

# L'apprentissage par renforcement

Mélanie Raymond

Université du Québec à Montréal

Séminaire d'été 2023

# Table des matières

- 1 Introduction
- 2 Processus décisionnel markovien fini
- 3 Méthodes tabulaires
- 4 Le problème du parieur
- 5 Méthodes d'estimation
- 6 Conclusion

Toutes les informations contenues dans cette présentation sont issues du livre « Reinforcement Learning, An Introduction » de Richard S. Sutton et Andrew G. Barto, 2e édition (2018).

# Introduction

L'idée : on apprend en interagissant avec notre environnement.

On veut apprendre les actions à poser dans une situation donnée de manière à maximiser une récompense numérique.

# Quelques exemples

- Apprendre le chemin le plus court pour sortir d'un labyrinthe.
- Apprendre à jouer (gagner) à un jeu (ex. : échecs, backgammon).
- Apprendre à un robot à ramasser le plus de déchets possible avant de retourner à sa charge.

# Principaux éléments du système

- L'agent : celui qui apprend.  
Ex. : Robot, machine, etc.
- L'environnement : tous les éléments en interaction avec l'agent. L'agent va choisir des actions qui vont modifier l'environnement.

# Sous-éléments du système

- La récompense : nombre que l'environnement envoie à l'agent. L'objectif de l'agent : maximiser ses récompenses.
- La valeur : somme des récompenses attendues en partant d'un état.
- La politique : plan des actions à prendre à chacun des états.
- Le modèle de l'environnement (optionnel) : permet de faire de l'inférence sur la manière dont l'environnement va se comporter.

# Processus décisionnel markovien fini

Dans un processus décisionnel markovien fini, l'ensemble des actions, des états et des récompenses sont de dimension finie.

- $\mathcal{S}$ , l'ensemble des états sans l'état final.
- $\mathcal{S}^+$ , l'ensemble des états incluant l'état final.
- $\mathcal{A}(s)$ , l'ensemble des actions possibles à l'état  $s$ .
- $\mathcal{R} \subset \mathbb{R}$ , l'ensemble des récompenses.



# Processus décisionnel markovien fini

À chaque temps  $t$ , l'agent reçoit une représentation de l'état de l'environnement,  $S_t \in \mathcal{S}$ .

L'agent choisit une action,  $A_t \in \mathcal{A}(s)$ .

Au temps suivant,  $t + 1$ , l'agent reçoit une récompense,  $R_{t+1} \in \mathcal{R}$ , et se retrouve dans un nouvel état,  $S_{t+1}$ .

Le processus peut donc être représenté de la façon suivante :

$$S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, \dots$$

# Dynamiques de l'environnement

Par la propriété markovienne, les variables aléatoires  $R_t$  et  $S_t$  dépendent seulement de l'état et l'action précédents.

On définit la probabilité de les observer au temps  $t$  :

$$p(s', r|s, a) \doteq P(S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a), \\ s, s' \in \mathcal{S}, r \in \mathcal{R}, a \in \mathcal{A}(s).$$

On définit les probabilités de transition d'un état  $s$  à un état  $s'$  :

$$p(s'|s, a) \doteq P(S_t = s' | S_{t-1} = s, A_{t-1} = a) = \sum_{r \in \mathcal{R}} p(s', r|s, a).$$

# Les récompenses et le retour

Le retour est défini comme la somme des récompenses. On le note  $G_t$  :

$$G_t \doteq R_{t+1} + R_{t+2} + \dots + R_T = \sum_{k=t+1}^T R_k,$$

où  $T$  est l'étape finale. Remarquons que :

$$G_t = R_{t+1} + G_{t+1}.$$

# La politique et la valeur

Si l'agent suit la politique  $\pi$  au temps  $t$ , on définit :

$$\pi(a|s) \doteq P(A_t = a | S_t = s).$$

La valeur d'un état  $s$  sous une politique  $\pi$  est :

$$v_\pi(s) \doteq E_\pi(G_t | S_t = s) = E_\pi\left(\sum_{k=t+1}^T R_k | S_t = s\right), \quad \forall s \in \mathcal{S}.$$

