DQN

- Q-learning : biaisé vers l'estimation du maximum dans la table Q
- Combinaison d'estimateur peut réduire le biais (\sim bagging)
- Double Q-Learning :
 - Combinaison de deux estimations de Q
 - Probabilité non nulle de sélectionner la meilleure option pour l'une ou pour l'autre des estimations
 - réduction du biais

Prioritized experience replay

- Soit i l'index du i -eme tuple dans le buffer d'expérience (s_i, a_i, r_i, s_{i+1})
- Échantillonnage des tuples pour mettre à jour suivant une distribution p
- Probabilité p_i de sélectionner le tuple i proportionnelle à l'erreur DQN e_i :

$$e_i = |r + \max_{a'} Q_{\theta}(s_{i+1}, a') - Q_{\theta}(s_i, a_i)|$$

- Mise à jour des p_i à chaque mise à jour de Q_{θ}
- une méthode :

$$p_i = \frac{p_i^{\beta}}{\sum_k p_k^{\beta}}$$

- β facteur de température

Table of Contents

1 Introduction

2 Approximation de la fonction de valeur

- Stochastic gradient descent
- Monte Carlo with value function approximation
- TD Learning with value function approximation
- Control with value function approximation

3 Deep Reinforcement Learning

- End to end learning of Q
- Double DQN
- Dueling DQN

4 Bibliographie

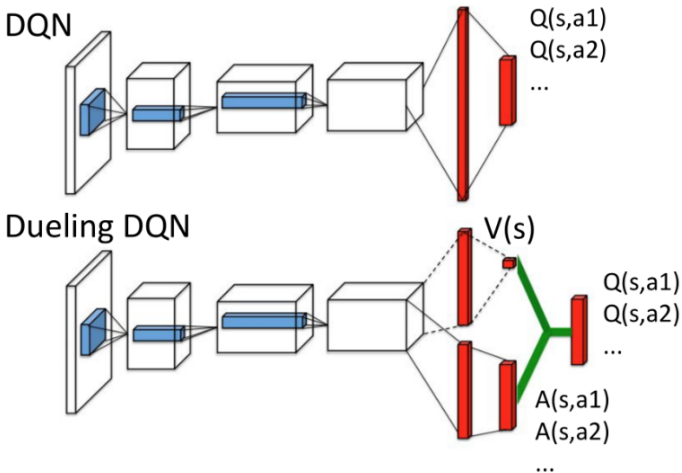
Fonctions valeur et avantage

- **Intuition** : représentations pour décrire les états pas forcément adaptées pour décrire les actions
- Exemple : le score dans un jeu indique l'intérêt d'un état s : $v(s)$ mais pas l'intérêt entre deux action $Q(s, a_1) \lesseqgtr Q(s, a_2)$
- Fonction d'avantage / *Advantage function* :

$$A_\pi(s, a) = Q_\pi(s, a) - v_\pi(s, a)$$

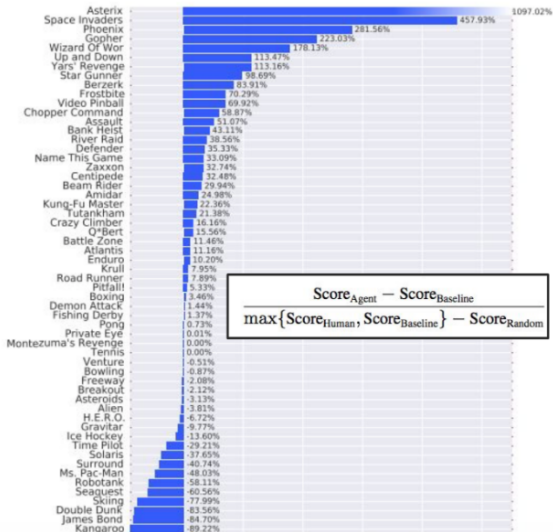
- A_π donne la valeur d'une action

Dueling DQN



Wang et.al., ICML, 2016

Dueling DQN



Non unicité de l'avantage

- Q est reconstruit à partir de l'avantage et la valeur d'état :

$$Q_{\theta}(s, a) = V_{\theta}(s, a) + A_{\theta}(s, a)$$

- Infinité de V et A donnant le même Q : solution à une constante prêt
- Nécessité de "punaiser" V et A en rajoutant une constante

Non unicité de l'avantage

- Q est reconstruit à partir de l'avantage et la valeur d'état :

$$Q_{\theta}(s, a) = V_{\theta}(s, a) + A_{\theta}(s, a)$$

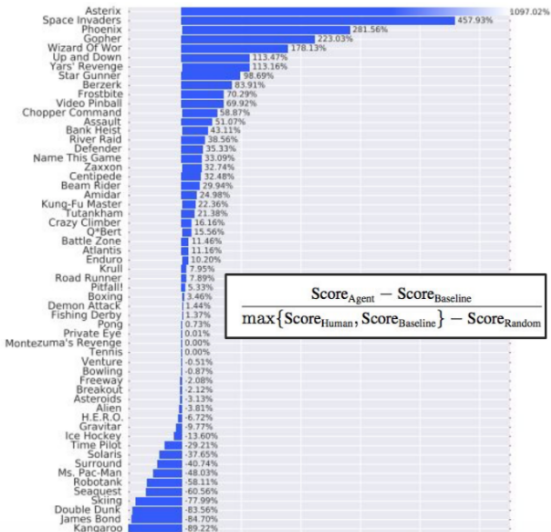
- Infinité de V et A donnant le même Q : solution à une constante prêt
- Nécessité de "punaiser" V et A en rajoutant une constante
- Méthode 1 : Forcer $Q_{\theta}(s, a) = V_{\theta}(s)$ pour la meilleure action possible

$$Q_{\theta}(s, a) = V_{\theta}(s, a) + \left(A_{\theta}(s, a) - \max_{a'} A_{\theta}(s, a') \right)$$

- Méthode 2 : Utiliser la moyenne sur les actions $\bar{A}_{\theta}(s, a')$ comme origine (plus stable)

$$Q_{\theta}(s, a) = V_{\theta}(s, a) + \left(A_{\theta}(s, a) - \bar{A}_{\theta}(s, a') \right)$$

Dueling DQN vs. Double DQN



Conseils pratiques

- Huber Loss sur l'erreur de Bellman (= l'erreur TD)
- Double DQN (peu d'effort par rapport à simple DQN)
- Essayer plusieurs seeds pendant les expériences
- Learning rate scheduling. Haut learning rate au départ pour l'exploration

Plan

- 1 Introduction
- 2 Approximation de la fonction de valeur
 - Stochastic gradient descent
 - Monte Carlo with value function approximation
 - TD Learning with value function approximation
 - Control with value function approximation
- 3 Deep Reinforcement Learning
 - End to end learning of Q
 - Double DQN
 - Dueling DQN
- 4 Bibliographie

Bibliographie

- Double DQN (Deep Reinforcement Learning with Double Q-Learning, Van Hasselt et al, AAAI 2016)
- Prioritized Replay (Prioritized Experience Replay, Schaul et al, ICLR 2016)
- Dueling DQN (best paper ICML 2016) (Dueling Network Architectures for Deep Reinforcement Learning, Wang et al, ICML 2016)