

Étape 4.

Apprentissage dans l'incertain

Rappel.

$$\text{MDP} = (\mathcal{S}, \mathcal{A}, p, r)$$

$$- \mathcal{S} = \{s_1, s_2, \dots, s_n\}$$

$$- \mathcal{A} = \{a_1, a_2, \dots, a_m\}$$

$$- p: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$$

$$- r: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$$

ce qu'on a vu jusque là :

p et r sont connues.

\Rightarrow Equations de Bellman

+ policy iteration / value iteration

Dans ce qui suit :

incertain = p et r ne sont pas connus

l'agent va devoir interagir avec l'environnement

- on veut pour apprendre q et r .
soit π une politique.

l'eq. de Bellman :

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) \left(r + \gamma V_{\pi}(s') \right)$$

si π , p et r sont déterministes et si on instancie l'eq. en une trajectoire :

$$V_{\pi}(s_t) = \frac{r_t}{\gamma} + \gamma V_{\pi}(s_{t+1})$$

$$\Rightarrow \frac{r_t}{\gamma} + \gamma V_{\pi}(s_{t+1}) - V_{\pi}(s_t) = 0$$

idée : V_{π} est inconnue mais au fil de l'apprentissage, on essaie s'approcher de sa vraie valeur...

On commence donc avec une valeur arbitraire \hat{V}

- Demand lagged effective new
interaction,

Algorithme Temporal Difference (TD)

Input: S, A, γ et π

$$\hat{V}_\pi(s) \leftarrow 0, \forall s \in S$$

$$m(s) \leftarrow 0, \forall s \in S$$

Répéter

initialiser l'état initial s_0

$$t \leftarrow 0$$

Répéter

émettre l'action $a_t = \pi(s_t)$

observer r_t et s_{t+1}

$$\hat{V}_\pi(s_t) \leftarrow \hat{V}_\pi(s_t)$$

$$+ \alpha(s_t, m(s_t)) \times$$

$$\left[r_t + \gamma \hat{V}_\pi(s_{t+1}) - \hat{V}_\pi(s_t) \right]$$

$$m(s) \leftarrow m(s) + 1$$

$$t \leftarrow t + 1$$

jusqu'à s_t état final

jusqu'à critère d'arrêt vérifié

Propriété. TD converge (asymptotiquement)

Si

- chaque état est visité une infinité de fois

- $\forall s \in \mathcal{S}$

$\sum_t \alpha_t^2 (r, m(s)) = +\infty$
à laquelle s est visité

$\sum_t \alpha_t^2 (r, m(s)) < +\infty$

Une solution:

$$\alpha_t (r, m(s)) = \frac{1}{n(s) + 1}$$

α est nommé le pas d'apprentissage.

