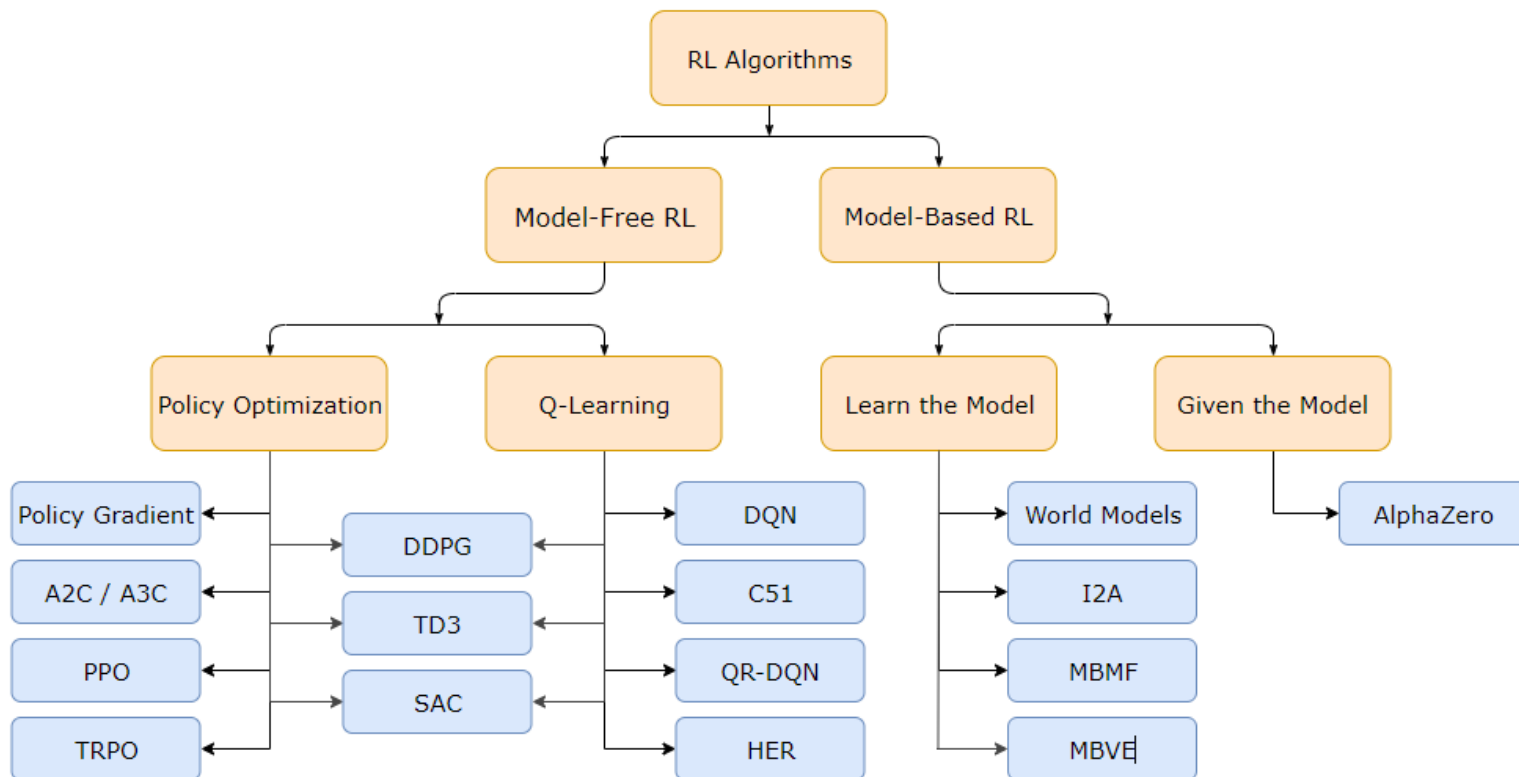




MC_SARSA_ACTEUR CRITIQUE



RL = RÉSOUDRE UN MDP





Model Based

VS



Model Free

La **Transition** est connue pour chaque état ainsi que la **récompense**. (Peut nécessiter une phase d'observation pour la déduire)

Construit ensuite un modèle pour définir la **policy** à partir de ces éléments (Bellman Equation)

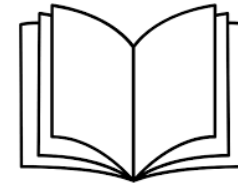
Calcule la policy (l'action) sans avoir besoin de connaître la transition et la récompense



Model Based

Décompose le modèle pour
calculer la transition et la reward.

