
Résolution du Problème de Labyrinthe

Donia Tekaya

Étudiante en Master 1 BDIA

et

Mohamed Elyes Maalel

Étudiant en Master

6 mai 2024

Table des matières

Introduction Générale	1
1 concept de l'apprentissage par renforcement	2
1.1 Introduction à l'Apprentissage par Renforcement	3
1.1.1 Les termes fondamentaux de l'apprentissage par renforcement : . . .	3
1.1.2 Processus d'Apprentissage :	5
1.1.2.1 Exploration et Exploitation	5
1.1.3 Description de l'Algorithme Q-learning	6
2 Résolution du problème du Labyrinthe	12
2.1 Résolution du problème :Labyrinthe :	13
2.1.1 Implémentation de l'algorithme de Qlearning	14
2.1.2 Interpretation des Comparaison des resultats :	25
2.1.3 Analyse du Graphique	31
2.1.3.1 Variabilité des Performances :	32
2.1.3.2 Comparaison des Tailles de Labyrinthe :	32
2.1.4 Visualisation du parcours du robot dans labyrinthe	33

Table des figures

1.1	Agent	3
1.2	l'algorithme de Qlearning	7
2.1	chemin du robot	20
2.2	resultat1	23
2.3	resultat2	23
2.4	resultat3	24
2.5	resultat4	24
2.6	resultat	31
2.7	le output	37

Introduction Générale

Le projet aborde le problème du labyrinthe dans le but de découvrir un chemin optimal permettant de surmonter les obstacles, en utilisant l'apprentissage par renforcement. Ce problème est pertinent en intelligence artificielle car il illustre de manière simplifiée les défis de navigation et de prise de décision rencontrés dans des environnements réels. Nous nous focalisons sur l'algorithme de Q-learning, en adoptant une stratégie d'action ϵ -greedy, où le facteur de dépréciation est fixé à 0,9. Les récompenses sont attribuées en fonction des résultats du chemin suivi : une case sortie vaut 100 points, un mur vaut -10 points, et les autres cases valent -1 point.

L'objectif de ce projet est d'appliquer cet algorithme au problème du labyrinthe, d'illustrer graphiquement le chemin trouvé, et d'explorer l'influence des paramètres ϵ (coefficient d'apprentissage) et γ (facteur de dépréciation) sur les performances. Les variations des paramètres seront analysées pour évaluer l'impact sur le nombre de cycles nécessaires pour trouver la sortie. De plus, nous étudierons les performances de l'algorithme lorsque la taille du labyrinthe est modifiée (grilles 20x20 et 30x30). Les résultats obtenus seront comparés et interprétés graphiquement pour offrir une meilleure compréhension des stratégies optimales.

CHAPITRE 1

concept de l'apprentissage par renforcement

1.1 Introduction à l'Apprentissage par Renforcement

L'apprentissage par renforcement (Reinforcement Learning, RL) est une branche de l'intelligence artificielle qui traite de la manière dont les agents logiciels devraient prendre des décisions dans un environnement afin de maximiser une notion de récompense cumulative. Contrairement à d'autres types d'apprentissage automatique, l'apprentissage par renforcement est caractérisé par une interaction entre un agent et son environnement, dans lequel l'agent apprend à atteindre un objectif en essayant différentes actions et en observant les résultats de ces actions.

1.1.1 Les termes fondamentaux de l'apprentissage par renforcement :

Dans le cadre de l'apprentissage par renforcement, les termes suivants sont fondamentaux pour comprendre comment les agents interagissent avec leur environnement pour apprendre des politiques efficaces :

Agent :



FIGURE 1.1 – Agent

L'agent est l'entité qui prend des décisions dans un environnement. Il effectue des actions en se basant sur son observation actuelle de l'état de l'environnement ou sur une politique prédéfinie. L'objectif de l'agent est de maximiser la

somme des récompenses qu'il reçoit de l'environnement. Dans un jeu de plateau, par exemple, l'agent pourrait être le joueur qui décide des mouvements à faire.

