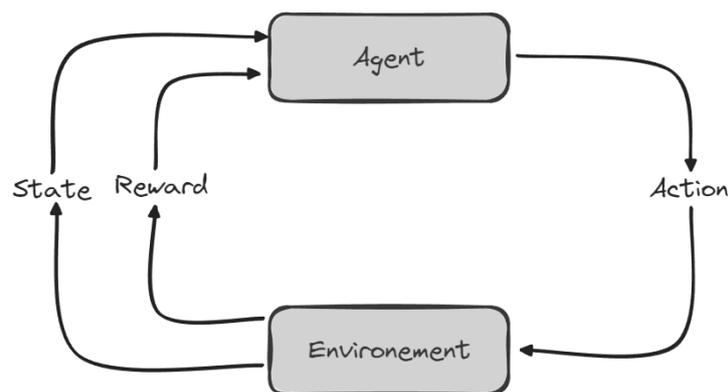


Apprentissage par renforcement

Dans ce chapitre, vous allez apprendre à développer des modèles d'apprentissage par renforcement. En machine learning, l'apprentissage par renforcement est une branche très différente de celles de l'apprentissage supervisé et non supervisé. En effet, dans cette discipline, on ne fournit pas de données (X, y) à notre machine, mais on la programme afin qu'elle puisse générer ses propres expériences et apprendre de ses erreurs.

L'idée centrale est donc de créer un **agent**, libre d'entreprendre des **actions** dans un **environnement**. Ces actions modifient l'**état** (*state*) de l'agent, et ce changement d'état s'accompagne d'une **récompense** (*reward*). Pour l'agent, le but du jeu est de maximiser ses récompenses, ce qui le pousse à apprendre quelles actions effectuer pour obtenir le plus de récompenses.



Cette discipline est très utilisée pour les applications suivantes :

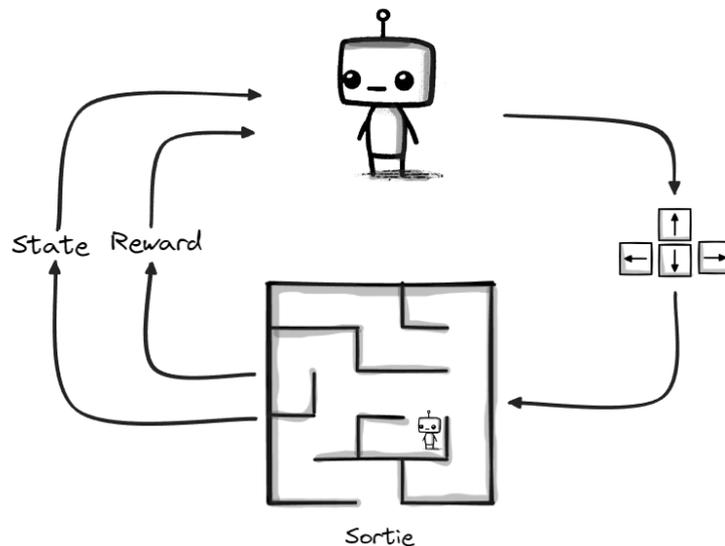
- Robotique
- Véhicules autonomes (voitures, drones)
- Trading
- Algorithmes de prise de décision

Mise en situation

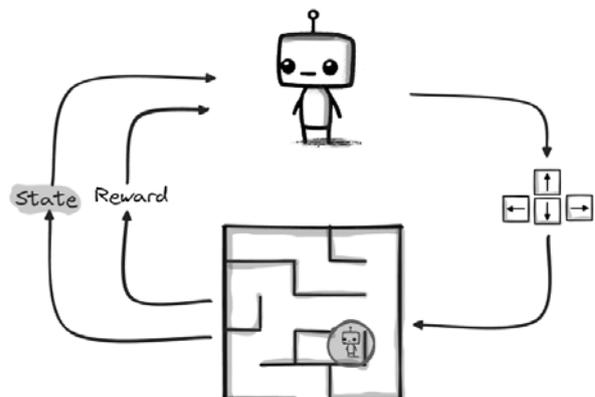
Imaginez : vous programmez un petit robot pour qu'il apprenne à sortir d'un labyrinthe le plus rapidement possible. Pour rendre les choses plus compliquées, le labyrinthe est généré aléatoirement.

