

SUJET D'ÉVALUATION  
**APPRENTISSAGE  
PAR RENFORCEMENT**  
UNIVERSITÉ PARIS–SACLAY

à rendre le 22 janvier 2024 au plus tard



Les deux premières parties du sujet doivent être traitées de façon intégrale. Dans la troisième et dernière partie, plusieurs pistes sont suggérées, et une grande liberté est laissée : un bon traitement d'une petite partie suffit pour obtenir une bonne note. Il doit s'agir d'un travail individuel. Il doit être envoyé par e-mail à l'adresse

joon.kwon@inrae.fr

et pourra comporter des fichiers sous différentes formes, par exemple : copie manuscrite scannée ou document rédigé en L<sup>A</sup>T<sub>E</sub>X, notebook Jupyter.

I. ÉTUDE THÉORIQUE D'UN NOUVEL ALGORITHME DE POINTS FIXES

Soit  $d \geq 1$ ,  $F : \mathbb{R}^d \rightarrow \mathbb{R}^d$  une application admettant un point fixe  $x_* \in \mathbb{R}^d$ ,  $\eta > 0$ , et  $x_0 \in \mathbb{R}^d$ . On définit l'algorithme Ada-FP comme suit :

$$x_{k+1} = x_k + \eta \frac{Fx_k - x_k}{\sqrt{\sum_{\ell=0}^k \|Fx_\ell - x_\ell\|_2^2}}, \quad k \geq 0, \quad (\text{Ada-FP})$$

avec la convention  $0/0 = 0$ . On note pour tout  $k \geq 0$ ,

$$\begin{aligned} u_k &= Fx_k - x_k, \\ \eta_k &= \frac{\eta}{\sqrt{\sum_{\ell=0}^k \|Fx_\ell - x_\ell\|_2^2}}, \\ D_k &= \max_{0 \leq \ell \leq k} \frac{1}{2} \|x_\ell - x_*\|_2^2. \end{aligned}$$

1) Montrer que pour tout  $k \geq 0$ ,

$$\|x_{k+1} - x_*\|_2^2 \leq \|x_k - x_*\|_2^2 + 2\eta_k u_k^\top (x_k - x_*) + \eta_k^2 \|u_k\|_2^2.$$

2) En déduire que pour tout  $k \geq 0$ ,

$$\sum_{\ell=0}^k u_\ell^\top (x_* - x_\ell) \leq \frac{D_k}{\eta_k} + \sum_{\ell=0}^k \frac{\eta_\ell \|u_\ell\|_2^2}{2}.$$

3) a) Soit  $(a_k)_{k \geq 0}$  une suite positive. Montrer que pour tout  $k \geq 0$ ,

$$\sum_{\ell=0}^k \frac{a_\ell}{\sqrt{\sum_{m=0}^{\ell} a_m}} \leq 2\sqrt{\sum_{\ell=0}^k a_\ell},$$

avec la convention  $0/0 = 0$ .

b) En déduire que pour tout  $k \geq 0$ ,

$$\sum_{\ell=0}^k u_\ell^\top (x_* - x_\ell) \leq \left( \eta + \frac{D_k}{\eta} \right) \sqrt{\sum_{\ell=0}^k \|u_\ell\|_2^2}.$$

4) On suppose que  $F$  est  $\gamma_F$ -lipschitzienne pour un certain  $0 \leq \gamma_F < 1$ .

a) Montrer que pour tout  $k \geq 0$ ,

$$\|Fx_k - x_k\|_2^2 \leq 2(Fx_k - x_k)^\top (x_* - x_k).$$

b) En déduire que pour tout  $k \geq 0$ ,

$$\min_{0 \leq \ell \leq k} \|Fx_\ell - x_\ell\|_2 \leq \frac{2}{\sqrt{k}} \left( \eta + \frac{D_k}{\eta} \right).$$

c) Qu'en déduire sur  $\min_{0 \leq \ell \leq k} \|x_\ell - x_*\|_2$  ?

On rappelle que dans un contexte de programmation dynamique (i.e. où la dynamique de transition  $p$  du MDP est disponible sous forme explicite), les itérations valeur (synchrones) pour l'évaluation d'une politique  $\pi \in \Pi_0$  (resp. pour le contrôle) sont données par

$$v_{k+1} = B_\pi v_k, \quad k \geq 1, \quad (\text{VI}_\pi^{(V)})$$

et

$$v_{k+1} = B_* v_k, \quad k \geq 1, \quad (\text{VI}_*^{(V)})$$

respectivement.

5) En utilisant (Ada-FP), définir des algorithmes analogues aux itérations valeur synchrones classiques  $(\text{VI}_\pi^{(V)})$  et  $(\text{VI}_*^{(V)})$ . On appellera  $(\text{Ada-VI}_\pi^{(V)})$  et  $(\text{Ada-VI}_*^{(V)})$  les algorithmes ainsi obtenus.

## 2. COMPARAISON EN PRATIQUE AVEC LES ALGORITHMES CLASSIQUES

Choisir un MDP de taille raisonnable, c'est-à-dire dont on puisse calculer les fonctions valeurs  $v_\pi, v_*$  ( $\pi \in \Pi_0$ ) en un temps raisonnable. On pourra soit reprendre un MDP vu en TP, soit en trouver un dans un livre, sur internet, dans une librairie (e.g. Gymnasium) ou encore en créer un soi-même, mais pour cette partie, il est nécessaire de connaître les transitions de façon explicite.

6) Se donner une politique stationnaire  $\pi \in \Pi_0$  quelconque, ainsi qu'une fonction valeur initiale  $v_0 \in \mathbb{R}^{\mathcal{S}}$  tirée aléatoirement une fois pour toutes. Comparer en pratique la vitesse de convergence de  $(\text{Ada-VI}_\pi^{(V)})$  avec celle de  $(\text{VI}_\pi^{(V)})$ . On pourra tracer, avec une échelle logarithmique en ordonnée, les quantités

$$\|v_k - v_\pi\|_\infty \quad \text{et} \quad \|v_k - B_\pi v_k\|_\infty$$

en fonction de  $k$ . Essayer différentes valeurs pour  $\eta > 0$ .

7) Même question pour  $(\text{Ada-VI}_*^{(V)})$  et  $(\text{VI}_*^{(V)})$ .

## 3. EXTENSIONS

Reprendre la démarche de la Section 2 en incorporant, par exemple, un ou plusieurs des aspects suivants :

- itérations de fonctions action-valeur,
- itérations asynchrones,
- méthode d'apprentissage par renforcement utilisant des estimateurs stochastiques des opérateurs de Bellman,
- approximation de la fonction valeur par une classe paramétrique,
- variante de (Ada-FP) définie composante par composante par :

$$x_{k+1,j} = x_{k,j} + \eta \frac{(F x_k)_j - x_{k,j}}{\sqrt{\sum_{\ell=0}^k (F(x_\ell)_j - x_{\ell,j})^2}}, \quad 1 \leq j \leq d, \quad k \geq 0.$$

- tout autre aspect potentiellement pertinent.

On commencera par définir précisément les algorithmes considérés avant de les implémenter. On pourra également tenter une étude théorique des algorithmes obtenus, en gardant à l'esprit qu'il n'est pas évident de démontrer des garanties dans tous les cas.

