

Processus de décisions markoviens.

Un PDM est défini par (S, A, T, R, γ) où

- S est un ensemble fini d'états
on notera $(s_1, \dots, s_n) = S$
- A est un ensemble fini d'actions
on notera $(a_1, \dots, a_m) = A$
- On notera S_t l'état à l'instant $t \in \mathbb{N}$
 A_t l'action prise à l'instant $t \in \mathbb{N}$
- T est la matrice de transition

$$T_{SS'}^a = p(S_{t+1} = s' \mid S_t = s, A_t = a)$$

probabilité de se trouver en s' sachant qu'on a pris l'action a dans l'état s .

- R est la matrice de récompense

$$R_s^a = \mathbb{E}(r_{t+1} \mid S_t = s, A_t = a)$$

espérance d'obtenir la récompense $r_{t+1} \in \mathbb{R}$ après avoir pris l'action a dans l'état s .

- γ est le taux d'escompte

Quelle est la valeur aujourd'hui d'obtenir 1 unité demain.

N.B. On trouve d'autres formulations pour modéliser la récompense :

- R_s : récompense quand on arrive dans l'état s
- $R_{SS'}$: récompense obtenue après avoir pris l'action a dans l'état s et être arrivé en s'

→ on peut ensuite reformuler et écrire des résultats équivalents.

def : la **fonction de valeur** pour la politique π est

$$v_\pi(s) = \mathbb{E} \left(\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid S_t = s, \pi \right)$$

def : la **fonction de valeur optimale** v^* est $\max_{\pi} v_\pi$.

prop : v_π satisfait l'équation de Bellman

$$v_\pi(s) = \sum_{a \in A} \pi(a|s) R_s^a + \gamma \sum_{s' \in S} \sum_{a \in A} \pi(a|s) T_{SS'}^a v_\pi(s')$$

- v^* satisfait l'équation de programmation dynamique

$$v^*(s) = \max_{a \in A} R_s^a + \gamma \sum_{s' \in S} T_{SS'}^a v^*(s')$$

$$\begin{aligned} v_\pi(s) &= \mathbb{E} \left(\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid S_t = s, \pi \right) \\ &= \mathbb{E} \left(r_{t+1} + \sum_{k=1}^{\infty} \gamma^k r_{t+k+1} \mid S_t = s, \pi \right) \\ &= \mathbb{E}(r_{t+1} \mid S_t = s, \pi) + \mathbb{E} \left(\sum_{k=1}^{\infty} \gamma^k r_{t+k+1} \mid S_t = s, \pi \right) \\ &= \sum_{a \in A} \pi(a|s) R_s^a + \gamma \sum_{a \in A} \pi(a|s) \sum_{s' \in S} T_{SS'}^a \sum_{k=1}^{\infty} \mathbb{E}(\gamma^{k-1} r_{t+k+1} \mid S_{t+1} = s', S_t = s, A_{t+1} = a) \\ &= \sum_{a \in A} \pi(a|s) R_s^a + \gamma \sum_{a \in A} \sum_{s' \in S} T_{SS'}^a \underbrace{\sum_{k=0}^{\infty} \mathbb{E}(\gamma^k r_{t+1+k+1} \mid S_{t+1} = s', A_{t+1} = a)}_{v_\pi(s')} \\ &= \sum_{a \in A} \pi(a|s) R_s^a + \gamma \sum_{s' \in S} \sum_{a \in A} T_{SS'}^a v_\pi(s') \quad \square \end{aligned}$$

$$\begin{aligned} v^*(s) &= \max_{\pi} \mathbb{E} \left(\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid S_t = s, \pi \right) \\ &= \max_{a, \pi'} R_s^a + \gamma \sum_{s' \in S} T_{SS'}^a v_{\pi'}(s') \\ &= \max_a \left[R_s^a + \max_{\pi'} \gamma \sum_{s' \in S} T_{SS'}^a v_{\pi'}(s') \right] \\ &= \max_a \left[R_s^a + \gamma \sum_{s' \in S} T_{SS'}^a \max_{\pi'} v_{\pi'}(s') \right] \\ &= \max_a \left[R_s^a + \gamma \sum_{s' \in S} T_{SS'}^a v^*(s') \right] \quad \square \end{aligned}$$

On peut ainsi définir une politique optimale $\pi^* = \operatorname{argmax}_{\pi} v_\pi$

On est maintenant prêt à définir deux opérateurs.

- $T^\pi : \mathbb{R}^n \rightarrow \mathbb{R}^n$
 $W \mapsto T^\pi W(s) = \sum_{s' \in S} \pi(a|s) T_{SS'}^a W(s')$
- $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$
 $W \mapsto TW(s) = \max_{a \in A} \left[R_s^a + \gamma \sum_{s' \in S} T_{SS'}^a W(s') \right]$

Ceci nous permet de revisiter v_π et v^* sous forme matricielle

$$v^\pi = (I - \gamma T^\pi)^{-1} r^\pi \quad \text{où} \quad \begin{cases} T^\pi(s, s') = \sum_{a \in A} \pi(a|s) T_{SS'}^a \\ r^\pi(s) = \sum_{a \in A} \pi(a|s) R_s^a \end{cases}$$

- v^π est l'unique point fixe de T^π
- v^* est l'unique point fixe de T^*
- $\pi^* = \operatorname{argmax}_{a \in A} \left[R_s^a + \gamma \sum_{s' \in S} T_{SS'}^a v^*(s') \right]$ est optimale
- $\forall W \in \mathbb{R}^n \quad \lim_{k \rightarrow +\infty} (T^\pi)^k W = v^\pi$
 $\lim_{k \rightarrow \infty} T^k W = v^*$
- Opérateurs contractants en $\|\cdot\|_\infty$

$$\forall W_1, W_2 \in \mathbb{R}^n \quad \begin{aligned} \|T^\pi W_1 - T^\pi W_2\|_\infty &\leq \gamma \|W_1 - W_2\|_\infty \\ \|TW_1 - TW_2\|_\infty &\leq \gamma \|W_1 - W_2\|_\infty \end{aligned}$$

$$\begin{aligned} |TW_1(s) - TW_2(s)| &= \left| \max_{a \in A} \left[R_s^a + \gamma \sum_{s' \in S} T_{SS'}^a W_1(s') \right] - \max_{a \in A} \left[R_s^a + \gamma \sum_{s' \in S} T_{SS'}^a W_2(s') \right] \right| \\ &\leq \gamma \max_{a \in A} \sum_{s' \in S} T_{SS'}^a |W_1(s') - W_2(s')| \\ &\leq \gamma \|W_1 - W_2\|_\infty \end{aligned}$$

Démonstrations

1. On a $v^\pi = r^\pi + \gamma T^\pi v^\pi \Leftrightarrow (I - \gamma T^\pi) v^\pi = r^\pi$
La matrice T^π est stochastique, donc ses valeurs propres sont de module ≤ 1
Donc les valeurs propres de $I - \gamma T^\pi$ ont de module $\geq 1 - \gamma$
et donc $I - \gamma T^\pi$ est inversible.

2 et 3. $\Rightarrow T$ et T^π sont contractants.