

Porte d'entrée	Porte d'oubli	Porte de sortie	Porte
Écrire ?	Supprimer ?	A quel point révéler ?	Combien écrire ?

□ **LSTM** – Un réseau de long court terme (en anglais *long sort-term memory*, LSTM) est un type de modèle RNN qui empêche le phénomène de *vanishing gradient* en ajoutant des portes d'oubli.

3.4 Reinforcement Learning

Le but du reinforcement learning est pour un agent d'apprendre comment évoluer dans un environnement.

□ **Processus de décision markovien** – Un processus de décision markovien (MDP) est décrite par 5 quantités $(S, \mathcal{A}, \{P_{sa}\}, \gamma, R)$, où :

- S est l'ensemble des états
- \mathcal{A} est l'ensemble des actions
- $\{P_{sa}\}$ sont les probabilités d'états de transition pour $s \in S$ et $a \in \mathcal{A}$
- $\gamma \in [0, 1[$ est le taux d'actualisation (en anglais *discount factor*)
- $R : S \times \mathcal{A} \rightarrow \mathbb{R}$ ou $R : S \rightarrow \mathbb{R}$ est la fonction de récompense que l'algorithme veut maximiser

□ **Politique** – Une politique π est une fonction $\pi : S \rightarrow \mathcal{A}$ qui lie les états aux actions.

Remarque : on dit que l'on effectue une politique donnée π si étant donné un état s , on prend l'action $a = \pi(s)$.

□ **Fonction de valeurs** – Pour une politique donnée π et un état donné s , on définit la fonction de valeurs V^π comme suit :

$$V^\pi(s) = E \left[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots | s_0 = s, \pi \right]$$

□ **Équation de Bellman** – Les équations de Bellman optimales caractérisent la fonction de valeurs V^{π^*} de la politique optimale π^* :

$$V^{\pi^*}(s) = R(s) + \max_{a \in \mathcal{A}} \gamma \sum_{s' \in S} P_{sa}(s') V^{\pi^*}(s')$$

Remarque : on note que la politique optimale π^ pour un état donné s est tel que :*

$$\pi^*(s) = \operatorname{argmax}_{a \in \mathcal{A}} \sum_{s' \in S} P_{sa}(s') V^*(s')$$

□ **Algorithme d'itération sur la valeur** – L'algorithme d'itération sur la valeur est faite de deux étapes :

- On initialise la valeur :

$$V_0(s) = 0$$