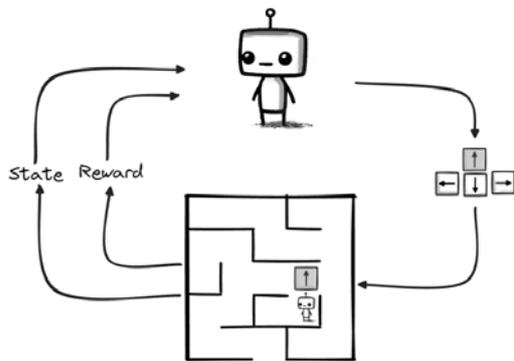
Bien sûr, nous pourrions pour cela construire un algorithme de recherche du chemin le plus court basé sur la théorie des graphes, mais cela reviendrait à coder explicitement le comportement de la machine. Et c'est là tout l'inverse du machine learning. À la place, nous voulons programmer la machine pour qu'elle développe elle-même la stratégie lui permettant de quitter le labyrinthe.

Nous allons donc utiliser une approche d'apprentissage par renforcement. Notre agent sera le petit robot, et le labyrinthe l'environnement dans lequel il évolue.



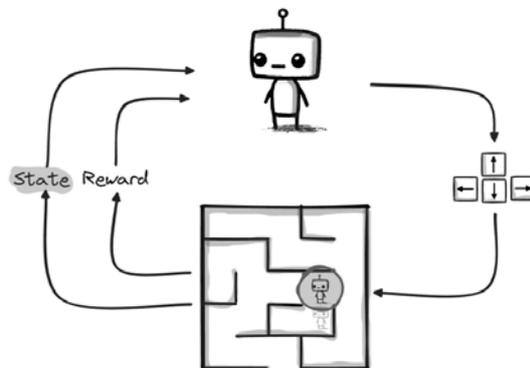
L'état (*State*) correspond à la position de l'agent dans le labyrinthe, et les actions à ses déplacements possibles (haut, bas, droite, gauche)



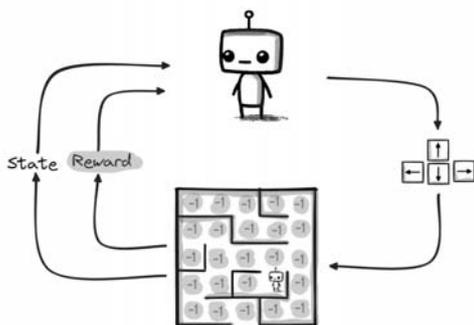


Depuis cet état, l'agent peut choisir une action, par exemple "aller en haut".

En effectuant cette action, l'agent change son état, ce qui s'accompagne d'une récompense. Ici, nous attribuons à la machine un score de -1 point.



Quoi?! Mais l'agent s'est pourtant rapproché de la sortie, pourquoi le sanctionner avec un score négatif? C'est vrai, cependant, rappelez-vous de notre objectif : "trouver la sortie le plus vite possible".



Ainsi, nous infligeons un point négatif chaque fois que la machine effectue un déplacement. Son but étant de maximiser son score final, elle cherchera donc le moyen le plus rapide de rejoindre la sortie, préférant marquer -10 points plutôt que -15. Ne vous en faites pas, la machine ne sera pas triste 😊

En résumé, notre objectif est d'apprendre à la machine, *quelle action effectuer lorsqu'elle se situe dans un état donné*. Autrement dit, nous cherchons à développer une fonction $f(s) = a$ qui prend en entrée l'état s et produit une action a . Cette fonction est couramment nommée **politique d'action**.

