

de passer de l'état s à l'état s' pour
chaque des décisions d , Il y a
également un coût C_d pour chaque des
décisions

- On note X_s la rentabilité de la
machine en état s

- Le jour $j+1$, notre décision de la
veille au soir a introduit un nouvel
état s' et une récompense

$$R_{d,s'} = C_d + X_{s'}$$

Question : Souvent optimiser la
décision que l'on doit prendre

Réponse : Cela dépend...

optimiser le gain (moyen) sur
un jour n'est pas du tout la
même chose que d'optimiser le
gain sur 10 ans ...

Formellement:

Définition: MDP = (S, A, T, R)

- S = ens. des états du système

- A = ens. des décisions/actions possibles

- $T: S \times A \times S \rightarrow [0, 1]$

$T(s, a, s') = \text{Pr}(\text{le syst. passe de l'état } s \text{ à l'état } s' \text{ lorsque la décision } a \text{ est prise})$

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

$R(s, a, s') = \text{récompense immédiate lorsque l'action entreprise est } a \text{ et que le syst. passe de } s \text{ à } s'$

Comme indiqué dans le chapitre précédent, on note S_t et R_t les v.o. indiquant l'état à l'instant t et la récompense à l'instant $t+1$ resp.

R_t et S_t vérifient la propriété de Markov :

$$\begin{aligned} \mathbb{P}(R_{t+1}=r' \mid S_t=s \wedge S_{t-1}=s_{t-1} \wedge \dots \wedge S_0=s_0) \\ = \mathbb{P}(R_{t+1}=r' \mid S_t=s) \end{aligned}$$

idem pour R_t .

La dynamique du MDP est alors décrite par la fonction p :

$$p(s', r \mid s, a) = \mathbb{P}(S_{t+1}=s' \wedge R_{t+1}=r \mid S_t=s \wedge A_t=a)$$

$\forall s, s', r, a \in \mathcal{S} \times \mathcal{S} \times \mathcal{R} \times \mathcal{A}$.

En effet, si on connaît $p(s', r | s, a)$:

$$- p(s' | s, a) = \Pr(s_{t+1} = s' | s_t = s \wedge A_t = a)$$

$$\Lambda(s | s, a) = \sum_{r \in \mathcal{R}} p(s', r | s, a)$$

$\Lambda \Rightarrow$ Prob. de transitions d'état

$$- r(s, a) = \mathbb{E}(R_{t+1} | s_t = s, A_t = a)$$

$$= \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s', r | s, a) \quad \triangle$$

$\underbrace{\hspace{10em}}_{p(s' | s, a)} \quad \triangle$

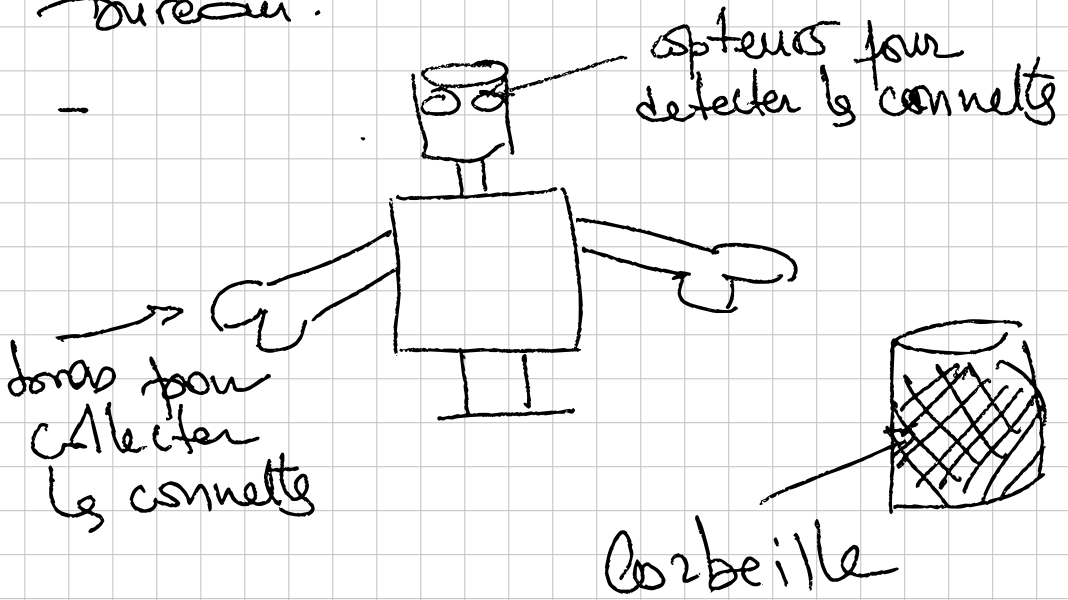
\Rightarrow Espérance de la récompense
pour s et a