VS



Model Free

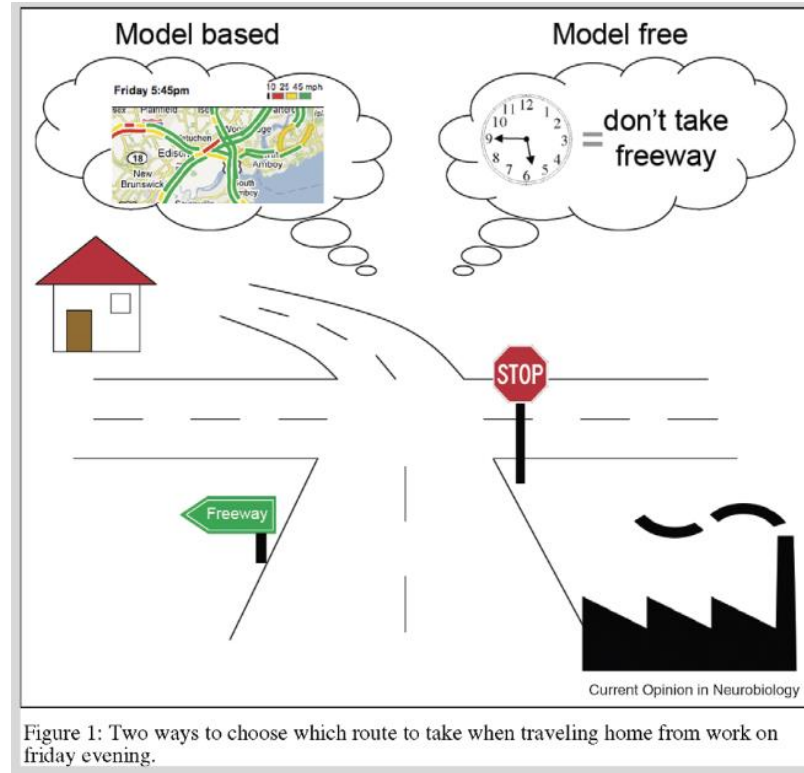
Cherche directement à calculer la
meilleure policy.

Question à se poser :

Après un apprentissage,

Est-ce que l'agent fait des prédictions sur le prochain état et sur la récompense engendrée pour prédire sa prochaine action ?

Si oui, c'est un "**model based**" sinon "**model free**".



OPTIMISATION UNE CAMPAGNE SUR LES RÉSEAUX SOCIAUX

EXEMPLE



Model Based

?

VS



Model Free

?

OPTIMISATION UNE CAMPAGNE SUR LES RÉSEAUX SOCIAUX

EXEMPLE



Model Based

Des études poussées en psychologies permettraient de bien comprendre l'environnement. On utiliserait ces connaissances pour poster des messages.

VS



Model Free

On part d'aucune hypothèse, on essaye et on regarde la récompense. On affine au fur et à mesure?

On part d'échantillons pour déduire une théorie générale (approximation) :



Faits particuliers : *la majorité des parties que j'ai gagné au monopoly était lorsque j'achetais les gares*