Environnement : L'environnement est le monde extérieur avec lequel l'agent interagit et sur lequel il n'a pas un contrôle total. Il peut être tout système ou cadre qui répond aux actions de l'agent par des changements d'états et des récompenses. L'environnement détermine les états et fournit des récompenses ou des pénalités en réponse aux actions de l'agent.

État : Un état est une description de la situation actuelle de l'environnement dans laquelle se trouve l'agent. Chaque action que l'agent prend résulte dans un nouvel état, qui peut être le même ou différent. Les états fournissent le contexte nécessaire à l'agent pour prendre des décisions. Par exemple, dans un jeu d'échecs, un état serait la configuration actuelle des pièces sur l'échiquier.

Action : Une action est une intervention spécifique que l'agent effectue dans un état donné. L'ensemble des actions disponibles peut varier d'un état à l'autre. Les actions de l'agent influencent l'état suivant de l'environnement et les récompenses ultérieures que l'agent reçoit.

Récompenses à l'instant t : Une récompense est un feedback immédiat fourni par l'environnement en réponse à une action prise par l'agent. La récompense à l'instant t , notée souvent R_t , est un indicateur de la valeur de l'action prise : elle peut être positive (encourageant l'action) ou négative (décourageant l'action). L'agent cherche à maximiser la somme totale des récompenses obtenues au cours du temps, ce qui peut impliquer des compromis entre les bénéfices immédiats et les récompenses futures.

Ces concepts forment la base de l'interaction dans l'apprentissage par renforcement, où l'agent apprend à naviguer dans l'environnement de manière à maximiser ses récompenses cumulatives.

1.1.2 Processus d'Apprentissage :

Le processus d'apprentissage dans le cadre de l'apprentissage par renforcement (RL) implique une série de décisions que l'agent doit prendre pour apprendre de son environnement de manière efficace. Ce processus est essentiellement itératif, où chaque étape ou décision contribue à l'amélioration progressive des connaissances de l'agent sur l'environnement. Un aspect central de ce processus est la dualité entre exploration et exploitation, deux stratégies qui doivent être équilibrées pour permettre un apprentissage optimal. Voici une explication plus détaillée de ces concepts :

1.1.2.1 Exploration et Exploitation

Exploration :

But : Découvrir de nouvelles connaissances sur l'environnement.

Comment : L'agent prend des actions au hasard, indépendamment de la politique actuelle. Cela permet à l'agent de visiter et de tester des états qui ne sont pas fréquemment explorés dans le cadre de la politique actuelle.

Pourquoi : Sans exploration suffisante, l'agent risque de ne jamais découvrir des actions qui pourraient être plus avantageuses que celles déjà connues. L'explora-

tion est cruciale, particulièrement dans les étapes initiales de l'apprentissage ou dans des environnements très dynamiques où les conditions changent et où de nouvelles stratégies peuvent devenir nécessaires.

Risque : Trop d'exploration peut conduire à des décisions inefficaces, car l'agent peut continuer à choisir des actions suboptimales au détriment de l'accumulation de récompenses immédiates.

Exploitation :

But : Maximiser les récompenses basées sur les connaissances actuelles.

Comment : L'agent choisit les actions qui, selon sa connaissance actuelle (souvent représentée par une table de valeurs Q ou une politique), semblent offrir la meilleure récompense.

Pourquoi : L'exploitation permet à l'agent d'accumuler des récompenses de manière efficace en utilisant les meilleures stratégies qu'il a apprises jusqu'à présent.

Risque : Si l'agent exploite exclusivement sans continuer à explorer, il peut rester piégé dans des comportements localement optimaux et ne jamais découvrir d'autres actions qui pourraient être plus bénéfiques à long terme.

1.1.3 Description de l'Algorithme Q-learning

Le Q-learning cherche à apprendre une fonction de valeur d'action, $Q(s, a)$, qui donne une estimation de la valeur totale des récompenses futures obtenues en prenant l'action a dans l'état s et en suivant ensuite la meilleure politique. Le but est de construire une table Q qui guide directement les décisions optimales.

