

# Problème du labyrinthe

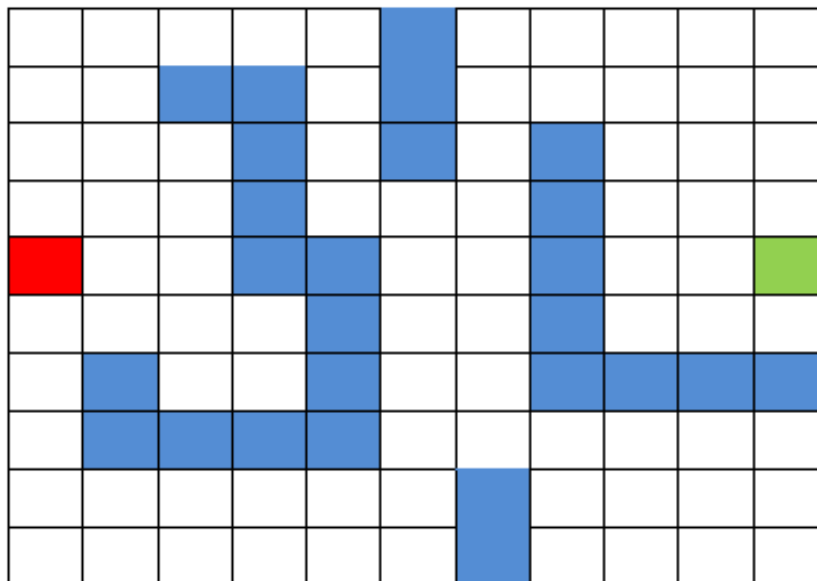
Youssef Abdelhedi et Aicha Ettriki

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Pseudo-code du Q-learning</b>	<b>2</b>
<b>3</b>	<b>Résolution du Problème du Labyrinthe avec le Q-learning</b>	<b>3</b>
3.1	Mise en place de notre grille . . . . .	3
3.2	Choix d'action (epsilon-greedy) . . . . .	4
3.3	Déplacement du robot dans la grille . . . . .	4
3.4	Implémentation de l'algorithme Q-learning . . . . .	5
<b>4</b>	<b>Illustration graphique</b>	<b>6</b>
4.1	Le chemin le plus optimale . . . . .	6
<b>5</b>	<b>Variation des Paramètres <math>\epsilon</math> et <math>\gamma</math></b>	<b>7</b>
5.1	Impact sur le Nombre de Cycles . . . . .	7
5.2	Quelques illustrations . . . . .	8

# 1 Introduction

Le Q-learning est une technique d'apprentissage par renforcement dans le domaine de l'intelligence artificielle. Il est utilisé pour résoudre des problèmes où un agent doit prendre des décisions séquentielles pour maximiser une récompense cumulée à long terme. Le Q-learning fonctionne en apprenant une fonction d'évaluation de l'action, appelée fonction Q, qui estime la récompense attendue pour prendre une certaine action dans un état donné. L'agent utilise cette fonction pour prendre des décisions en choisissant l'action qui maximise la valeur Q. Le Q-learning est basé sur une méthode itérative où l'agent ajuste progressivement ses estimations de la fonction Q en explorant et en exploitant différentes actions dans son environnement. Dans le cadre de ce rapport, nous explorerons le Q-learning, en le mettant en contexte avec le problème du labyrinthe. Imaginez un labyrinthe avec des passages, des murs et une sortie. Dans ce labyrinthe, un agent virtuel doit trouver son chemin de l'entrée à la sortie en prenant une série de décisions à chaque intersection. Chaque fois que l'agent atteint la sortie, il reçoit une récompense positive, tandis que s'il heurte un mur, il reçoit une récompense négative.



Utilisant le Q-learning, l'agent commence par explorer le labyrinthe de manière aléatoire, en prenant des décisions au hasard à chaque intersection. À mesure qu'il explore, il met à jour ses estimations des valeurs Q, qui représentent la récompense attendue pour chaque action dans chaque état du labyrinthe. En apprenant par essais et erreurs, l'agent découvre progressivement les chemins les plus efficaces pour atteindre la sortie.