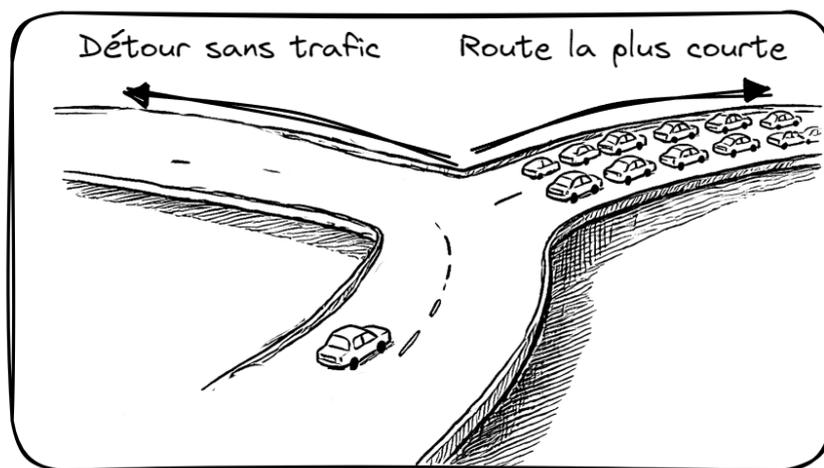
En apprentissage par renforcement, il existe plusieurs méthodes permettant d'apprendre à développer cette politique d'action. Celle que nous allons voir ici est l'approche appelée **Q-Learning**.

Le fonctionnement du Q-Learning

Le Q-Learning consiste à apprendre quelle est la **valeur** d'une action A dans un état S donné.

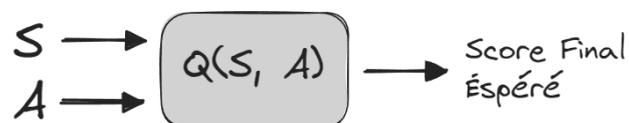
Par exemple, je mets habituellement 30 minutes en voiture pour aller à mon travail. Aujourd'hui, cette route est encombrée par le trafic. J'aperçois cependant qu'une autre route, un peu plus longue, est également possible. Ainsi, deux choix s'offrent à moi : l'action de tourner à gauche ou l'action de tourner à droite.

Pour choisir laquelle emprunter, je vais évaluer la valeur de chacune. À mon avis, quel sera le temps de trajet total si je tourne à gauche? Et si je tourne à droite?

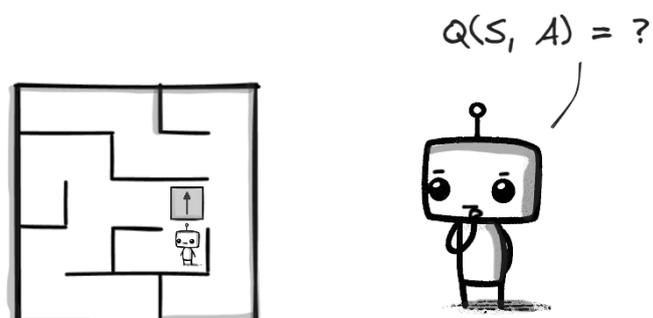


Voilà le but du Q-Learning : apprendre à prédire la valeur de mes actions, en tenant compte de la situation dans laquelle je me trouve.

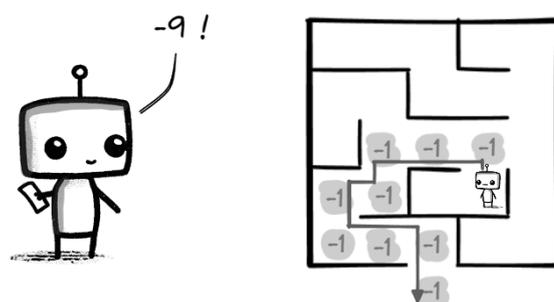
Plus formellement, il s'agit donc d'apprendre une fonction Q qui prend en entrée un état S et une action A, et qui prédit le **score final** que l'on peut espérer obtenir si tout se passe pour le mieux par la suite.



