

CORRECTION DES TRAVAUX DIRIGÉS D'
APPRENTISSAGE PAR RENFORCEMENT
UNIVERSITÉ PARIS–SACLAY

Joon Kwon

lundi 25 novembre 2024



EXERCICE 1 (*Bornes d'erreur et critère d'arrêt*). —

1) On écrit

$$v - v_\pi = v - B_\pi v_\pi = v - B_\pi v + B_\pi v - B_\pi v_\pi.$$

Puis, en prenant la norme ℓ^∞ , et en utilisant l'inégalité triangulaire, et le fait que v_π est l'unique point fixe de l'opérateur B_π qui est une γ -contraction, on obtient

$$\|v - v_\pi\|_\infty \leq \|v - B_\pi v\|_\infty + \gamma \|v - v_\pi\|_\infty,$$

d'où le résultat. On procède de même pour les fonctions état-valeur, et pour les fonctions valeur optimales.

2) Pour une itération valeur $v_{k+1} = B_\pi v_k$, on peut vérifier le critère d'arrêt

$$\|v_{k+1} - B_\pi v_k\|_\infty \leq \varepsilon(1 - \gamma),$$

ce qui implique, d'après ce qui précède, que $\|v_k - v_\pi\|_\infty \leq \varepsilon$.

3) On considère une itération action-valeur pour le contrôle, autrement dit $q_{k+1} = B_* q_k$. On pose

$$\varepsilon = \frac{1}{3} \cdot \min_{s \in \mathcal{S}} \min_{\substack{a, a' \in \mathcal{A} \\ q_*(s, a) \neq q_*(s, a')}} |q_*(s, a) - q_*(s, a')|.$$

Puisque $q_k \rightarrow q_*$, il existe un rang k_0 à partir duquel $\|q_k - q_*\|_\infty \leq \varepsilon$. Soit $k \geq k_0$ et $\pi \in \Pi_g [q_k]$. Montrons que $\pi \in \Pi_g [q_*]$, cela montrera qu'il s'agit d'une politique optimale. Soit $s \in \mathcal{S}$. Montrons que $\pi(s) \in \text{Arg max}_{a \in \mathcal{A}} q_*(s, a)$, autrement dit que pour tout $a \in \mathcal{A}$,

$$q_*(s, \pi(s)) \geq q_*(s, a).$$

Soit $a \in \mathcal{A}$. On a

$$q_*(s, \pi(s)) \geq q_k(s, \pi(s)) - \varepsilon \geq q_k(s, a) - \varepsilon \geq q_*(s, a) - 2\varepsilon,$$

où on a utilisé la définition de π pour la deuxième inégalité. Si $q_*(s, \pi(s)) \leq q_*(s, a)$, on a

$$0 \leq q_*(s, a) - q_*(s, \pi(s)) \leq 2\varepsilon < \min_{\substack{a', a'' \in \mathcal{A} \\ q_*(s, a') \neq q_*(s, a'')}} |q_*(s, a') - q_*(s, a'')|,$$

et donc nécessairement $q_*(s, \pi(s)) = q_*(s, a)$. Dans tous les cas, on a

$$q_*(s, \pi(s)) \geq q_*(s, a),$$

et finalement que $\pi \in \Pi_g [q_*]$. On a bien obtenu une politique optimale au bout d'un nombre fini d'itérations.

EXERCICE 2 (Convergence en temps fini de l'itération de politiques). — L'itération de politiques produit une suite $(\pi^{(k)})_{k \geq 0}$ de politiques stationnaires et déterministes. Le nombre de telles politiques est $|\mathcal{A}|^{|\mathcal{S}|}$.

Notons $v_k = v_{\pi_k}$ pour tout $k \geq 0$. Par propriété d'amélioration gloutonne de politique, on a $v_{k+1} \geq v_k$ pour tout $k \geq 0$. Soit $k_0 < |\mathcal{A}|^{|\mathcal{S}|}$ le plus petit entier $k \geq 0$ tel que $v_{k+1} = v_k$. Un tel entier existe car sinon les fonctions valeur v_k seraient toutes différentes pour $0 \leq k \leq |\mathcal{A}|^{|\mathcal{S}|}$, ce qui impliquerait que les politiques $(\pi^{(k)})_{0 \leq k \leq |\mathcal{A}|^{|\mathcal{S}|}}$ seraient toutes différentes, ce qui est impossible.

Puisque $v_{k_0+1} = v_{k_0}$, et qu'une fonction valeur inchangée après amélioration gloutonne implique l'optimalité, on a $v_{k_0} = v_*$. De plus, pour tout $k \geq k_0 + 1$, par définition de v_* ,

$$v_* \geq v_k \geq v_{k_0} = v_*,$$

donc $v_k = v_*$ et $\pi^{(k)}$ est optimale.

EXERCICE 3 (*Amélioration gloutonne par rapport à plusieurs politiques*). — Pour tout $1 \leq m \leq M$, on écrit

$$v_{\pi_m} = B_{\pi_m} v_{\pi_m} \leq \max_{\pi' \in \Pi_0} B_{\pi'} v_{\pi_m}.$$

En prenant le maximum sur $1 \leq m \leq M$, on obtient

$$\begin{aligned} \left(\max_{1 \leq m \leq M} v_{\pi_m} \right) &\leq \max_{1 \leq m \leq M} \max_{\pi' \in \Pi_0} B_{\pi'} v_{\pi_m} = \max_{\pi' \in \Pi_0} B_{\pi'} \left(\max_{1 \leq m \leq M} v_{\pi_m} \right) \\ &= B_* \left(\max_{1 \leq m \leq M} v_{\pi_m} \right) = B_{\pi} \left(\max_{1 \leq m \leq M} v_{\pi_m} \right), \end{aligned}$$

où on a utilisé la définition de π . En appliquant l'opérateur B_{π} qui est monotone, on a

$$B_{\pi} \left(\max_{1 \leq m \leq M} v_{\pi_m} \right) \leq B_{\pi}^2 \left(\max_{1 \leq m \leq M} v_{\pi_m} \right).$$

Par récurrence immédiate, on obtient que pour tout $k \geq 1$,

$$\max_{1 \leq m \leq M} v_{\pi_m} \leq B_{\pi}^k \left(\max_{1 \leq m \leq M} v_{\pi_m} \right) \xrightarrow{k \rightarrow +\infty} v_{\pi}.$$

D'où le résultat. On en déduit immédiatement l'inégalité correspondante pour les fonction action-valeur en appliquant l'opérateur monotone D .