## 2 Pseudo-code du Q-learning

---

### Algorithm 1 Q-Learning

---

```

1: Initialiser  $Q(s, a)$  arbitrairement pour toutes les paires état-action
2: for chaque épisode do
3:   Initialiser l'état  $s$ 
4:   while  $s$  n'est pas l'état terminal do
5:     Choisir une action  $a$  pour  $s$  en utilisant une politique dérivée de  $Q$  (par exemple,  $\epsilon$ -greedy)
6:     Exécuter l'action  $a$ , observer la récompense  $r$  et l'état suivant  $s'$ 
7:     Mettre à jour la valeur  $Q$  de l'état-action actuel :
8:        $Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$ 
9:        $s \leftarrow s'$ 
10:   end while
11: end for

```

---

Voici une explication des étapes clés de cet algorithme :

1. **Initialisation** : Les valeurs  $Q(s, a)$  sont initialisées arbitrairement pour toutes les paires état-action.
2. **Boucle principale** : Pour chaque épisode :
  - (a) **Initialisation de l'état** : L'agent commence dans un état initial.
  - (b) **Boucle de l'épisode** :
    - i. **Choix de l'action** : L'agent choisit une action à prendre dans cet état.
    - ii. **Exécution de l'action** : L'agent exécute l'action choisie dans l'environnement et observe la récompense résultante ainsi que l'état suivant.
    - iii. **Mise à jour de Q** : Les valeurs  $Q$  sont mises à jour en fonction de la récompense reçue, de la valeur de  $Q$  actuelle pour l'état-action, de la meilleure valeur  $Q$  pour l'état suivant, et des paramètres d'apprentissage  $\alpha$  (taux d'apprentissage) et  $\gamma$  (facteur d'actualisation).
    - iv. **Transition d'état** : L'agent passe à l'état suivant et la boucle continue jusqu'à ce qu'un état terminal soit atteint.
  - (c) **Fin de l'épisode** : Une fois l'épisode terminé, l'algorithme revient à l'étape précédente et commence un nouvel épisode.

### 3 Résolution du Problème du Labyrinthe avec le Q-learning

#### Problème du Labyrinthe et Q-learning :

Le problème du labyrinthe est un cas, où un agent doit naviguer à travers un labyrinthe pour atteindre une destination tout en évitant les obstacles. L'objectif est d'apprendre une politique de décision optimale qui guide l'agent vers la sortie du labyrinthe tout en minimisant le nombre de pas et en évitant les chemins sans issue.

L'algorithme du Q-learning est une approche pour résoudre ce problème. Dans ce contexte, **les états représentent les positions possibles de l'agent dans le labyrinthe, les actions sont les mouvements disponibles** (par exemple, haut, bas, gauche, droite...), et **les récompenses sont définies pour chaque état-action**, avec une récompense positive pour atteindre la sortie et une récompense négative pour heurter un mur.

L'agent utilise l'algorithme du Q-learning pour explorer le labyrinthe, apprendre les valeurs de  $Q$  pour chaque paire état-action, et déterminer la meilleure action à prendre dans chaque situation pour maximiser les récompenses cumulatives. Au fil du temps, l'agent converge vers une politique optimale qui lui permet de naviguer efficacement à travers le labyrinthe et d'atteindre la sortie de manière optimale.

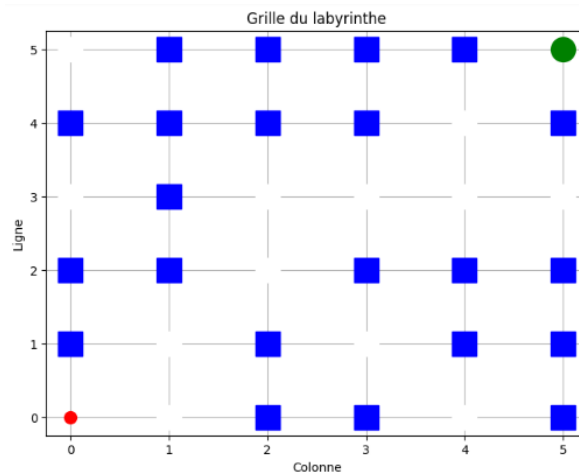
#### 3.1 Mise en place de notre grille

Nous allons tout d'abord, définir notre grid que nous allons utiliser comme support pour le labyrinthe.

```
1 # 0: case vide
2 # -10: mur infranchissable
3 # 100: case de sortie
4
5 grid = [
6     [-10, 0, -10, -10, 0, -10],
7     [-10, 0, -10, 0, -10, -10],
8     [-10, -10, 0, -10, -10, -10],
9     [0, -10, 0, 0, 0, 0],
10    [-10, -10, -10, -10, 0, -10],
11    [0, -10, -10, -10, -10, 100]
12 ]
```

Nous aurons comme résultat ce graph :

Les carrés bleus sont des murs infranchissables , le point rouge c'est le point de départ et le point vert c'est la sortie.