$$- r(s, a, s') = \mathbb{E}(R_{t+1} | s_t = s \wedge A_t = a \wedge s'_{t+1} = s')$$

$$= \sum_{r \in \mathcal{R}} r \frac{p(s', r | s, a)}{p(s' | s, a)}$$

En gros, si on connaît $p(s', r | s, a)$,
on connaît tout (du moins localement)

Exemple: Un robot pour le recyclage.
- Le robot collecte des cannettes vides dans un bureau.



En fonction du niveau de charge de la batterie, l'agent qui guide le comportement du robot peut choisir:

- chercher des cannettes pdt une période
- rester sur place et attendre que quelqu'un lui apporte une cannette
- retourner à la base pour recharger la batterie

$$= \{ \text{high} \}$$

$$= \{ \text{search} \}$$

$$\text{act}(h) =$$

Avant de continuer, un jeu de formalisme

$$P = \{ \text{high}, \text{low} \}$$

$$\text{act} = \{ \text{search}, \text{wait}, \text{recharge} \}.$$

$$\text{avec } \text{act}(h) = \{ s, w \}; \text{act}(e) = \{ s, w, r \}$$

Le but du robot est de collecter un maximum de cannettes \Rightarrow "févihaliger" la recherche.



Cela décharge la batterie alors que rester sur place l'conservée.

La batterie finit se décharger complètement ce qui implique une intervention humaine

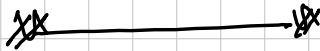
Hypothèses

- Si la batterie est "high" et que le robot fait une période de recherche, elle passe à "low" avec prob. α .
- Si la batterie est "low", alors elle y reste avec prob β .
- Chaque cannette récupérée rapporte 1
- si le robot doit être secouru alors la pénalité est de -3

On note

$\Gamma_S =$ l'espérance du nombre de courtes passées dans une période de recherche

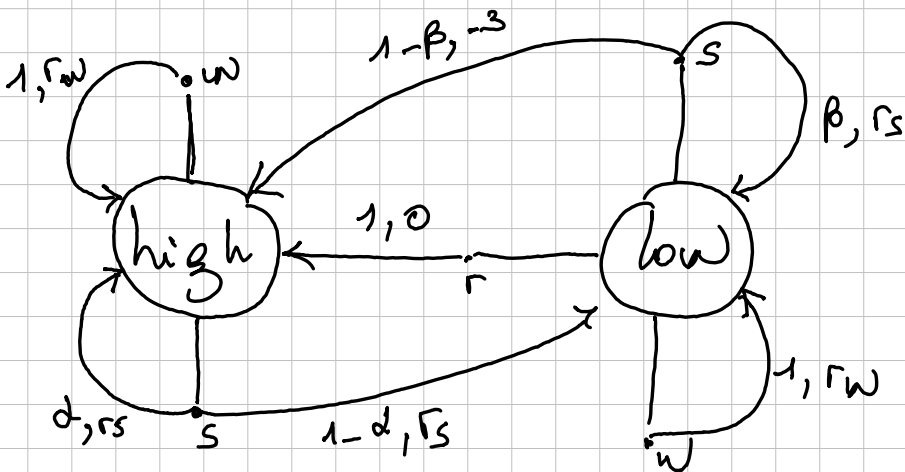
$\Gamma_W =$ α α α α α
 α α α d'attente



MDP:

s	a	s'	$\phi(s' s,a)$	$r(s,a,s')$
h	s	h	$1-\alpha$	r_s
	w	h	α	r_w
l	s	h	$1-\beta$	r_s
	w	l	β	r_w

Graph de transition



Rappelons que le but de l'agent est
d'optimiser sa récompense (récompense cumulée)
Soit $t \geq 0$, l'agent recevra des récompenses
 R_{t+1}, R_{t+2}, \dots

\Rightarrow Il doit chercher à maximiser
l'espérance G_t du gain.

$$G_t = f(R_{t+1}, R_{t+2}, \dots)$$

Exemple: $f = \bar{\Sigma}$:

$$G_t = R_{t+1} + R_{t+2} + \dots + R_T$$

T étant le dernier pas, il faut
être un épisode ou encore correspondre
à un état final (ou absorbant).

