

On considère les notations suivantes

$$P(s', r | s, a) = P[S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a]$$

$$R(s, a) = E[R_{t+1} | S_t = s, A_t = a] = \sum_{r \in \mathcal{R}} r P(r | S_t = s, A_t = a)$$

$$V_{\pi}(s) = E_{\pi} \left[\sum_{t=0}^{\infty} R(S_t) \mid S_0 = s \right]$$

On souhaite montrer que $V_{\pi}(s) = E_{\pi} \left[\sum_{t=0}^{\infty} R(S_t) \mid S_0 = s \right]$

$$= E_{\pi} [R(s, \pi(s))] + \sum_{s' \in \mathcal{S}} P(s' | s, \pi(s)) V_{\pi}(s')$$

On note G_t la somme des futur rewards :

$$G_t = \sum_{R=t+1}^{\infty} R_R \quad \text{avec } R(S_t) = R_{t+1}$$

On peut alors réécrire $V_{\pi}(s)$ de la façon suivante :

$$V_{\pi}(s) = E_{\pi} [G_t | S_t = s]$$

$$V_{\pi}(s) = E_{\pi} [R_{t+1} + G_{t+1} | S_t = s]$$

$$= E_{\pi} [R_{t+1} | S_t = s] + E_{\pi} [G_{t+1} | S_t = s] \quad (a)$$

En développant le terme $E_{\pi} [R_{t+1} | S_t = s]$:

$$E_{\pi} [R_{t+1} | S_t = s] = \sum_{r \in \mathcal{R}} r P(r | S_t = s)$$

$$= \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} P(r, s' | s) = \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} \frac{P(R_{t+1}=r, S_{t+1}=s', S_t=s)}{P(S_t=s)}$$

$$= \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}} \frac{P(R_{t+1}=r, S_{t+1}=s', S_t=s, A_t=a)}{P(S_t=s)}$$

$$= \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}} \frac{P(R_{t+1}=r, S_{t+1}=s', S_t=s, A_t=a)}{P(S_t=s)} \times \frac{P(S_t=s, A_t=a)}{P(S_t=s, A_t=a)}$$

$$= \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}} P(s', r | s, a) \times P(A_t=a | S_t=s)$$

On note $\pi(a | s) = P(A_t=a | S_t=s)$, Ainsi on peut écrire

$$E_{\pi} [R_{t+1} | S_t = s] = \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}} P(s', r | s, a) \pi(a | s)$$

En regardant désormais le second terme de (1) : $E_{\pi}[G_{t+n} | S_t = s]$
 (On pourrait bien le réécrire selon $V_{t+n, \pi}(s') = E_{\pi}[G_{t+n} | S_{t+n} = s']$)
 Pour avoir une expression de récurrence entre chaque $V_{t, \pi}(\cdot)$

$$\begin{aligned}
 E_{\pi}[G_{t+n} | S_t = s] &= \sum_{g \in G} g P[g | S_t = s] = \sum_{g \in G} g \sum_{s' \in S} P(g, S_{t+n} = s' | S_t = s) \\
 &= \sum_{g \in G} g \sum_{s' \in S} \sum_{a \in A} P(g, S_{t+n} = s', A_t = a | S_t = s) \\
 &= \sum_{g \in G} g \sum_{s' \in S} \sum_{a \in A} \sum_{r \in R} P(g, S_{t+n} = s', A_t = a, R_{t+n} = r | S_t = s) \\
 &= \sum_{g \in G} g \sum_{s'} \sum_a \sum_r \frac{P(G_{t+n} = g, S_{t+n} = s', A_t = a, R_{t+n} = r, S_t = s)}{P(S_t = s)} \\
 &= \sum_{g \in G} g \sum_{s'} \sum_a \sum_r \frac{P(g, s', a, r, s)}{P(s)} \times \frac{P(s', a, r, s)}{P(s', a, r, s)} = \sum_{g \in G} g \sum_{s'} \sum_a \sum_r \frac{P(g, s', a, r, s)}{P(s', a, r, s) \times P(s)} \\
 &= \sum_{g \in G} g \sum_{s'} \sum_a \sum_r \underbrace{P(g | s', a, r, s)} \times P(s', a, r | s)
 \end{aligned}$$

$$P(g | s', a, r, s) = P(G_{t+n} = g | S_{t+n} = s', A_t = a, R_{t+n} = r, S_t = s)$$