Étapes de l'Algorithme Q-learning

```
procedure QLEARNING
   $\forall s \forall a Q(s, a) \leftarrow 0$ 
  for  $n \leftarrow 1, nbCycles$  do                                 $\triangleright$  nbCycles d'apprentissage
     $\lambda \leftarrow 1; \epsilon \leftarrow 1;$ 
     $etatCourant \leftarrow etatInitial$ 
    for  $i \leftarrow 1, nbMaxActions$  do                         $\triangleright$  max supposé d'actions à tester
       $s \leftarrow etatCourant$ 
       $nb \leftarrow random(0,1)$ 
      if  $(nb < \epsilon)$  then
         $a \leftarrow randomAction(s)$   $\triangleright$  choix aléatoire d'une action à partir de s
      else
         $a \leftarrow argMax_{a'}(Q(s, a'))$   $\triangleright$  choix de l'action de s avec Q maximum
      end if
       $s' \leftarrow a(s)$   $\triangleright$  s' est l'état d'arrivé après exécution de a en s
       $Q(s, a) \leftarrow \lambda \times (r + \gamma \times max_{a'}(Q(s', a'))) + (1 - \lambda) \times Q(s, a)$ 
       $\lambda \leftarrow 0.99 \times \lambda$   $\triangleright$  décrémenter les coefficients
       $\epsilon \leftarrow 0.99 \times \epsilon$ 
      if  $(s' = etatFinal)$  then
        Sortie boucle i
      end if
    end for
  end for
```

FIGURE1.2 – l'algorithme de Qlearning

Initialisation :

Initialisez la table Q à des valeurs initiales pour chaque paire état-action. Ces valeurs peuvent être mises à zéro ou à toute petite valeur aléatoire, selon la méthode préférée.

Boucles d'Apprentissage :

Pour chaque épisode d'apprentissage (un épisode commence généralement dans un état initial et se termine lorsque l'état terminal est atteint ou après un nombre maximum d'étapes) :

Choix d'une Action (Stratégie $\epsilon - greedy$) :

Sélectionnez une action a pour l'état courant s en utilisant une politique dérivée

de Q (par exemple, $\epsilon - greedy$). Avec une probabilité ϵ , choisissez une action aléatoire (exploration), et avec une probabilité $1 - \epsilon$, choisissez l'action qui maximise la valeur de Q pour l'état courant (exploitation).

Exécution de l'Action et Observation de la Récompense :

Exécutez l'action a et observez la récompense r et le nouvel état s .

Mise à Jour de la Table Q :

Mettez à jour la valeur de Q pour l'état-action précédent (s, a) en utilisant la règle de mise à jour suivante :

$$Q(s, a) = \lambda \times (r + \gamma \times \max_{a'} Q(s', a')) + (1 - \lambda) \times Q(s, a)$$

La formule décrit comment mettre à jour la valeur $Q(s, a)$, qui est l'estimation de la valeur de prendre une action a dans un état s , en utilisant la récompense obtenue r et la meilleure estimation de la valeur future maximale $\max_{a'} Q(s', a')$ dans le nouvel état s' . Les composants de la formule sont les suivants :

$Q(s, a)$: La valeur actuelle de l'action a dans l'état s .

λ (**lambda**) : Le taux d'apprentissage. Ce coefficient détermine à quel point la nouvelle information (la récompense plus la récompense future estimée) affectera la valeur Q actuelle. Une valeur proche de 0 fait que l'apprentissage est très lent (la valeur actuelle change peu), tandis qu'une valeur proche de 1 fait que l'apprentissage est rapide (la valeur actuelle change beaucoup à chaque mise à jour).

r : La récompense obtenue après avoir pris l'action a dans l'état s .

γ (**gamma**) : L'ajustement du paramètre gamma se fait lors de **l'exploitation** dans les algorithmes d'apprentissage par renforcement. Le facteur de dépréciation. Il pondère l'importance des récompenses futures. Un γ proche de 0 signifie que l'agent ne valorise que les récompenses immédiates, tandis qu'un γ proche de 1 signifie que l'agent valorise davantage les récompenses futures.

$\max_{a'} Q(s', a')$: La meilleure valeur Q estimée pour le prochain état s' , maximisée sur toutes les actions possibles a' . Cela représente l'estimation de la meilleure récompense future que l'agent peut obtenir à partir du prochain état.