⚠ Si la tâche ne peut pas être découpée
en épisodes, la définition de G_t
devient problématique car $T = \infty$ et
 G_t non convergente...

⇒ On introduit un facteur de dépréciation $0 < \gamma \leq 1$:

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k \geq 0} \gamma^k R_{t+k+1} \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned}$$

Remarques:

- La nouvelle déf. de G_t englobe la précédente ($G_t = 0 \forall t > T$ et $\gamma = 1$)
- Si la récompense est finie et non nulle et si $\gamma < 1$, alors G_t est finie.

$$\begin{aligned} (G_t &= R_{t+1} + \gamma(R_{t+2} + \gamma G_{t+2})) \\ &= R_{t+1} + \gamma R_{t+2} + \gamma^2 G_{t+2}. \end{aligned}$$

???

Politiques et Fonctions de valeurs

But ultime de tout algo de RL:
trouver une politique $\pi : \mathcal{S} \rightarrow \mathcal{A}$ qui
soit optimale.

i.e. une f^{π} qui à chaque état s ,
indique l'action à exécuter

La politique π ainsi trouvée est censée
maximiser les récompenses.

Requ Une politique peut également
être stochastique :

$$\pi : \mathcal{A}, \mathcal{S} \rightarrow [0, 1]$$

$$\pi(a, s) = \mathbb{P}(A_t = a \mid S_t = s)$$

La fonction de valeur d'un état s , sous
la politique π est le gain si on commence
à l'état s et que l'on suit la
politique π :

$$\begin{aligned}
 \frac{Q}{\pi}(s) &= \mathbb{E}(G_t \mid s_t = s) \\
 &= \mathbb{E}\left(\sum_{k \geq 0} \gamma^k R_{t+k+1} \mid s_t = s\right)
 \end{aligned}$$

Fonction valeur de actions :

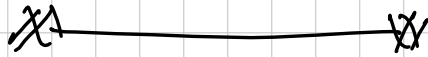
$$\begin{aligned}
 q(s, a) &= \mathbb{E}(G_t \mid s_t = s \wedge A_t = a) \\
 &= \mathbb{E}\left(\sum_{k \geq 0} \gamma^k R_{t+k+1} \mid s_t = s\right)
 \end{aligned}$$

De manière récursive :

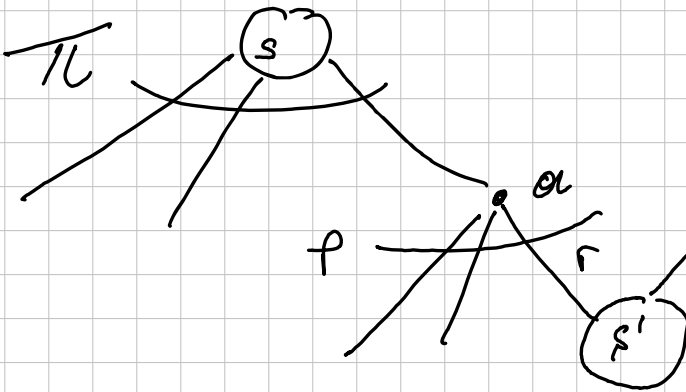
$$\begin{aligned}
 \frac{Q}{\pi}(s) &= \mathbb{E}_{\pi}(G_t \mid s_t = s) \\
 &= \mathbb{E}_{\pi}(R_{t+1} + \gamma G_{t+1} \mid s_t = s) \\
 &= \sum_a \pi(a|s) \sum_{s', r} p(s', r \mid s, a) \left[r + \gamma \mathbb{E}(G_{t+1} \mid s_{t+1} = s') \right] \\
 &= \sum_a \pi(a|s) \sum_{s', r} p(s', r \mid s, a) \left(r + \gamma V_{\pi}(s') \right)
 \end{aligned}$$

Eq. de Bellmann :

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



back-up diagram



→
- Première séance arrêtée là. -
(le 11/10/2022)

Exemple : retoune à l'exemple jacob.