Nous créons ensuite ce que l'on appelle une table ou une matrice de type Q qui suit la forme [état, action] et nous initialisons nos valeurs à zéro. Nous mettons après à jour et stockons nos q-valeurs après chaque épisode. Ce tableau de valeurs devient une table de référence pour notre agent qui sélectionne la meilleure action en fonction des valeurs de cette matrice.

```
1 Q = np.zeros((len(grid), len(grid[0]), 8))
```

En paramètre on a `len(grid)` qui renvoie la longueur de la liste 'grid', qui correspond au nombre de lignes dans la grille, `len(grid[0])` qui renvoie la longueur de la première sous-liste de 'grid', qui correspond au nombre de colonnes dans la grille et le troisième paramètre c'est la dimension supplémentaire pour stocker les valeurs Q pour chaque paire (état, action). Dans ce contexte, le robot peut se déplacer de 8 façons différentes, donc il y a 8 actions possibles pour chaque état.

### 3.2 Choix d'action (epsilon-greedy)

L'agent interagit avec son environnement de deux façons distinctes. Tout d'abord, il explore en agissant de manière aléatoire. C'est ce qu'on appelle l'exploration. Au lieu de sélectionner des actions basées sur la récompense future maximale, il choisit une action au hasard. Cette démarche est cruciale car elle permet à l'agent d'explorer et de découvrir de nouveaux états qui pourraient autrement être négligés lors du processus d'exploitation. Ensuite, l'agent se tourne vers l'exploitation en utilisant la table Q comme référence et en évaluant toutes les actions possibles pour un état donné. Il sélectionne alors l'action qui offre la plus grande valeur. Ce processus d'exploitation utilise les informations disponibles pour prendre des décisions.

Voici un code d'une fonction illustrant ce concept :

```
1 def choose_action(state, epsilon):
2     if random.random() < epsilon:
3         return random.randint(0, 7) #exploration
4     else:
5         return np.argmax(Q[state[0]][state[1]]) #exploitation
```

### 3.3 Déplacement du robot dans la grille

Nous allons créer la fonction `get_new_state(state, action)` qui représente une fonction de transition d'état pour le robot. Il se déplace dans la grille en fonction de l'action qu'il prend à partir de son état actuel. Cette fonction prend en entrée l'état actuel du robot (ses coordonnées dans la grille) ainsi que l'action qu'il souhaite entreprendre. Elle calcule ensuite les nouvelles coordonnées du robot en fonction de cette action, en vérifiant les limites de la grille et en s'assurant qu'il ne se déplace pas vers un mur (représenté par une valeur de -10 dans la grille). Si le nouvel état mène à un mur, la fonction retourne l'état actuel du robot ; sinon, elle retourne les nouvelles coordonnées calculées. En résumé, cette fonction permet de déplacer le robot dans la grille tout en évitant les obstacles.

```

1  def get_new_state(state, action):
2      i, j = state # Les coordonnées actuelles du robot
3      rows = len(grid)
4      cols = len(grid[0])
5
6      # Calcul des nouvelles coordonnées en fonction de l'action
7      if action == 0 and i > 0 and j > 0: # Haut gauche
8          new_i, new_j = i - 1, j - 1
9      elif action == 1 and i > 0: # Haut
10         new_i, new_j = i - 1, j
11     elif action == 2 and i > 0 and j < cols - 1: # Haut droite
12         new_i, new_j = i - 1, j + 1
13     elif action == 3 and j > 0: # Gauche
14         new_i, new_j = i, j - 1
15     elif action == 4 and j < cols - 1: # Droite
16         new_i, new_j = i, j + 1
17     elif action == 5 and i < rows - 1 and j > 0: # Bas gauche
18         new_i, new_j = i + 1, j - 1
19     elif action == 6 and i < rows - 1: # Bas
20         new_i, new_j = i + 1, j
21     elif action == 7 and i < rows - 1 and j < cols - 1: # Bas droite
22         new_i, new_j = i + 1, j + 1
23     else: # Pas dans la grille
24         return state
25
26     # Vérification si le nouvel état n'est pas un mur
27     if grid[new_i][new_j] == -10:
28         # Si le nouvel état est un mur, retourner l'état actuel
29         return state
30     else:
31         # Sinon, retourner le nouvel état
32         return (new_i, new_j)