La valeur d'une action  $a$  à l'état  $s$  sous une politique  $\pi$  est :

$$q_\pi(s, a) \doteq E_\pi(G_t | S_t = s, A_t = a) = E_\pi\left(\sum_{k=t+1}^T R_k | S_t = s, A_t = a\right).$$

# Les équations de Bellman

On peut définir ces équations récursivement.

$$\begin{aligned}v_{\pi}(s) &= E_{\pi}(G_t | S_t = s) \\&= E_{\pi}(R_{t+1} + G_{t+1} | S_t = s) \\&= \sum_a E_{\pi}(R_{t+1} + G_{t+1} | S_t = s, A_t = a) \cdot P(A_t = a | S_t = s) \\&= \sum_a E_{\pi}(R_{t+1} | S_t = s, A_t = a) \cdot \pi(a | s) \\&\quad + \sum_a E_{\pi}(G_{t+1} | S_t = s, A_t = a) \cdot \pi(a | s) \\&= \sum_a \pi(a | s) \sum_{s'} \sum_r r \cdot p(s', r | s, a) \\&\quad + \sum_a E_{\pi}(G_{t+1} | S_t = s, A_t = a) \cdot \pi(a | s)\end{aligned}$$

# Les équations de Bellman

$$\begin{aligned} &= \sum_a \pi(a|s) \sum_{s'} \sum_r r \cdot p(s', r|s, a) \\ &+ \sum_a \sum_{s'} E_\pi(G_{t+1}|S_{t+1} = s') \cdot p(s'|s, a) \cdot \pi(a|s) \\ &= \sum_a \pi(a|s) \sum_{s'} \sum_r r \cdot p(s', r|s, a) \\ &+ \sum_a \sum_{s'} \sum_r E_\pi(G_{t+1}|S_{t+1} = s') \cdot p(s', r|s, a) \cdot \pi(a|s) \\ &= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r|s, a) \left( r + E_\pi(G_{t+1}|S_{t+1} = s') \right) \\ &= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r|s, a) \left( r + v_\pi(s') \right). \end{aligned}$$

# Les équations de Bellman

$$\begin{aligned} q_{\pi}(s, a) &= E_{\pi}(G_t \mid S_t = s, A_t = a) \\ &= E_{\pi}(R_{t+1} + G_{t+1} \mid S_t = s, A_t = a) \\ &= \sum_r r \sum_{s'} p(s', r \mid s, a) + E_{\pi}(G_{t+1} \mid S_t = s, A_t = a) \\ &= \sum_r r \sum_{s'} p(s', r \mid s, a) + \sum_{s'} E_{\pi}(G_{t+1} \mid S_{t+1} = s') \cdot p(s' \mid s, a) \\ &= \sum_r r \sum_{s'} p(s', r \mid s, a) + \sum_{s'} E_{\pi}(G_{t+1} \mid S_{t+1} = s') \sum_r p(s', r \mid s, a) \\ &= \sum_r r \sum_{s'} p(s', r \mid s, a) \\ &\quad + \sum_{s'} \sum_{a'} E_{\pi}(G_{t+1} \mid S_{t+1} = s', A_{t+1} = a') \cdot \pi(a' \mid s') \sum_r p(s', r \mid s, a) \\ &= \sum_r \sum_{s'} p(s', r \mid s, a) \left( r + \sum_{a'} \pi(a' \mid s') \cdot q_{\pi}(s', a') \right) \end{aligned}$$

Une politique  $\pi$  est meilleure qu'une politique  $\pi'$  si :

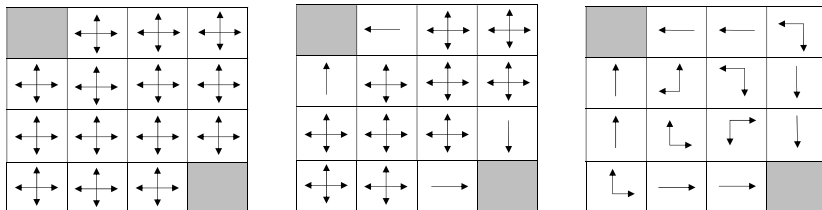
$$\pi \geq \pi' \Leftrightarrow v_{\pi}(s) \geq v_{\pi'}(s), \quad \forall s \in \mathcal{S}.$$

Il existe toujours une politique meilleure ou égale aux autres, on l'appelle la politique optimale. On la note  $\pi_*$ .