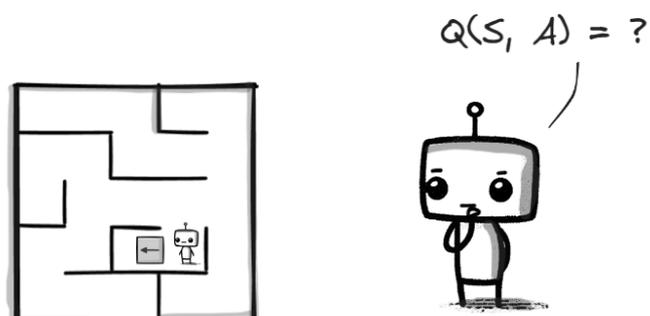
Prenons un exemple : dans la position actuelle, quel est le meilleur score que la machine puisse espérer obtenir si elle choisit l'action "allez en haut" ?



En tenant compte du fait que la machine perd un point pour chaque déplacement, le meilleur score qu'elle puisse espérer obtenir est -9.



Et dans le cas où la machine choisit d'aller à gauche ?



$Q(S, A)$

$S \backslash A$	↑	→	↓	←
	0.34	-0.18	1.69	-0.45
	0.81	-0.74	0.98	0.49
	-0.63	-0.90	0.29	0.73
⋮	⋮	⋮	⋮	⋮

Une formule, connue sous le nom de **l'équation de Bellman**, permet d'exprimer le score total que l'on peut espérer obtenir à l'avenir. Ce score est égal à la récompense obtenue immédiatement dans l'état S , à laquelle on ajoute le même score pour l'état suivant S' , si l'on choisit la meilleure action possible dans ce nouvel état.

$$Q(S, A) = E [R(S') + \gamma \max_{A'} Q(S', A')]$$

Récompense immédiate
obtenue en allant à l'état S'
dévaluation

Récompenses
totales espérées
Espérance
Mathématique
Récompenses
totales espérées dans
le nouvel état S' , en prenant
la meilleure action A'

Pour bien comprendre, reprenons l'analogie avec la voiture.

L'équation de Bellman vous permet d'estimer que la valeur de tourner à gauche vous donne la récompense immédiate d'éviter le trafic, en plus de la valeur que vous allez obtenir dans votre prochain état. Si vous tombez sur des bouchons 1 km plus loin, le fait de tourner à gauche n'était peut-être pas la meilleure option.

À chaque fois que la machine réalise une action, elle peut mettre à jour sa table Q, car elle reçoit une récompense associée à son changement d'état. D'après la formule de Bellman, cette récompense $R(S')$ peut être utilisée pour recalculer la valeur de $Q(S,A)$ en se basant sur sa valeur actuelle et cette nouvelle donnée :

$$\boxed{Q(S, A)} = (1 - \eta) \boxed{Q(S, A)} + \overset{\text{learning rate}}{\eta} \left(\boxed{R(S') + \gamma \max Q(S', A')} \right)$$

Nouvelle valeur Q Valeur Q actuelle Equation de Bellman

Exploration et exploitation

Cependant, dès que la table-Q commence à donner de bons résultats, la machine peut tomber dans le piège de suivre une politique d'action limitée, ignorant au passage les autres possibilités encore inexplorées.

De la même manière qu'après avoir trouvé un bon restaurant dans votre ville, vous choisissez systématiquement d'aller dîner dans ce restaurant car vous pensez qu'il n'y en a pas de meilleurs.

