

1 Configuration triviale

1.1 Processus de décision markovien

ENTRÉE : $S(\{i | i = 0..N^2\})$, $\pi(S)$, γ , ϵ

SORTIE : $\pi^*(s)$, $v_\pi^*(S)$

Initialisation

$v_{\pi,0}^*(S) \leftarrow 0$

POUR $t = 1$ à t_{obs} FAIRE

POUR s dans S FAIRE

$v_{\pi,t}(S) \leftarrow \left(\sum_{s' \in S} p(s, a, s') \left(r(s, a, s') + \gamma v_{\pi,t-1}^*(s') \right) \right)_{a \in \pi(s)}$

$v_{\pi,t}^*(S) \leftarrow \max_{a \in \pi(S)} v_{\pi,t}(S)$

$\pi^*(S) \leftarrow \left(\max_{a \in \pi(s)} v_{\pi,t}(s) \right)_{s \in S}$

Si $\max_{s \in S} \left(v_{\pi,t}^*(s) - v_{\pi,t-1}^*(s) \right) \leq \epsilon$ ALORS

RETOURNER $\pi^*(S)$, $v_{\pi,t}^*(S)$

FIN SI

FIN POUR

FIN POUR

RETOURNER $\pi^*(S)$, $v_{\pi,t}^*(S)$

1.2 SARSA

1.3 Q-learning

2 Intégration de pièges