Déduction



Théorie : *Il faut acheter les gares pour gagner au monopoly*



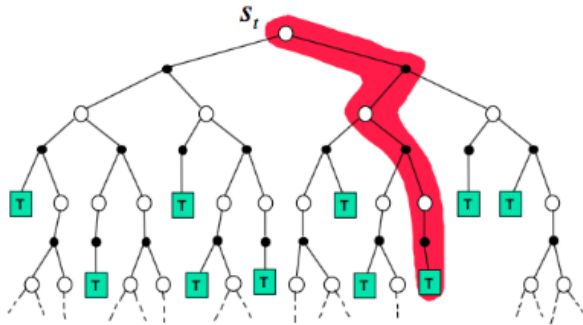
SARSA

*STATE-ACTION-REWARD-STATE-
ACTION*

Model-free

Monte-Carlo

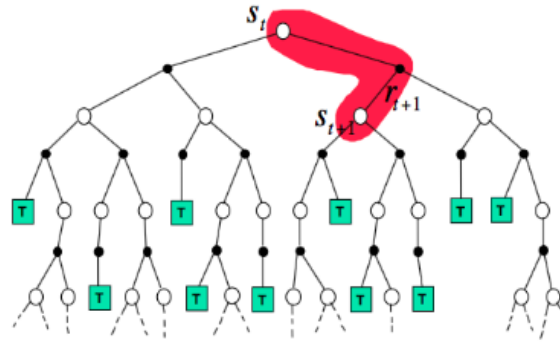
$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$$



runs, expérimentations
complètes

Temporal-Difference

$$V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$

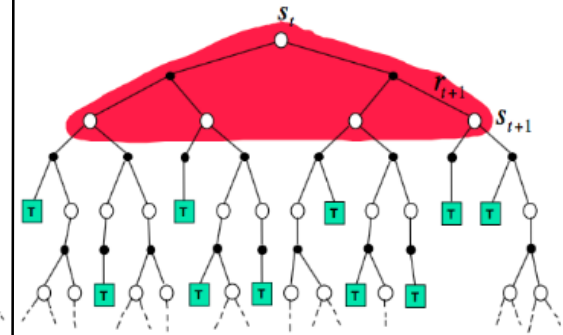


runs, expérimentations
incomplètes

Model-based

Dynamic Programming

$$V(S_t) \leftarrow \mathbb{E}_{\pi} [R_{t+1} + \gamma V(S_{t+1})]$$



Target policy :

C'est la politique que l'on souhaite apprendre.

Behavior policy:

C'est la politique utilisée par l'agent pour sélectionner les actions à faire dans l'environnement

SARSA

Dans l'apprentissage **on-policy**, la fonction $Q(s,a)$ est apprise à partir des actions que nous avons prises en utilisant notre politique actuelle $\pi(a|s)$.

apprendre à quel point il est bon de faire quelque chose en le faisant

Target policy == behavior policy

Q-Learning

Dans l'apprentissage **off-policy**, la fonction $Q(s,a)$ est apprise en prenant différentes actions (par exemple, des actions aléatoires).

apprendre à quel point il est bon de faire quelque chose en faisant autre chose.

Target policy != behavior policy

SARSA VS Q-LEARNING

Questions	SARSA	Q-Learning
Model free / update pendant l'épisode	oui	oui
Policy	On policy	Off policy
Solution trouvée	Proche de l'optimale	optimale
Variance	oui	Forte variance
Comportement	Prudent	+ explorateur
Durée d'entraînement	longue	Très longue

Prend en compte la meilleure action possible (même si elle ne sera pas prise)

Q Learning:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t)]$$

SARSA:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

Prend en compte l'action "réelle" prise

Q – Learning vs SARSA (State Action Reward State Action) Algorithm

Q – Learning (Off policy)

Updated Q Value Current Q Value Target Q Value Current Q Value

$$Q(s, a) = Q(s, a) + \alpha \left[r + \underbrace{\max_{a'} \gamma Q(s', a')}_{\text{Target policy is always Greedy Policy}} - Q(s, a) \right]$$

α = Learning Rate

Target policy is always Greedy Policy

SARSA (State Action Reward State Action) Algorithm (On policy)

Updated Q Value Current Q Value Target Q Value Current Q Value

$$Q(s, a) = Q(s, a) + \alpha \left[r + \underbrace{\gamma Q(s', a')}_{\text{Target Policy is always same as Behaviour Policy}} - Q(s, a) \right]$$