Répétition des Étapes Boucle d'Apprentissage : Répétition des étapes 3 à 5 : Dans le contexte du Q-learning, une fois que vous avez établi la manière dont l'agent choisit une action, observe les conséquences, et met à jour ses estimations de valeurs Q , ces étapes sont répétées continuellement à chaque pas de chaque épisode. Un "pas" ici correspond à un cycle unique de choix d'action, observation de la réaction de l'environnement, et mise à jour de la table Q . L'ensemble de ces pas forme un "épisode", qui commence généralement dans un état initial et se termine quand un état terminal est atteint ou après un nombre prédéfini de pas.

Fin d'un Épisode :

Un épisode peut se terminer pour plusieurs raisons, comme atteindre un objectif, échouer de manière irréversible, ou après un certain nombre de pas pour éviter des épisodes excessivement longs. À la fin de chaque épisode, l'agent a potentiellement acquis de nouvelles connaissances sur les conséquences des actions dans divers états, ce qui affinera ses décisions futures.

Réduction de ϵ

Stratégie $\epsilon - greedy$:

La stratégie $\epsilon - greedy$ joue un rôle crucial dans l'équilibre entre exploration (choisir des actions aléatoires) et exploitation (choisir les meilleures actions selon la table Q actuelle).

Au début de l'apprentissage, lorsque l'agent ne possède que peu ou pas de connaissance sur l'environnement, il est important de favoriser l'exploration pour découvrir de nouvelles actions et apprendre des récompenses associées à ces actions.

Par conséquent, un petit epsilon est choisi, ce qui signifie que l'agent a une faible probabilité de choisir une action de manière aléatoire plutôt que d'exploiter les connaissances acquises ($1 - \epsilon$).

Au fur et à mesure que l'agent acquiert de l'expérience et que sa connaissance de l'environnement s'améliore, il devient de moins en moins bénéfique de prendre des actions aléatoires. À ce stade, il est plus efficace pour l'agent d'exploiter les connaissances qu'il a acquises pour maximiser les récompenses à long terme. Par conséquent, l'epsilon peut être progressivement augmenté au fil du temps, ce qui réduit la probabilité d'exploration aléatoire et favorise davantage l'exploitation des connaissances actuelles.

Cependant, il est important de noter que même lorsque l'epsilon est augmenté, il ne devrait pas nécessairement être fixé à zéro. En laissant une petite probabilité d'exploration aléatoire même après que l'agent ait acquis une connaissance approfondie de l'environnement, on évite le risque de rester coincé dans des comportements sous-optimaux et on laisse une certaine marge de ma-

nœuvre pour découvrir de nouvelles stratégies potentiellement plus efficaces. Cette approche permet à l'agent de maintenir un équilibre entre exploration et exploitation, même à des stades avancés de l'apprentissage.

Méthode de réduction :

ϵ peut être diminué progressivement, souvent en fonction du nombre d'épisodes ou de pas, afin d'augmenter la fréquence d'exploitation des connaissances acquises. Cette diminution peut être linéaire ou exponentielle, selon le comportement désiré.

CHAPITRE 2

Résolution du problème du Labyrinthe

2.1 Résolution du problème :Labyrinthe :

1. États :

Les états dans le labyrinthe sont définis par les coordonnées de chaque case accessible sur la grille. Chaque case du labyrinthe que le robot peut occuper (toutes les cases sauf celles qui sont des murs) représente un état unique dans l'espace d'état.

2. Actions :

Le robot peut se déplacer dans huit directions possibles : haut, bas, gauche, droite, ainsi que les quatre diagonales (haut-gauche, haut-droite, bas-gauche, bas-droite). Chaque action est représentée par un déplacement de la position actuelle du robot :

Haut : $(-1, 0)$

Bas : $(1, 0)$

Gauche : $(0, -1)$

Droite : $(0, 1)$

Haut-gauche : $(-1, -1)$

Haut-droite : $(-1, 1)$

Bas-gauche : $(1, -1)$

Bas-droite : $(1, 1)$

3. Récompenses :

Case sortie : +100 points.

Mur : -10 points (le robot ne peut pas traverser les murs, donc cette récompense peut être utilisée pour indiquer une tentative infructueuse de traverser un mur).