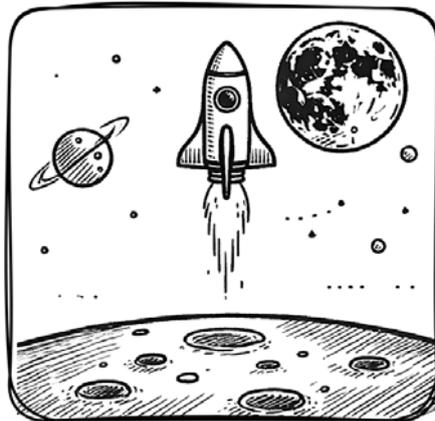
Pour éviter ce phénomène, on définit dans l'algorithme du Q-learning une probabilité **epsilon** qui pousse la machine à continuer **d'explorer** l'environnement des possibles, plutôt que de simplement **exploiter** sa politique d'action déjà apprise.

C'est là le dernier élément qui compose l'algorithme du Q-Learning.

Pour finir ce chapitre, je vous propose donc de passer à l'action avec une implémentation des plus amusantes 😊

Implémentation Python

Dans cette partie, nous allons développer un algorithme de Q-Learning dont le but sera d'apprendre à poser un vaisseau spatial sur la lune sans qu'il ne s'écrase.



Pour cela, nous aurons besoin d'un environnement (la lune) ainsi que d'un agent (le vaisseau spatial). Il existe de nombreuses bibliothèques permettant de générer de tels environnements, l'une des plus populaires étant la bibliothèque **gymnasium** de OpenAI.

Avant de commencer, il vous faudra donc créer un nouvel **environnement virtuel** pour y installer le package *gymnasium* dans sa version 0.29.1 ainsi que son extension *gymnasium[box2d]*.

Si vous rencontrez des difficultés lors de l'installation, les lignes de commandes suivantes devraient permettre de tout installer pour la plupart des utilisateurs :

1. `pip install wheel setuptools pip --upgrade`
2. `pip install swig`
3. `pip install gymnasium==0.29.1`
4. `pip install gymnasium[box2d]`

Une fois que vos packages sont bien installés, le code ci-dessous devrait fonctionner et vous montrer un joli petit vaisseau spatial.

```

1 import gymnasium as gym
2 env = gym.make("LunarLander-v2", render_mode="human")
3 observation, info = env.reset(seed=42)
4
5 for _ in range(1000):
6     action = env.action_space.sample() # this is where you would insert your policy
7     observation, reward, terminated, truncated, info = env.step(action)
8
9     if terminated or truncated:
10        observation, info = env.reset()
11
12 env.close()
```

Dans ce projet, le vaisseau spatial dispose de **quatre actions** possibles :

- Ne rien faire
- Activer le réacteur du bas, afin de ralentir sa descente
- Activer le réacteur de gauche, afin de s'incliner à droite
- Activer le réacteur de droite, afin de s'incliner à gauche

Son **état** lui est décrit par sa position dans l'espace, ainsi que sa vitesse, son angle, et son contact avec le sol.

Pour finir, le système de **récompenses** est défini par plusieurs règles, parmi lesquelles :

- Le vaisseau gagne 100 points s'il réussit à se poser sans dégâts
- Le vaisseau perd 100 points s'il s'écrase au sol (ne se pose pas sur ses pieds)
- Le vaisseau perd des points chaque fois qu'il active un moteur

Nous allons donc commencer par initialiser notre environnement lunaire, ainsi que les différents hyper-paramètres de notre modèle : le *learning rate*, le facteur de dévalorisation, la valeur epsilon, et le nombre d'épisodes que la machine va jouer durant son apprentissage.

```
1 env = gym.make("LunarLander-v2", render_mode="human")
2
3 # hyper-parametres
4 alpha = 0.1
5 gamma = 0.99
6 epsilon = 1.0
7 epsilon_decay = 0.995
8 epsilon_min = 0.01
9 episodes = 1000
10 max_steps = 1000
```