α = Learning Rate

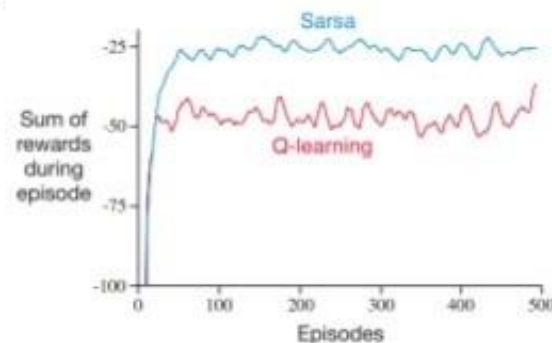
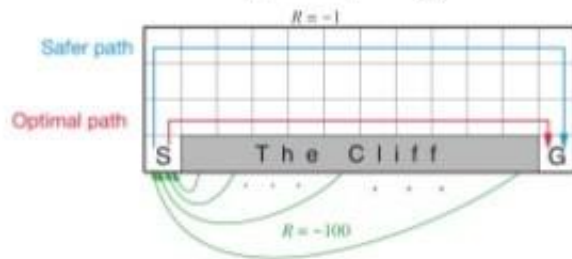
Target Policy is always same as
Behaviour Policy

by Dr. Pankaj Kumar Porwal (BTech - IIT Mumbai, PhD - Cornell University) : Principal, Techno India NJR Institute of Technology, Udaipur

[source](#)

Cliff Walking Example: Sarsa vs. Q-learning

- Q-learning learns optimal policy
- Sarsa learns safe policy
- Q-learning has worse online performance
- Both reach optimal policy with ϵ -decay



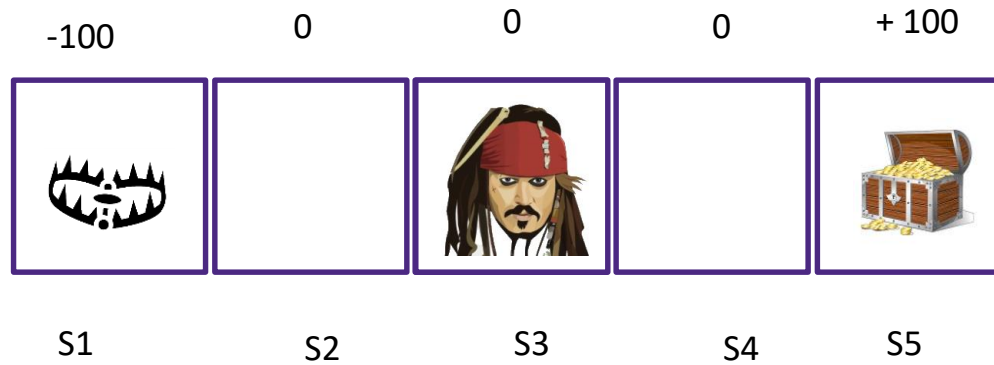
Sutton & Barto book



MONTE CARLO

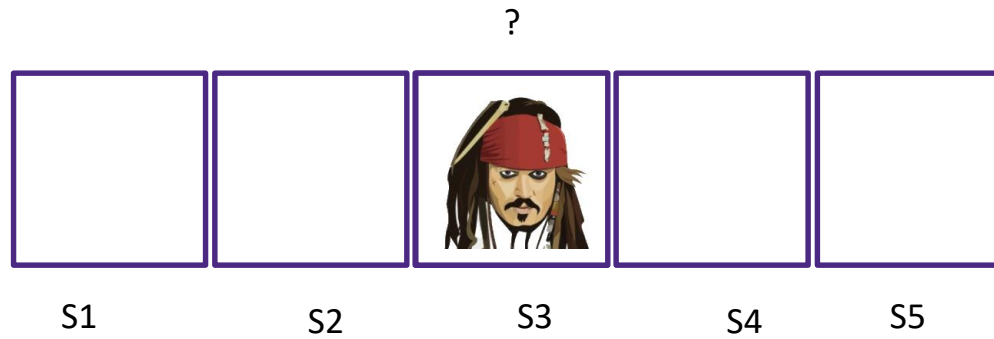


Monte Carlo = Casino = hasard





Comment déduire la valeur pour l'ensemble des states S ?
Avez-vous une idée ?