```

### 3.4 Implémentation de l'algorithme Q-learning

1. La fonction `Q_learning(epsilon, gamma, grid)` prend trois paramètres en entrée : `epsilon`, qui contrôle l'exploration par rapport à l'exploitation, `gamma`, qui représente le facteur de remise pour les récompenses futures, et `grid`, qui est la grille de l'environnement dans lequel l'agent évolue.
2. La boucle `for` principale exécute l'algorithme d'apprentissage par renforcement pendant un certain nombre d'itérations (dans cet exemple, 1000 itérations).
3. À chaque itération, l'état de départ du robot est initialisé à (0, 0), représentant la position de départ dans la grille.
4. La boucle `while` est utilisée pour continuer à exécuter l'algorithme jusqu'à ce que l'agent atteigne la sortie de la grille, représentée par les coordonnées `(len(grid)-1, len(grid[0])-1)`.
5. À chaque itération de la boucle `while`, l'agent choisit une action pour l'état actuel en utilisant notre fonction `choose_action(state, epsilon)`.
6. Ensuite, l'agent obtient le nouvel état après avoir pris cette action en appelant la fonction `get_new_state(state, action)`. Cela met à jour la position de l'agent dans la grille en fonction de l'action choisie.
7. L'agent reçoit une récompense pour avoir atteint le nouvel état, récupérée à partir de la grille à l'aide des coordonnées du nouvel état.
8. Ensuite, l'agent calcule la meilleure récompense future pour le nouvel état en prenant la valeur maximale de la table `Q` associée à cet état.

9. Enfin, l'agent met à jour la valeur Q pour l'état et l'action actuels en utilisant la formule de mise à jour de la valeur Q du Q-learning.
10. L'agent met à jour l'état actuel pour le prochain pas et répète le processus jusqu'à ce qu'il atteigne la sortie de la grille ou que le nombre d'itérations défini soit atteint.

```

1 def Q_learning(epsilon, gamma, grid):
2     for n in range(1000):
3         state = (0, 0)
4         while state != (len(grid)-1, len(grid[0])-1):
5             action = choose_action(state, epsilon)
6             new_state = get_new_state(state, action)
7             reward = grid[new_state[0]][new_state[1]]
8             best_future_reward = np.max(Q[new_state[0]][new_state[1]])
9             Q[state[0]][state[1]][action] = reward + gamma * best_future_reward
10            state = new_state