# Politique optimale

Quelle politique est la meilleure ?



**Figure** – Exemple tiré de SUTTON, R. S., & BARTO, A. G. (2018). *Reinforcement learning : An introduction* (2nd ed.). The MIT Press.

# Valeurs optimales

Valeur optimale des états :

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

Valeur optimale des couples état-action :

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

$$q_*(s, a) = E_{\pi_*}(G_t | S_t = s, A_t = a)$$

$$= E(R_{t+1} + v_*(S_{t+1}) | S_t = s, A_t = a).$$

# Équations d'optimalité de Bellman

Valeurs optimales définies récursivement :

$$\begin{aligned}v_*(s) &= \max_{a \in \mathcal{A}(s)} q_{\pi_*}(s, a) \\&= \max_a E_{\pi_*}(G_t | S_t = s, A_t = a) \\&= \max_a E_{\pi_*}(R_{t+1} + G_{t+1} | S_t = s, A_t = a) \\&= \max_a E(R_{t+1} + v_*(S_{t+1}) | S_t = s, A_t = a) \\&= \max_a \sum_{s'} \sum_r p(s', r | s, a) (r + v_*(s')).\end{aligned}$$

$$\begin{aligned}q_*(s, a) &= E(R_{t+1} + \max_{a'} q_*(S_{t+1}, a') | S_t = s, A_t = a) \\&= \sum_{s'} \sum_r p(s', r | s, a) (r + \max_{a'} q_*(s', a')).\end{aligned}$$

L'idée pour résoudre les équations d'optimalité de Bellman :

- Évaluer la politique : déterminer la valeur des états sous une certaine politique  $\pi$ .
- Améliorer la politique : vérifier s'il est préférable de suivre cette politique  $\pi$  ou une nouvelle politique  $\pi'$ .

## Algorithme d'itération de la valeur

❶ On initialise  $V(s)$  arbitrairement  $\forall s \in \mathcal{S}$ ,  $V(\text{final}) = 0$ ,  
 $\theta > 0$ ,  $\Delta \leftarrow 0$ .

❷ Faire

• Pour chaque  $s \in \mathcal{S}$ ,

❶  $v \leftarrow V(s)$

❷  $V(s) \leftarrow \max_a \sum_{s',r} p(s', r | s, a) (r + V(s'))$

❸  $\Delta \leftarrow \max(\Delta, |v - V(s)|)$

Tant que  $\Delta > \theta$ .

❸ Retourner  $\pi \approx \pi_*$  telle que

$$\pi(s) = \arg \max_a \sum_{s',r} p(s', r | s, a) (r + V(s'))$$

# Problème du parieur

Une personne lance une pièce de monnaie.

- Face : gagne sa mise.
- Pile : perd sa mise.

Fin du jeu : le parieur a 100\$ ou il a tout perdu.

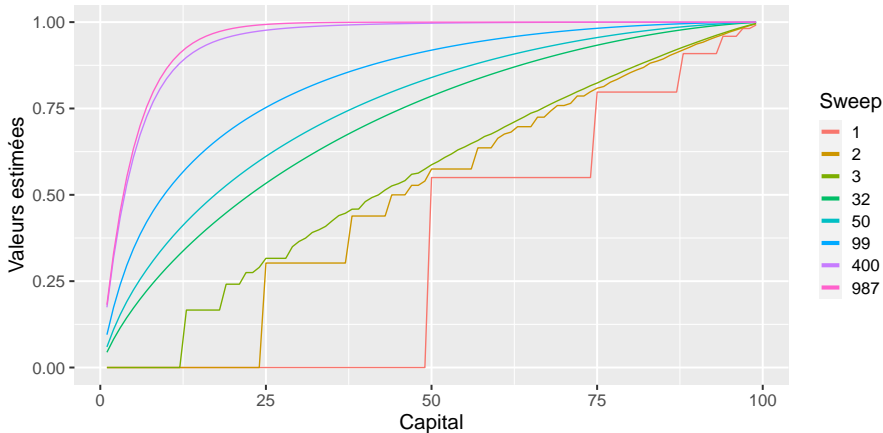
Objectif : déterminer le montant à miser.