Calculons la valeur pour le State 3

PRÉDICTION => POUR APPROXIMER $V(S)$

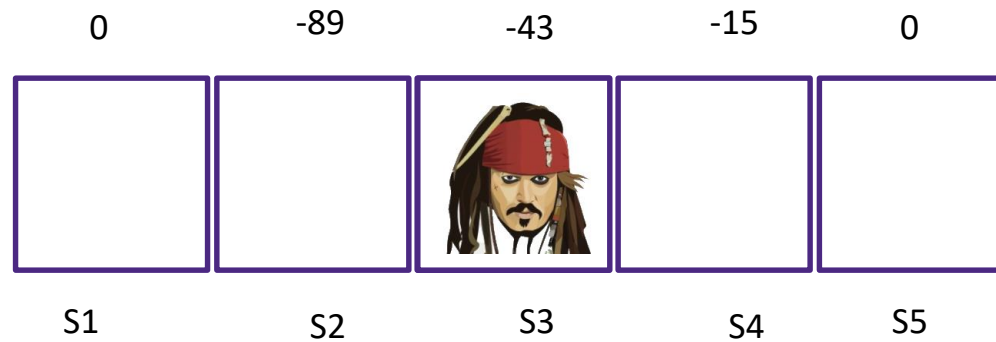
Algorithm 1: Monte-Carlo Prédiction

```
1 Input:  $\pi$ 
2 Init  $V(s) = 0$  pour chaque  $s \in S$ 
3
4 repeat
5   collecte d'expérience avec  $\pi$  (un épisode)
6   foreach  $s$  do
7      $V(s) \leftarrow \text{moyenne}(\text{returns}(s))$ 
8   end
9 until  $V$  ne change presque plus
10
11 return  $V \approx v_\pi$ 
```

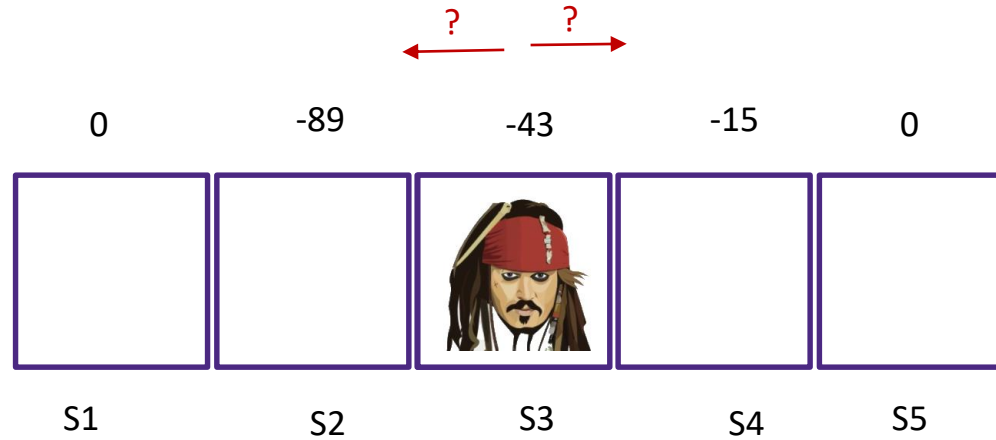
2 contraintes :

Il faut aller au bout d'un "run" pour connaître la finalité du jeu, les environnements doivent pouvoir se terminer (perdu / gagné)

Il faut pouvoir lancer suffisamment de runs pour que ça fonctionne bien)



Nous avons approximé $V(s)$ pouvant nous faire de même avec $Q(s,A)$?