```

## 4 Illustration graphique

### 4.1 Le chemin le plus optimale

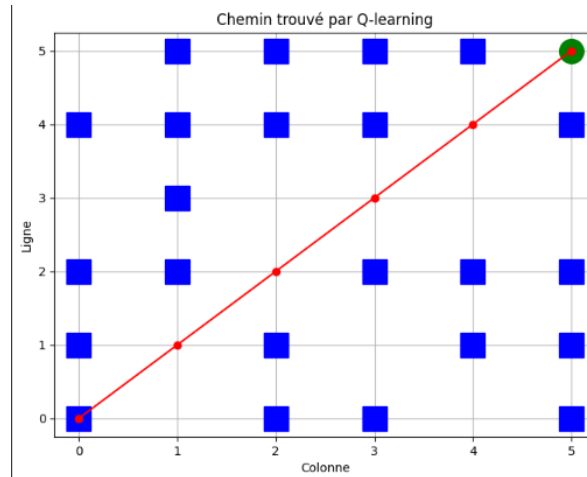
La fonction `Q_learning(0.1,0.9)` est appelée ensuite, la fonction `plot_path(grid, Q)` est appelée pour tracer le chemin optimal trouvé par l'algorithme Q-learning dans la grille. À l'intérieur de la fonction `plot_path`, on définit un état initial (0,0) et un chemin est initialisé avec cet état. En utilisant la politique apprise représentée par la table Q, l'algorithme sélectionne itérativement l'action avec la valeur Q maximale pour chaque état. Cela permet de déterminer le prochain état à visiter dans le chemin. Le chemin complet est stocké dans la liste `path`. Ensuite, le chemin est tracé sur la grille à l'aide de la bibliothèque Matplotlib. Le chemin optimal est tracé en reliant les états successifs du chemin avec des ronds rouges.

```

1 Q_learning(0.1,0.9)
2
3 def plot_path(grid, Q):
4     current_state = (0, 0)
5     path = [current_state]
6     while current_state != (len(grid)-1, len(grid[0])-1):
7         action = np.argmax(Q[current_state[0], current_state[1]])
8         new_state = get_new_state(current_state, action)
9         path.append(new_state)
10        current_state = new_state
11
12    # Tracer le chemin
13    path_x = [step[1] for step in path]
14    path_y = [step[0] for step in path]
15    plt.figure(figsize=(8, 6))
16    for i in range(len(grid)):
17        for j in range(len(grid[0])):
18            if grid[i][j] == -1: # Case blanche
19                plt.plot(j, i, 'ws', markersize=20)
20            elif grid[i][j] == -10: # Mur
21                plt.plot(j, i, 'bs', markersize=20)
22            elif grid[i][j] == 100: # Sortie
23                plt.plot(j, i, 'go', markersize=20)
24    plt.plot(path_x, path_y, marker='o', color='red')
25    plt.title("Chemin trouvé par Q-learning")
26    plt.xlabel('Colonne')
27    plt.ylabel('Ligne')
28    plt.grid(True)
29    plt.show()
30    plot_path(grid, Q)

```

Voici la sortie de ce code :



## 5 Variation des Paramètres $\epsilon$ et $\gamma$

### 5.1 Impact sur le Nombre de Cycles

Nous apportons une modification à notre fonction `Q_learning` afin de calculer le nombre de cycles. À chaque itération du cycle, nous enregistrons les valeurs d' $\epsilon$ ,  $\gamma$  et le nombre de cycles pour chaque combinaison de  $\epsilon$  et  $\gamma$ .

```

1  def q_learning_cycle(epsilon, gamma):
2      for n in range(1000):
3          state = (0, 0)
4          cycles = 0
5          while state != (len(grid)-1, len(grid[0])-1):
6              action = choose_action(state, epsilon)
7              new_state = get_new_state(state, action)
8              reward = grid[new_state[0]][new_state[1]]
9              best_future_reward = np.max(Q[new_state[0]][new_state[1]])
10             Q[state[0]][state[1]][action] = reward + gamma *
                best_future_reward
11             state = new_state
12             cycles += 1
13
14     return {'epsilon': epsilon, 'gamma': gamma, 'cycles': cycles}

```

Ce code explore différentes combinaisons de valeurs pour les paramètres  $\epsilon$  et  $\gamma$ , en enregistrant le nombre de cycles nécessaires.

```

1  epsilon_values = [0.01, 0.1, 0.5, 0.9]
2  gamma_values = [0.01, 0.1, 0.5, 0.9]
3  results = []
4  for epsilon in epsilon_values:
5      for gamma in gamma_values:
6          cycles = q_learning_cycle(epsilon, gamma)
7          results.append(cycles)
8
9
10 print(results)

```

Le code ci-dessus nous donne cette sortie :

```

[{'epsilon': 0.01, 'gamma': 0.01, 'cycles': 5},
{'epsilon': 0.01, 'gamma': 0.1, 'cycles': 5},
{'epsilon': 0.01, 'gamma': 0.5, 'cycles': 5},

```

```
{'epsilon': 0.01, 'gamma': 0.9, 'cycles': 5},
{'epsilon': 0.1, 'gamma': 0.01, 'cycles': 5},
{'epsilon': 0.1, 'gamma': 0.1, 'cycles': 5},
{'epsilon': 0.1, 'gamma': 0.5, 'cycles': 5},
{'epsilon': 0.1, 'gamma': 0.9, 'cycles': 5},
{'epsilon': 0.5, 'gamma': 0.01, 'cycles': 12},
{'epsilon': 0.5, 'gamma': 0.1, 'cycles': 9},
{'epsilon': 0.5, 'gamma': 0.5, 'cycles': 7},
{'epsilon': 0.5, 'gamma': 0.9, 'cycles': 7},
{'epsilon': 0.9, 'gamma': 0.01, 'cycles': 61},
{'epsilon': 0.9, 'gamma': 0.1, 'cycles': 59},
{'epsilon': 0.9, 'gamma': 0.5, 'cycles': 80},
{'epsilon': 0.9, 'gamma': 0.9, 'cycles': 20}]
```