- États : le capital détenu par l'agent.
- Actions : la mise faite par l'agent.
- Récompense : 1 si l'agent gagne, 0 sinon.

# Estimation de la valeur des états - 1er cas

## Estimation de la valeur des états

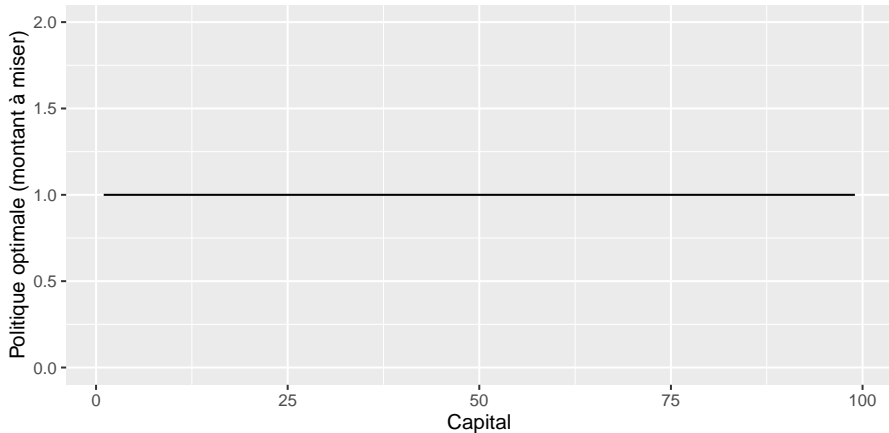
probabilité d'avoir face = 0.55



# Politique optimale

Actions optimales en fonction du capital

probabilité d'avoir face = 0.55

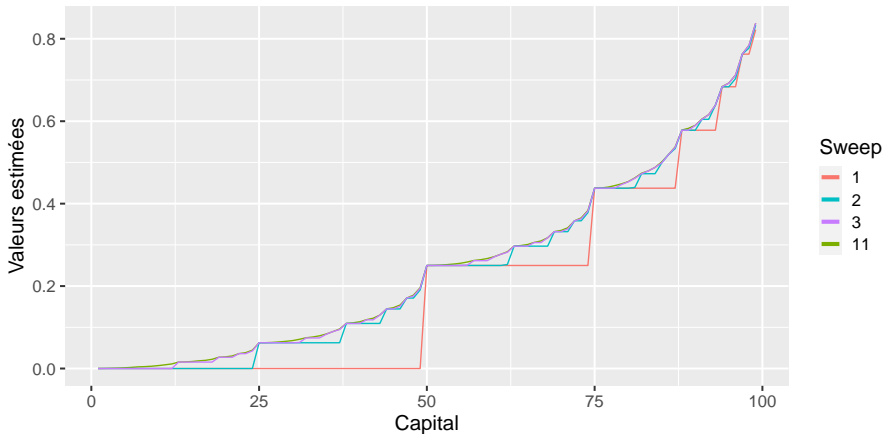




# Estimation de la valeur des états - 2e cas

Estimation de la valeur des états

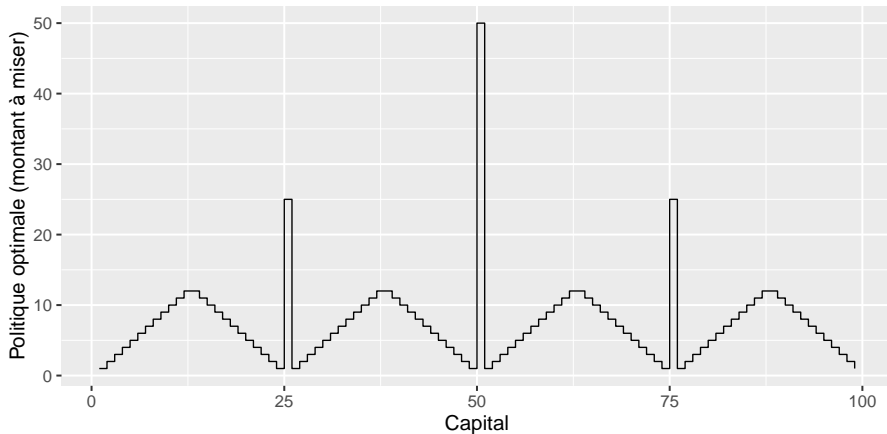
probabilité d'avoir face = 0.25



# Politique optimale - 1re version

Actions optimales en fonction du capital

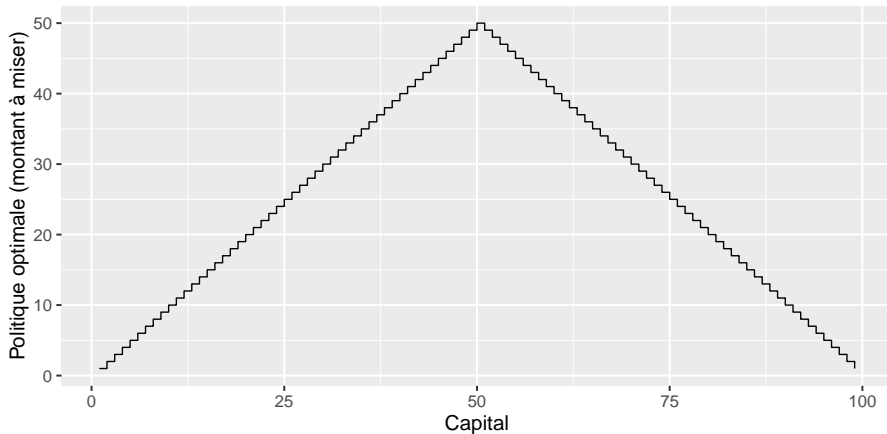
probabilité d'avoir face = 0.25



# Politique optimale - 2e version

Actions optimales en fonction du capital

probabilité d'avoir face = 0.25



# Quelques problèmes...

Pour résoudre les équations d'optimalité de Bellman, il faut :

- 1 Connaître les dynamiques de l'environnement :  $p(s', r|s, a)$ .

Solutions :

- méthodes Monte-Carlo,
- méthodes avec la différence temporelle (TD).

# Quelques problèmes...

Pour résoudre les équations d'optimalité de Bellman, il faut :

- 2 Avoir les ressources informatiques pour faire les calculs.

## Solutions :

- méthodes d'estimation.

# Les méthodes d'estimation

Idée : On veut généraliser les apprentissages.