$Q(S,A)$ = motenne (récompense (S,A))

CONTRÔLE => POUR APPROXIMER $Q(S,A)$

Algorithm 1: Monte-Carlo Contrôle

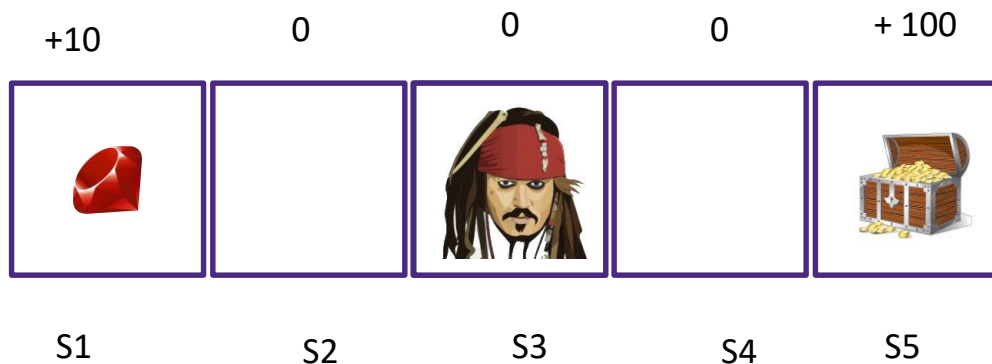
```
1 Init  $Q(s, a) = 0$  pour chaque  $s \in S$  et  $a \in A$ ,  $\pi(a|s)$  arbitrairement
2
3 repeat
4   collecte d'expérience avec  $\pi$  (un épisode)
5   foreach  $s, a$  do
6      $Q(s, a) \leftarrow$  moyenne(returns( $s, a$ ))
7   end
8   foreach  $s$  do
9      $\pi(s) \leftarrow \text{argmax}_a Q(s, a)$ 
10  end
11 until  $Q$  ne change presque plus
12
13 return  $Q \approx q_*$ ,  $\pi \approx \pi_*$ 
```

2 contraintes (identiques) :

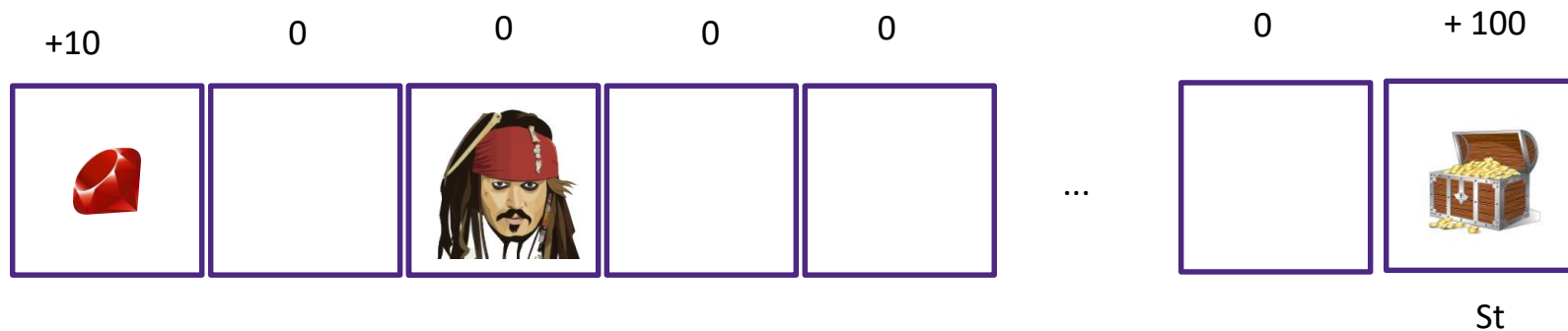
Il faut aller au bout d'un "run" pour connaître la finalité du jeu, les environnements doivent pouvoir se terminer (perdu / gagné)

Il faut pouvoir lancer suffisamment de runs pour que ça fonctionne bien)

Comment s'assurer d'être dans la récompense optimale ?