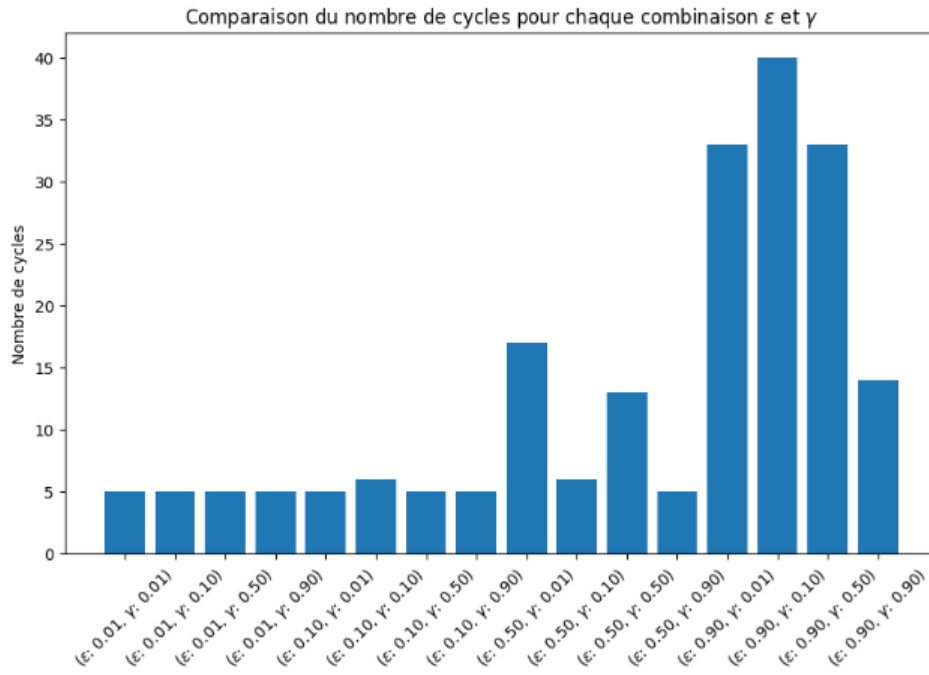
#### Effet de $\epsilon$ :

Pour de petites valeurs de  $\epsilon$  (par exemple,  $\epsilon = 0.01$ ), l'agent privilégie fortement l'exploitation plutôt que l'exploration. Cela signifie qu'il choisit principalement les actions basées sur ce qu'il connaît déjà, ce qui peut conduire à une convergence rapide vers une solution, comme indiqué par le faible nombre de cycles (par exemple, 5 cycles pour  $\epsilon = 0.01$ ). En revanche, pour des valeurs plus élevées de  $\epsilon$  (par exemple,  $\epsilon = 0.9$ ), l'agent privilégie davantage l'exploration, ce qui l'amène à explorer plus activement de nouvelles actions. Cela peut conduire à une convergence plus lente vers une solution optimale, comme indiqué par un nombre plus élevé de cycles nécessaires (par exemple, 61 cycles pour  $\epsilon = 0.9$ ).

#### Effet de $\gamma$ :

Pour de petites valeurs de  $\gamma$  (par exemple,  $\gamma = 0.01$ ), l'agent accorde peu d'importance aux récompenses futures, ce qui peut entraîner une convergence rapide vers une solution locale mais potentiellement sous-optimale. Cependant, cela peut nécessiter moins de cycles pour converger car l'agent est moins attiré à explorer de nouvelles options.