Soit π la politique suivante :

$$\pi = \left\{ \begin{array}{l} \text{high} \rightarrow \text{search;} \\ \text{low} \rightarrow \text{wait} \end{array} \right\}.$$



$$\left\{ \begin{array}{l} v_{\pi}(h) = \alpha(r_s + \delta v_{\pi}(h)) + (1-\alpha)(r_s + \delta v_{\pi}(l)) \\ v_{\pi}(l) = 1(r_w + \delta v_{\pi}(l)) \end{array} \right.$$

avec des variables : $x = v_{\pi}(h)$; $y = v_{\pi}(l)$

$$x = \alpha(r_s + \delta x) + (1-\alpha)(r_s + \delta y) \quad (1)$$

$$y = r_w + \delta y \quad (2)$$

$$(2) \Rightarrow \boxed{y = r_w / (1-\delta)}$$

$$\text{Dans (1)} : x = r_s + \alpha\delta x + (1-\alpha)\delta \frac{r_w}{1-\delta}$$

$$\Rightarrow (1-\alpha\delta)x = (1-\alpha)r_s + (1-\alpha)\delta \frac{r_w}{1-\delta}$$

$$\Rightarrow (1-\alpha\delta - (1-\alpha)\delta)x = (1-\alpha)r_s + (1-\alpha)\delta \frac{r_w}{1-\delta}$$

$$\Rightarrow \boxed{x = \frac{(1-\alpha)r_s + (1-\alpha)\delta \frac{r_w}{1-\delta}}{1-\alpha\delta - (1-\alpha)\delta}}$$

Application numérique :

$$\delta = 10^{-1} ; \alpha = \frac{1}{2} ; \beta = \frac{1}{2}$$

$$r_s = 10 ; r_w = 2$$

$$\begin{aligned} \Rightarrow \frac{\sigma}{\pi}(h) = \alpha &= \frac{(1 - 10^{-1})10 + (1 - 1/2)10^{-1}}{1 - 10^{-1} - \frac{1}{2}10^{-1}} \\ &= \frac{(10 - 1) + 0.05}{0.9 - 0.05} = \frac{9.05}{0.85} = \frac{905}{85} \\ &= 10.65 \end{aligned}$$

$$\frac{\sigma}{\pi}(e) = \gamma = \frac{2}{1 - 10^{-1}} = \frac{2}{0.9} = \frac{20}{9} = 2.22$$

Politiques optimales / fonctions de val. optimale

Déf. π et π' deux politiques

$$\left\{ \pi \succcurlyeq \pi' \Leftrightarrow v_{\pi}(s) \geq v_{\pi'}(s) \quad \forall s \in \mathcal{S} \right.$$

\exists au moins une politique π_* meilleure que toutes les autres politiques

π_* : politique optimale

A ces politiques est associée une fonction de valeurs optimale

$$v_*(s) = \max_{\pi} v_{\pi}(s) \quad \forall s \in \mathcal{S}$$

On y associe également une fonction de valeurs état-action optimale :

$$q_*(s, a) = \max_{\pi} q_{\pi}(s, a) \quad \forall s, a \in \mathcal{S} \times \mathcal{A}$$

Remarque,

$$q_*(s, a) = \mathbb{E}(R_{t+1} + \gamma v_x^*(s_{t+1}) \mid s_t = s \wedge A_t = a)$$

Eq. de Bellman pour v_* :

$$v_*(s) = \max_a q_*(s, a)$$

$$= \max_a E_{\pi} (G_t \mid s_t = s \wedge A_t = a)$$

$$= \max_a E_{\pi} (R_{t+1} + \gamma G_{t+1} \mid s_t = s \wedge A_t = a)$$

$$= \max_a E_{\pi} (R_{t+1} + \gamma v_*(s_{t+1}) \mid s_t = s \wedge A_t = a)$$

$$v_*(s) = \max_a \sum_{s', r} p(s', r \mid s, a) (r + \gamma v_*(s'))$$

idem pour q_* :

$$q_*(s, a) = \sum_{s', r} p(s', r \mid s, a) (r + \max_{a'} q_*(s', a'))$$