On veut prendre une bonne décision dans un état jamais visité, en se basant sur ceux déjà rencontrés.

# Estimer la valeur des états

On suppose qu'il existe :

$$\hat{v}(s, \mathbf{w}) \approx v_{\pi}(s),$$

une fonction dérivable en  $\mathbf{w} \in \mathbb{R}^d$ ,  $\forall s \in \mathcal{S}$ , avec  $d \ll |\mathcal{S}|$ .

Par exemple, pour chaque état  $s$ , on suppose que :

$$\hat{v}(s, \mathbf{w}) = \mathbf{w}^T \mathbf{x}(s) = \sum_{i=1}^d w_i x_i(s),$$

où  $\mathbf{x}(s) = (x_1(s), x_2(s), \dots, x_d(s))^T$  est un vecteur d'attributs qui représente l'état  $s$ .



# Estimer la valeur des états

On cherche à minimiser :

$$\overline{VE} = \sum_{s \in \mathcal{S}} \mu(s) [v_{\pi}(s) - \hat{v}(s, \mathbf{w})]^2,$$

où  $\mu(s) \geq 0$ ,  $\sum_s \mu(s) = 1$ , représente la distribution des états.

# Algorithme du gradient stochastique

$$\begin{aligned}\mathbf{w}_{t+1} &= \mathbf{w}_t - \frac{1}{2}\alpha \nabla [v_\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t)]^2 \\ &= \mathbf{w}_t + \alpha [v_\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t)] \nabla \hat{v}(S_t, \mathbf{w}_t),\end{aligned}$$

où  $\alpha > 0$ , contrôle le taux d'apprentissage et

$$\nabla \hat{v}(S_t, \mathbf{w}_t) = \left( \frac{\partial \hat{v}(S_t, \mathbf{w}_t)}{\partial w_{t_1}}, \frac{\partial \hat{v}(S_t, \mathbf{w}_t)}{\partial w_{t_2}}, \dots, \frac{\partial \hat{v}(S_t, \mathbf{w}_t)}{\partial w_{t_d}} \right)^T.$$