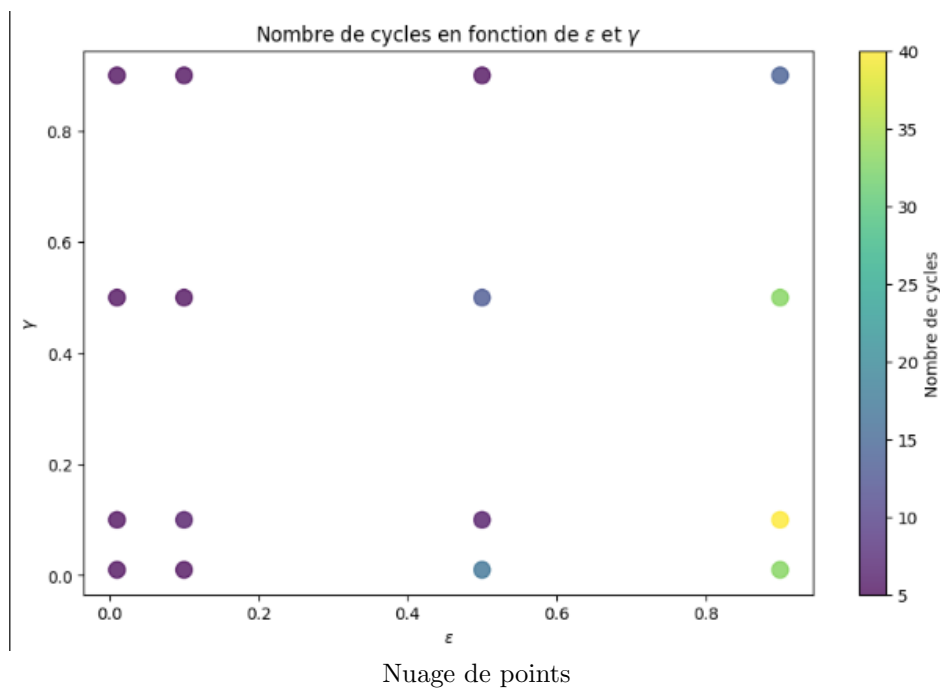
## 5.2 Quelques illustrations



Graphique à barres

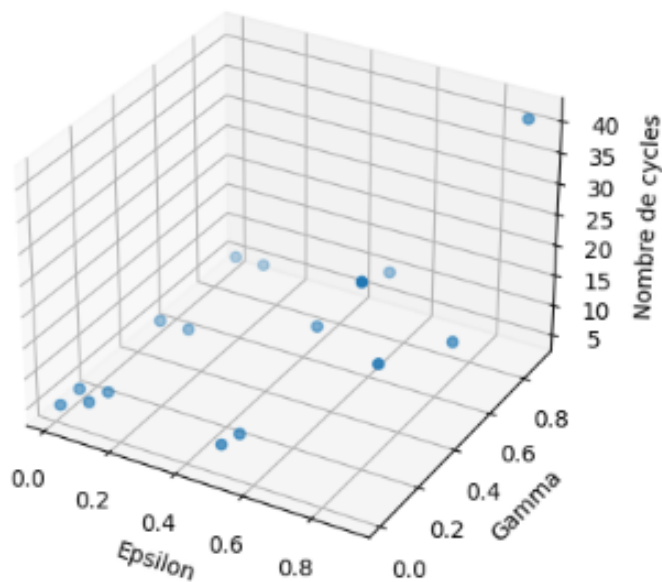
Ce graphique compare le nombre de cycles nécessaires pour différentes combinaisons d' $\epsilon$  et de  $\gamma$  utilisées. On peut remarquer que lorsque  $\gamma$  est augmenté, le nombre de cycles augmente également.





En général, en se déplaçant de gauche à droite sur l'axe des  $x$  ( $\epsilon$  plus élevé), le nombre de cycles augmente pour toutes les valeurs de  $\gamma$ . En effet, avec plus d'exploration ( $\epsilon$  plus élevé), l'agent met plus de temps à converger vers une politique optimale. Pour chaque valeur de  $\epsilon$ , à mesure que  $\gamma$  augmente, le nombre de cycles tend à diminuer.

### Nombre de cycles pour chaque combinaison de Epsilon et Gamma



Visualisation en 3D

Voici ci-dessus un graph en 3D, l'objectif idéal est de trouver un équilibre entre l'exploration et l'exploitation. Le graphique suggère que, pour ce cas spécifique, une combinaison d'épsilon modéré et de gamma plus élevé pourrait conduire au plus faible nombre de cycles nécessaires à l'apprentissage de l'agent