

$$Q_{\pi}^{(t-1)}(s, \pi(s)) = \sum_{s' \in \text{States}} T(s, \pi(s), s') \left[ \text{Reward}(s, \pi(s), s') + \gamma V_{\pi}^{(t-1)}(s') \right]$$

*Remarque : en notant  $S$  le nombre d'états,  $A$  le nombre d'actions par états,  $S'$  le nombre de successeurs et  $T$  le nombre d'itérations, la complexité en temps est alors de  $\mathcal{O}(T_P E S S')$ .*

□ **Q-value optimale** – La  $Q$ -value optimale  $Q_{\text{opt}}(s, a)$  d'un état  $s$  avec l'action  $a$  est définie comme étant la  $Q$ -value maximale atteinte avec n'importe quelle politique. Elle est calculée avec la formule :

$$Q_{\text{opt}}(s, a) = \sum_{s' \in \text{States}} T(s, a, s') \left[ \text{Reward}(s, a, s') + \gamma V_{\text{opt}}(s') \right]$$

□ **Valeur optimale** – La valeur optimale  $V_{\text{opt}}(s)$  d'un état  $s$  est définie comme étant la valeur maximum atteinte par n'importe quelle politique. Elle est calculée avec la formule :

$$V_{\text{opt}}(s) = \max_{a \in \text{Actions}(s)} Q_{\text{opt}}(s, a)$$

□ **Politique optimale** – La politique optimale  $\pi_{\text{opt}}$  est définie comme étant la politique liée aux valeurs optimales. Elle est définie par :

$$\forall s, \quad \pi_{\text{opt}}(s) = \operatorname{argmax}_{a \in \text{Actions}(s)} Q_{\text{opt}}(s, a)$$

□ **Itération sur la valeur** – L'algorithme d'itération sur la valeur (en anglais *value iteration*) vise à trouver la valeur optimale  $V_{\text{opt}}$  ainsi que la politique optimale  $\pi_{\text{opt}}$  en deux temps :

— Initialisation : pour tout état  $s$ , on a

$$V_{\text{opt}}^{(0)}(s) \leftarrow 0$$

— Itération : pour  $t$  allant de 1 à  $T_{\text{VI}}$ , on a

$$\forall s, \quad V_{\text{opt}}^{(t)}(s) \leftarrow \max_{a \in \text{Actions}(s)} Q_{\text{opt}}^{(t-1)}(s, a)$$

avec

$$Q_{\text{opt}}^{(t-1)}(s, a) = \sum_{s' \in \text{States}} T(s, a, s') \left[ \text{Reward}(s, a, s') + \gamma V_{\text{opt}}^{(t-1)}(s') \right]$$

*Remarque : si  $\gamma < 1$  ou si le graphe associé au processus de décision markovien est acyclique, alors l'algorithme d'itération sur la valeur est garanti de converger vers la bonne solution.*