

1 Markov Decision Process (MDP)

D'après Wikipédia: Un processus de décision markovien est un processus de contrôle stochastique discret. À chaque étape, le processus est dans un certain état \mathbf{s} et l'agent choisit une action \mathbf{a} . La probabilité que le processus arrive à l'état \mathbf{s}' est déterminée par l'action choisie. Plus précisément, elle est décrite par la fonction de transition d'états $\mathbf{T}(\mathbf{s}, \mathbf{a}, \mathbf{s}')$. Donc, l'état \mathbf{s}' dépend de l'état actuel \mathbf{s} et de l'action \mathbf{a} sélectionnée par le décideur. Cependant, pour un état \mathbf{s} et une action \mathbf{a} , le prochain état est indépendant des actions et états précédents. On dit alors que le processus satisfait la propriété de Markov.

Quand le processus passe de l'état \mathbf{s} à l'état \mathbf{s}' avec l'action \mathbf{a} , l'agent gagne une récompense \mathbf{r} .

2 Q-Learning

Le Q-Learning, basé sur le processus de décision markovien, est une technique d'apprentissage automatique utilisée en intelligence artificielle, plus particulièrement en apprentissage par renforcement.

2.1 Q-Function

La Q-Function est la base du Q-Learning. Pour un état \mathbf{s} et une action \mathbf{a} , elle donne la récompense esperée $\mathbf{Q}(\mathbf{s}, \mathbf{a})$.

2.2 Policy

La Politique π a l'état \mathbf{s} est la façon de choisir l'action \mathbf{a} qui maximise la récompense esperée $\mathbf{Q}(\mathbf{s}, \mathbf{a})$. C'est à dire, si on est à l'état \mathbf{s} et qu'on choisit l'action \mathbf{a} selon la politique π , alors, on aura une récompense optimale.

$$\pi(s) = \operatorname{argmax}_a(Q(s,a))$$

2.3 Rewards

La récompense esperée R_t à l'instant t pour un état s_t et une action a_t est la somme des récompenses futures.

$$R_t = Q(s_t, a_t)$$

$$R_t = r_t + r_{t+1} + r_{t+2} + r_{t+3} + \dots + r_n$$

Mais, sachant que le processus est stochastique, alors plus on va dans le future, moins les récompenses sont évidentes, ainsi, on introduit un **facteur de discontinuité** γ .

$$R_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots + \gamma^{n-t} r_n$$

Ainsi,

$$\begin{aligned}
R_t &= r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots + \gamma^{n-t} r_n \\
R_t &= r_t + \gamma(r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots + \gamma^{n-(t+1)} r_n) \\
R_t &= r_t + \gamma R_{t+1}
\end{aligned} \tag{1}$$

2.4 Optimal Reward

La recompense à l'instant t et à l'état s_t , en choisissant l'action a_t suivant la politique π est donc optimale et a la valeur $Q^\pi(s_t, a_t)$.

$$Q^\pi(s_t, a_t) = r_t + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})$$

On peut aussi écrire:

$$Q(s_t, a_t) = r_t + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1}), \quad a_t = \pi(s_t)$$

Note,

$$\begin{aligned}
Q^\pi(s_t, a_t) &= \max R_t \\
&= \max(r_t + \gamma R_{t+1}) \\
&= r_t + \gamma \max(R_{t+1}) \\
Q^\pi(s_t, a_t) &= r_t + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})
\end{aligned} \tag{2}$$

2.5 Loss

$$Q(s,a) = \max_{a'} R(s,a')$$

$$Q(s_t, a_t) = r_t + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})$$

En effet,

$$\begin{aligned} R_t &= r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots + \gamma^{n-t} r_n \\ &= r_t + \gamma(r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots + \gamma^{n-(t+1)} r_n) \\ &= r_t + \gamma R_{t+1} \end{aligned}$$

$$\begin{aligned} \max R_t &= \max(r_t + \gamma R_{t+1}) \\ &= r_t + \gamma \max(R_{t+1}) \end{aligned} \tag{3}$$

$$Q(s_t, a_t) = r_t + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})$$

A l'instant \mathbf{t} , le processus étant à l'état s_t , en choisissant l'action \mathbf{a} , le processus arrive à l'état s_{t+1} et gagne une récompense R_t .

$$R_t = r_t + r_{t+1}$$

Il utilise le modèle MDP à n-étapes (fini) telle pour un état \mathbf{s} et une action \mathbf{a} , $\mathbf{Q}(\mathbf{s}, \mathbf{a})$ représente l'optimum des récompenses futurs.

$$loss = \underbrace{(r + \gamma \max_{a'} Q'(s', a'))}_{target} - \underbrace{Q(s, a)}_{prediction})^2$$