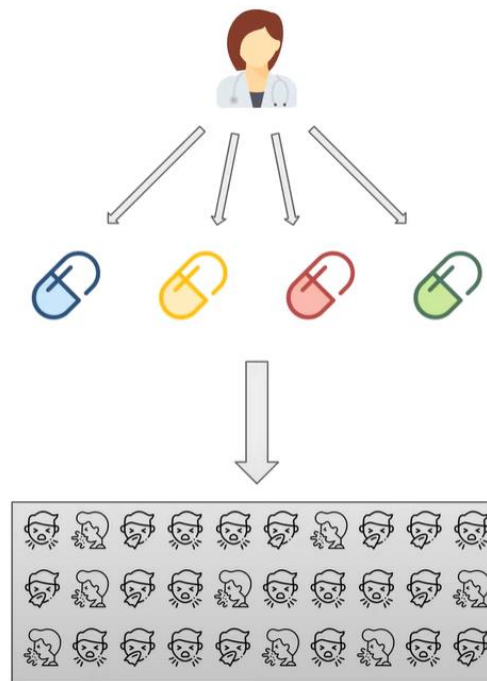
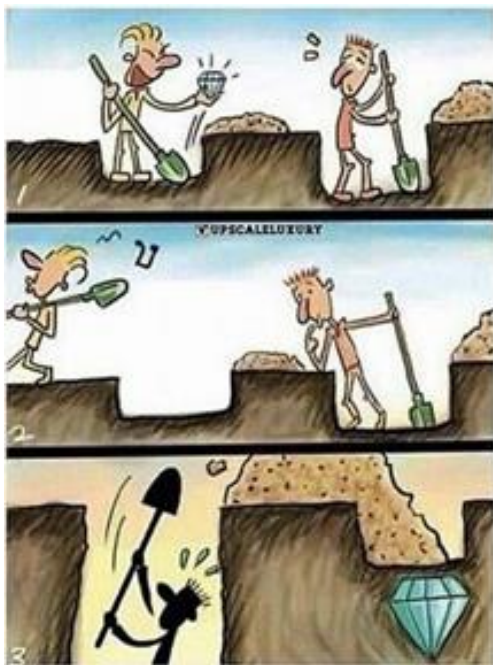


Comment s'assurer d'être dans la récompense optimale ?



Quelle problématique cela soulève ?

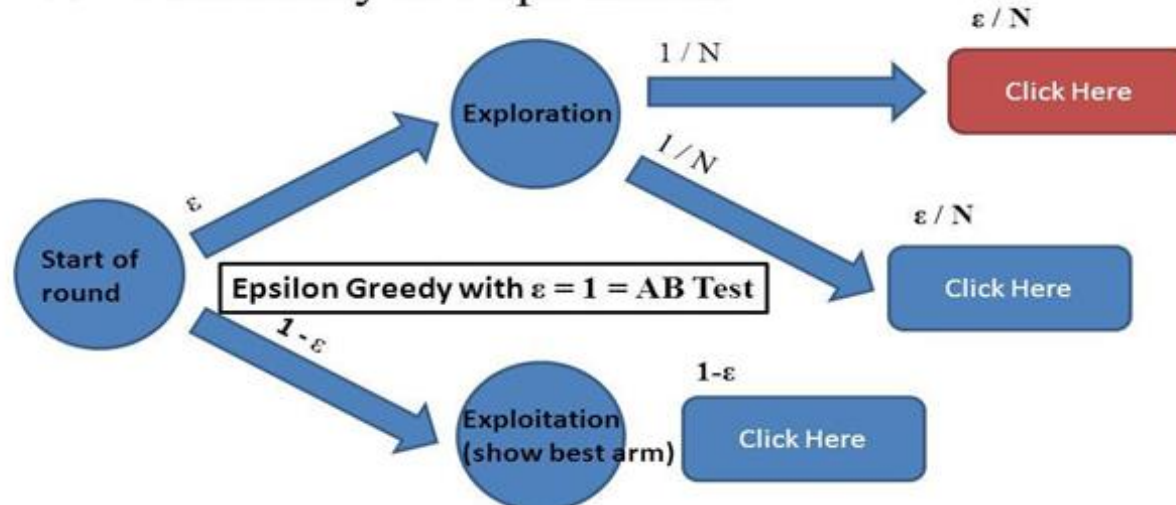
Comment s'assurer d'être dans la récompense optimale ?