# Algorithme du gradient stochastique

Dans le cas linéaire, on a :

$$\hat{v}(s, \mathbf{w}) = \mathbf{w}^T \mathbf{x}(s),$$

donc, on a :

$$\nabla \hat{v}(S_t, \mathbf{w}_t) = \mathbf{x}(s).$$

# Algorithme du gradient stochastique

Comme on ne connaît pas  $v_\pi(S_t)$ , on utilise :

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha [U_t - \hat{v}(S_t, \mathbf{w}_t)] \nabla \hat{v}(S_t, \mathbf{w}_t),$$

où  $U_t$  est un estimateur non biaisé de  $v_\pi(s)$ .

Avec une méthode Monte-Carlo, on utilise  $U_t = G_t$ , car  $E(G_t | S_t = s) = v_\pi(s)$ .

# Estimer la valeur des états sous $\pi$ , $\hat{v} \approx v_\pi$

## Algorithme du gradient Monte Carlo

- ➊ Entrée : une politique  $\pi$  à évaluer
- ➋ Entrée : une fonction dérivable  $\hat{v} : \mathcal{S} \times \mathbb{R}^d \rightarrow \mathbb{R}$
- ➌ Paramètre :  $\alpha > 0$
- ➍ Initialisation :  $\mathbf{w} \in \mathbb{R}^d$ , arbitrairement (p. ex.  $\mathbf{w} = 0$ )
- ➎ Boucle infinie (pour chaque épisode) :
  - ➊ Générer un épisode  $S_0, A_0, R_1, S_1, A_1, \dots, R_T, S_T$  en suivant  $\pi$
  - ➋ Pour  $t = 0, 1, 2, \dots, T - 1$ , faire :
    - $\mathbf{w} \leftarrow \mathbf{w} + \alpha [G_t - \hat{v}(S_t, \mathbf{w})] \nabla \hat{v}(S_t, \mathbf{w})$

# Exploitation vs exploration

On **exploite** les apprentissages déjà faits  
versus  
on **explore** les états jamais visités.

# Estimer la valeur optimale des états

On suit une politique  $\epsilon - greedy$ .

- Exploitation : on choisit les actions maximisant les récompenses  $(1 - \epsilon)\%$  du temps.
- Exploration : on choisit une autre action  $\epsilon\%$  du temps.

On cherche alors  $\hat{q}(s, \mathbf{w}, a) \approx q_*(s, a)$ . On utilise :

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha [U_t - \hat{q}(S_t, A_t, \mathbf{w}_t)] \nabla \hat{q}(S_t, A_t, \mathbf{w}_t).$$

# Estimer $\hat{q}(s, \mathbf{w}, a) \approx q_*(s, a)$

## Algorithme du semi-gradient avec Sarsa

- 1 Entrée : une fonction dérivable  $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$ .
- 2 Paramètre :  $\alpha > 0$ , un petit  $\epsilon > 0$ .
- 3 Initialisation : poids  $\mathbf{w} \in \mathbb{R}^d$  arbitrairement (p. ex.  $\mathbf{w} = 0$ ).



# Estimer $\hat{q}(s, \mathbf{w}, a) \approx q_*(s, a)$

## Algorithme du semi-gradient avec Sarsa - suite

Boucle infinie (pour chaque épisode) :

- ➊ Initialiser  $S, A$  ( $\epsilon - greedy$ ).
- ➋ Pour chaque temps  $t$  de l'épisode :
  - ➊ Faire l'action  $A$ , observer  $R, S'$
  - ➋ Si  $S'$  est l'état final :
$$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w}).$$
Aller à l'épisode suivant.
  - ➌ Sinon :
$$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w}).$$
$$S \leftarrow S'.$$
$$A \leftarrow A'.$$

# Conclusion

- On cherche les actions pour maximiser les récompenses.
- L'agent apprend en interagissant avec son environnement.
- Pour résoudre les problèmes d'apprentissage par renforcement :
  - méthodes tabulaires (programmation dynamique, méthodes Monte Carlo, différence temporelle),
  - méthodes d'estimation (apprentissage supervisé).

Merci !  
Des questions ?