On rappelle que $G_{t+n} = R_{t+2} + \dots + R_T$

De manière explicite, on sait que R_{t+2} ne dépend que de A_{t+1} , S_{t+1}
 De manière plus formelle on utilise les propriétés d'un processus Markovien pour réécrire :

$$P(G_{t+n} = g | S_{t+n} = s', A_t = a, R_{t+n} = r, S_t = s) = P(G_{t+n} = g | S_{t+n} = s')$$

$$\text{Ainsi on a : } E_{\pi}[G_{t+n} | S_t = s] = \sum_g g \sum_{s'} \sum_a \sum_r P(g | s') P(s', a, r | s)$$

$$= \sum_g g \sum_{s'} \sum_a \sum_r P(g | s') \times \frac{P(s', a, r, s)}{P(s)} = \sum_g g \sum_{s'} \sum_a \sum_r P(g | s') \times \frac{P(s', a, r, s)}{P(s)} \times \frac{P(s)}{P(s)}$$

$$= \sum_g g \sum_{s'} \sum_a \sum_r P(g | s') \times \frac{P(s', a, r, s)}{P(a, s)} \times \frac{P(a, s)}{P(s)}$$

$$= \sum_g g \sum_{s'} \sum_d \sum_r P(g|s') \cdot P(s', r | d, s) \cdot \pi(d|s)$$

$$= \sum_r \sum_{s'} \sum_d \left(\sum_{g \in G} g P(g|s') \right) \cdot P(s', r | d, s) \cdot \pi(d|s)$$

$$= \sum_r \sum_{s'} \sum_d E_{\pi}[G_{t+1} | S_{t+1}=s'] \cdot P(s', r | d, s) \cdot \pi(d|s)$$

$$= \sum_r \sum_{s'} \sum_d V_{(t+1)\pi}(s') \cdot P(s', r | d, s) \cdot \pi(d|s)$$

Ainsi on a donc l'expression de $V_{t\pi}(s)$ suivant:

$$V_{t\pi}(s) = \sum_{r \in R} \sum_{s' \in S} \sum_{d \in A} \left[r P(r, s' | s, d) \pi(d|s) + V_{(t+1)\pi}(s') P(r, s' | s, d) \right] \pi(d|s)$$

$$V_{t\pi}(s) = \sum_{r \in R} \sum_{s' \in S} \sum_{d \in A} \pi(d|s) \left[r P(r, s' | s, d) + V_{(t+1)\pi}(s') P(r, s' | s, d) \right]$$

$$V_{t\pi}(s) = \sum_{d \in A} \pi(d|s) \sum_{s' \in S} \sum_{r \in R} P(r, s' | s, d) [r + V_{(t+1)\pi}(s')]$$

Pour retrouver la forme voulue: $E_{\pi}[R(A\pi(s))] + \sum_{s' \in S} P(s' | s, \pi(s)) V_{(t+1)\pi}(s')$

$$V_{t+1}(s) = \sum_{r \in R} \sum_{s' \in S} \sum_{d \in A} \pi(d|s) r \cdot P(r, s' | s, d) + \sum_{r \in R} \sum_{s' \in S} \sum_{d \in A} P(r, s' | s, d) \pi(d|s) V_{t+1}(s')$$

$$= \sum_{r \in R} \sum_{s' \in S} \sum_{d \in A} r \cdot P(r, s' | s, d)$$

$$+ \sum_{s' \in S} \sum_{d \in A} \pi(d|s) V_{t+1}(s') \sum_{r \in R} P(r, s' | s, d) \quad \text{PT}$$

$$= \sum_{r \in R} \sum_{d \in A} \sum_{s' \in S} P(r, s' | s, d) \quad \text{PT}$$

$$+ \sum_{s' \in S} \sum_{d \in A} \pi(d|s) V_{t+1}(s') P(s' | s, d)$$

$$= \sum_r \sum_d P(r|s, d) + \sum_{s'} V_{t+1}(s') P(s' | s, d) \pi(d|s)$$

$$= \sum_d E_{\pi}[R_{t+1} | S_t=s, A_t=d] + \sum_{s'} V_{t+1}(s') \sum_d P(s' | s, d) \pi(d|s)$$

(3)

PT: proba totale

En considérant que $(A_t)_t$ est pris par la policy $\pi(a)$, on a donc

$$\sum_a \mathbb{E}(R_{t+1} | S_t = s, A_t = \pi(a)) = \sum_a R(s, \pi(a)) = \mathbb{E}_\pi[R(s, \pi(a))]$$

De plus :

$$\begin{aligned} \sum_{s'} V_{\pi(t+1)}(s') \sum_a P(s' | s, a) \pi(a) \\ = \sum_{s'} V_{\pi(t+1)}(s') P(s' | s, \pi(s)) \end{aligned}$$

Ainsi on retrouve bien la "formule"

$$V_{\pi(t+1)}(s) = \mathbb{E}_\pi[R(s, \pi(s))] + \sum_{s' \in \mathcal{S}} V_{\pi(t+1)}(s') P(s' | s, \pi(s))$$