Tout comme notre petit robot qui se déplaçait dans le labyrinthe d'une case à l'autre, nous allons considérer que le vaisseau spatial se déplace sur une grille de coordonnées. C'est pourquoi nous allons découper l'espace des états du vaisseau spatial en plusieurs parties, afin de regrouper tous ces états dans la table-Q (initialement remplie de 0).

```
1 state_bins = [np.linspace(-1.0, 1.0, 10) for _ in range(env.observation_space.shape
   ↪ [0])]
2 n_bins = tuple(len(bins) + 1 for bins in state_bins)
3 q_table = np.zeros(n_bins + (env.action_space.n,))
4
5 def discretize_state(state):
6     return tuple(np.digitize(state[i], state_bins[i]) for i in range(len(state)))
```

Ensuite, nous allons créer une boucle pour entraîner notre vaisseau spatial à se poser sans dégât. Cette boucle comprendra 1000 épisodes, et chaque épisode comprend au maximum 1000 pas.

À chaque pas, le vaisseau spatial "lancer un dé", et en fonction du résultat, il choisira soit **d'explorer** son environnement en choisissant une action au hasard, soit **d'exploiter** sa politique d'action en choisissant l'action qui, selon lui, a la plus grande valeur (cette valeur est définie par la table-Q)

Une fois son action effectuée, le vaisseau recevra une récompense, qu'il utilisera pour mettre à jour la table-Q en appliquant la formule de mise à jour et l'équation de Bellman.

```

1  for episode in range(episodes):
2      state, _ = env.reset(seed=42)
3      state = discretize_state(state)
4      total_reward = 0
5
6      for step in range(max_steps):
7          # Epsilon-greedy policy
8          if random.uniform(0, 1) < epsilon:
9              action = env.action_space.sample() # Exploration
10             else:
11                 action = np.argmax(q_table[state]) # Exploitation
12
13             next_state, reward, terminated, truncated, _ = env.step(action)
14             next_state = discretize_state(next_state)
15
16
17             old_value = q_table[state][action] # valeur Q actuelle
18             next_max = np.max(q_table[next_state]) # valeur Q max de l'état S'
19
20             # mise a jour de la table q avec la formule de Bellman
21             new_value = (1 - alpha) * old_value + alpha * (reward + gamma * next_max)
22             q_table[state][action] = new_value
23
24             state = next_state
25             total_reward += reward
26
27             if terminated or truncated:
28                 break
29
30         # Reduction de Epsilon
31         if epsilon > epsilon_min:
32             epsilon *= epsilon_decay
33
34         print(f"Episode {episode + 1}, Total Reward: {total_reward}")
35
36     env.close()

```



The image shows a Pygame window on the left displaying a spaceship in a dark space with stars and a planet. On the right, an Anaconda Prompt terminal window shows the following output:

```

Episode 956, Total Reward: -52.537119506962654
Episode 957, Total Reward: -274.1815661708278
Episode 958, Total Reward: -129.50405391817526
Episode 959, Total Reward: -123.99879361127579
Episode 960, Total Reward: -81.33572377785472
Episode 961, Total Reward: -30.63434420758678
Episode 962, Total Reward: -171.53399114231794
Episode 963, Total Reward: 195.52354335914873
Episode 964, Total Reward: -49.5905441421678
Episode 965, Total Reward: -131.4994323341388
Episode 966, Total Reward: -78.00692424369342
Episode 967, Total Reward: -116.34641539245031
Episode 968, Total Reward: -82.41889417200659
Episode 969, Total Reward: -339.8177120713906
Episode 970, Total Reward: -138.8188739891168
Episode 971, Total Reward: -106.56025654175505
Episode 972, Total Reward: -37.52645273161238
Episode 973, Total Reward: -125.09824367756268
Episode 974, Total Reward: -66.37850965810063
Episode 975, Total Reward: -166.43691235106746
Episode 976, Total Reward: 196.03895432404119
Episode 977, Total Reward: -166.26412792807378
Episode 978, Total Reward: -40.064445149469485
Episode 979, Total Reward: -79.61670759987423

```

Bravo Commandant! Mission accomplie! 😊