Autres cases : -1 point pour chaque déplacement, incitant le robot à trouver le chemin le plus court possible.

2.1.1 Implémentation de l'algorithme de Qlearning

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import matplotlib.patches as patches
4
5 class Maze:
6     def __init__(self, grid, start, goal, wall_penalty=-10,
7         step_penalty=-1, goal_reward=100):
8
9         self.grid = np.array(grid) #Représente le labyrinthe ou 1
10            indique un mur et 0 un espace libre
11
12         self.start = start
13         self.goal = goal
14
15         """start et goal : Définissent les coordonnées de départ et
16            de l'objectif dans le labyrinthe"""
17
18         self.wall_penalty = wall_penalty
19         self.step_penalty = step_penalty
20         self.goal_reward = goal_reward
21
22         """wall_penalty, step_penalty, goal_reward : Définissent les
23            récompenses (ou pénalités) pour heurter un mur, faire
24            un pas, et atteindre l'objectif, respectivement"""
25
26         self.actions = [(0, 1), (1, 0), (0, -1), (-1, 0)] # Right,
27            Down, Left, Up
28
29     def step(self, state, action):
```

```
19     """Prend un etat(l'etat initial c est (0,0) et une action,
    calcule le nouvel etat resultant de cette action.
20 Gere les collisions avec les murs et verifie si l'objectif est
    atteint.
21 Retourne le nouvel etat, la recompense/penalite correspondante,
    et un booleen indiquant si l'objectif est atteint."""
22     next_state = tuple(np.array(state) + np.array(action))
23     """Par exemple, si l'etat est (0, 0) et l'action est (1, 0)
    , cela donnera (0+1, 0+0), ce qui equivaut (1, 0)."""
24     if (next_state[0] < 0 or next_state[0] >= self.grid.shape[0]
        or
25         next_state[1] < 0 or next_state[1] >= self.grid.shape[1]
        or
26         self.grid[next_state] == 1): # Wall check
27         return state, self.wall_penalty, False
28     elif next_state == self.goal:
29         return next_state, self.goal_reward, True
30     else:
31         return next_state, self.step_penalty, False
32     """si le nouvel etat resulte en une position en dehors des
    limites de la grille ou s'il rencontre un mur (represente par
    la valeur 1 dans la grille). Si l'une de ces conditions est
    vraie, cela signifie que l'agent a rencontre un mur ou a
    depasse les limites de la grille. Dans ce cas, la methode
    retourne l'etat actuel (state), une penalite pour avoir
    rencontr un mur (self.wall_penalty), et un booleen
    indiquant que l'objectif n'a pas ete atteint (False)."""
33 def render(self, path):
34     """Visualise le labyrinthe en dessinant les murs, le chemin
    parcouru, et marque les positions speciales (depart, objectif
    ) avec des couleurs distinctes.
```

```
35     Utilise matplotlib pour la visualisation, avec des carres pour
    les murs et des points pour le chemin."""
36     fig, ax = plt.subplots()
37     # Set up the grid for visualization
38     for y in range(self.grid.shape[0]):
39         for x in range(self.grid.shape[1]):
40             if self.grid[y, x] == 1: # Wall
41                 ax.add_patch(patches.Rectangle((x, y), 1, 1,
42                     color='blue'))
43             if (y, x) == self.start: # Start
44                 ax.add_patch(patches.Rectangle((x, y), 1, 1,
45                     color='red'))
46             if (y, x) == self.goal: # Goal
47                 ax.add_patch(patches.Rectangle((x, y), 1, 1,
48                     color='green'))
49
50     # Draw path with small dots
51     for (y, x) in path:
52         ax.plot(x + 0.5, y + 0.5, 'yo', markersize=10, marker='o
53             ') # Centered yellow dots
54
55
56     plt.xlim(0, self.grid.shape[1])
57     plt.ylim(0, self.grid.shape[0])
58     plt.gca().invert_yaxis()
59     plt.grid(True)
60     plt.xticks(np.arange(0, self.grid.shape[1], 1))
61     plt.yticks(np.arange(0, self.grid.shape[0], 1))
62     plt.show()
63
64     def render1(self, path, title, num_cycles):
65         fig, ax = plt.subplots()
66         for y in range(self.grid.shape[0]):
67             for x in range(self.grid.shape[1]):
```