ALGORITHME E(EPSILON)-GREEDY

Epsilon Greedy Algorithm

ϵ = Probability of Exploration



On détermine une valeur pour Epsilon, valeur utilisée par la suite pour déterminer si dans le prochain "run" on va explorer ou exploiter. Epsilon est initialisé avec une valeur élevée au début puis va diminuer au fur et à mesure.

- L'algorithme de Monte Carlo appartient à la famille des algorithmes "model free"
- Il permet d'approximer v_π ou q_π
- Pour cela on interagit avec l'environnement (exploration) et on garde en mémoire les récompenses obtenues lorsqu'on est passé par une étape S et par une action + une étape, pour ensuite utiliser la moyenne des récompenses :

$$V(s) \leftarrow \text{moyenne}(\text{returns}(s))$$

$$Q(s, a) \leftarrow \text{moyenne}(\text{returns}(s, a))$$

- Nécessite que les runs, expérimentations soient complètes



ACTEUR CRITIQUE

VS VALUE FUNCTION APPROXIMATION

Value-Based:

Les méthodes vues auparavant comme **Q-Learning** et **SARSA** essayent d'apprendre la fonction de value/action.

Policy-Based:

Cherche à apprendre la police optimale. Cette approche ne nécessite pas de connaître directement la "Value" d'un State ou d'un State + action ($Q(s,a)$) : $\pi:s \rightarrow a$

Dans les algorithmes de Policy Gradient on est face à un problème de maximisation ou l'on va chercher à **maximiser la fonction Reward** (à l'opposé de la fonction de coût utilisée en Deep Learning, à l'aide de l'ascension de gradient).

2 ÉLÉMENTS DIFFÉRENTS : ACTOR-CRITIC

Acteur : l'enfant qui explore le monde

Critique : la mère qui évalue l'enfant (bien / pas bien)

L'acteur va évoluer et explorer. Il va apprendre à réaliser des actions plus intéressantes au fur et à mesure du temps, en fonction des retours de sa mère.

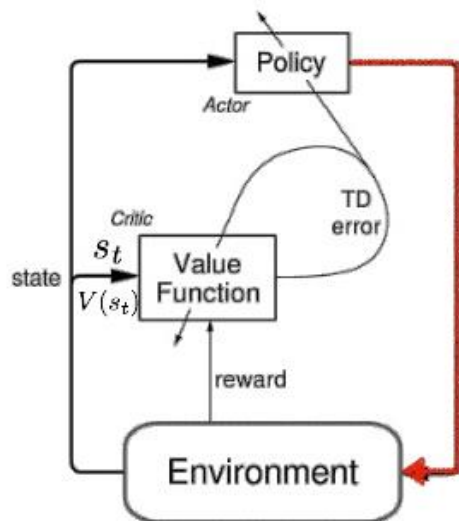
Sa mère va être "surprise" par son comportement et va au fur à mesure améliorer son évaluation sur la situation.

Dans quelques années (fin de l'entraînement) l'enfant aura de "bonnes" actions et la mère saura évaluer parfaitement ses actions.

2 ÉLÉMENTS DIFFÉRENTS : ACTOR-CRITIC

Acteur : met à jour les paramètres de la politique θ , dans la direction suggérée par le critique, $\pi(a|s;\theta)$. Cela peut être une fonction d'approximation qui prédit l'action à faire.

Critique : met à jour les paramètres de la fonction de valeur w et, selon l'algorithme, il peut s'agir d'une valeur d'action $Q(a|s;w)$ ou d'une valeur d'état $V(s;w)$. Cela peut être une fonction d'approximation qui évalue l'action faite.



- Actor: decides which action to take
- Critic: tells the actor how good its action was and how it should adjust

Ce modèle combine les approches "Value based" et "policy based". Séparer la tâche en deux tâches facilite l'apprentissage

Les deux réseaux s'entraînent séparément en optimisant les paramètres à l'aide de l'ascension de gradient. Les mises à jour sont effectuées lors de chaque étape (TD, donc pas besoin d'attendre la fin du run))



EXERCICES