

Chapitre 5 . Q-Learning

- Q-Learning = principal algorithme de RL
- permet d'apprendre la politique optimale d'une MDP sans connaître ni π^* , ni r .
- L'apprentissage est réalisé par interaction avec l'environnement.
~~**~~ *

C'est ce qu'on cherche à apprendre :
les faire (état, action) de la politique optimale.

Deux ingrédients :

- Équation de Bellman pour la fonction q^*
- Différence Temporelle

Eq. d'optimalité de Bellman
pour \hat{q} :

$$\hat{q}^*(s, a) = \sum_{s' \in S} P(s'|s, a)(r(s'|s, a) + \gamma \max_{a' \in A} \hat{q}^*(s', a'))$$

Comme pour l'algo. TD, on considère une trajectoire (indice far t) et avec une estimation \hat{q} de q^* , le terme $\gamma + \max_{a'} \hat{q}(s_{t+1}, a') - \hat{q}(s_t, a_t)$ est une différence temporelle que l'on peut utiliser comme correction pour $\hat{q}(s_t, a_t)$.

L'algorithme peut alors être défini par la formule:

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) + \alpha_t(s_t, a_t) \left[r_{t+1} + \gamma \max_a q_t(s_{t+1}, a) - q_t(s_t, a_t) \right]$$

Algorithme (α -Learning):

Input: S, A, γ

$$\hat{q}(s, a) \leftarrow 0 \quad \forall s, a \in S \cup A$$

$$n(s, a) \leftarrow 0 \quad \forall s, a \in S \cup A$$

Repéter

initialiser l'état initial s_0

$$t \leftarrow 0$$

Tant que Episode non terminé Faire

selectionner l'action a_t et l'enregistrer

observer r_t et s_{t+1}

$$\hat{q}(s_t, a_t) \leftarrow \hat{q}(s_t, a_t)$$

$$+ \alpha (r_t, a_t, n(s_t, a_t))$$

$$\left[\hat{q}(s_{t+1}, \max_{a'} \hat{q}(s_{t+1}, a')) - \hat{q}(s_t, a_t) \right]$$

$$n(s_t, a_t) \leftarrow n(s_t, a_t) + 1$$

$$t \leftarrow t + 1$$

Fin Tant que
jusqu'à critère d'arrêt vérifié

Un exemple : Retour au Grid World.

On reprend la classe Game.

→ mais nous ferons cette fois-ci de l'apprentissage par la récompense.

SARSA: du autre algorithme d'apprentissage mais on-policy
 (l'action utilisée dans la mise à jour de \hat{q} est celle qui est exécutée)

$$\begin{aligned} q(s_t, a_t) &= q(s_t, a_t) \\ &\quad + \gamma [R_{t+1} \\ &\quad + \gamma q(s_{t+1}, a_{t+1}) \\ &\quad - q(s_t, a_t)] \end{aligned}$$

SARSA