```
61         if self.grid[y, x] == 1:
62             ax.add_patch(patches.Rectangle((x, y), 1, 1,
63                                     color='blue'))
64         elif (y, x) == self.start:
65             ax.add_patch(patches.Rectangle((x, y), 1, 1,
66                                     color='red'))
67         elif (y, x) == self.goal:
68             ax.add_patch(patches.Rectangle((x, y), 1, 1,
69                                     color='green'))
70
71     for (y, x) in path:
72         ax.plot(x + 0.5, y + 0.5, 'yo', markersize=10)
73     ax.set_xlim(0, self.grid.shape[1])
74     ax.set_ylim(0, self.grid.shape[0])
75     ax.set_title(f'{title} - Cycles: {num_cycles}')
76     plt.gca().invert_yaxis()
77     plt.grid(True)
78     plt.show()
79
80 def q_learning(maze, episodes, alpha=0.1, gamma=0.9, epsilon=0.1):
81     """Implements l'algorithme Q-learning pour apprendre une politique
82     optimale de navigation dans le labyrinthe.
83
84     episodes : Nombre total de fois que l'agent tentera de parcourir le
85     labyrinthe pour apprendre.
86
87     alpha, gamma, epsilon : Parametres de l'algorithme (taux d'
88     apprentissage, facteur de depreciation, et probabilite d'
89     exploration).
90
91     Utilise une strategie epsilon-greedy pour equilibrer exploration et
92     exploitation.
93
94     Met a jour la table Q en utilisant les recompenses obtenues et les
95     estimations des meilleures recompenses futures."""
96     Q = np.zeros(maze.grid.shape + (4,))
```

```
83     for episode in range(epochs):
84         state = maze.start
85         done = False
86         #Extrait de Code Implementant epsilon-Greedy
87         while not done:
88             if np.random.random() < epsilon:
89                 action = np.random.choice(len(maze.actions)) #
90                 Explore
91             else:
92                 action = np.argmax(Q[state]) # Exploit
93                 next_state, reward, done = maze.step(state, maze.actions
94                 [action])
95                 old_value = Q[state][action]
96                 future_value = np.max(Q[next_state])
97                 Q[state][action] = old_value + alpha * (reward + gamma *
98                 future_value - old_value)
99                 state = next_state
100         return Q
101
102 """Le labyrinthe est initialise avec une grille specifiee, un point
103 de depart, et un point d arrivee.
104 La fonction q_learning est appelee pour entrainer l agent avec un
105 nombre defini d'episodes.
106 Apres l apprentissage, le chemin optimal est extrait en suivant la
107 politique optimale derivee des valeurs Q apprises.
108 Le chemin optimal est visualise a l aide de la methode render."""
109 # Define the grid (0: free, 1: wall)
110 grid = [
111     [0, 1, 0, 0, 0, 0, 0, 0, 0, 0],
112     [0, 1, 0, 1, 1, 1, 1, 1, 1, 0],
113     [0, 0, 0, 1, 0, 0, 0, 0, 0, 0],
114     [0, 1, 1, 1, 0, 1, 1, 1, 1, 0],
```

```
108     [0, 0, 0, 0, 0, 1, 0, 0, 0, 0]
109 ]
110 """0 et 1 : Dans ce tableau, 0 represente un espace libre ou le
    robot peut se deplacer, et 1 represente un mur, o le
    deplacement est bloque."""
111 start = (0, 0) #Coordonnes de la position de d part
112 goal = (4, 9) #Coordonnes de la position d'arrive
113
114
115 # Initialize the maze
116 maze = Maze(grid, start, goal)
117
118 # Run Q-learning
119 Q = q_learning(maze, 1000) #nbr episode=1000
120
121 # Extract the optimal path from Q-values
122 state = start
123 path = [state]
124 while state != goal:
125     action = np.argmax(Q[state])
126     state = tuple(np.array(state) + np.array(maze.actions[action]))
127     path.append(state)
128
129 # Visualize the path
130 maze.render(path)
```

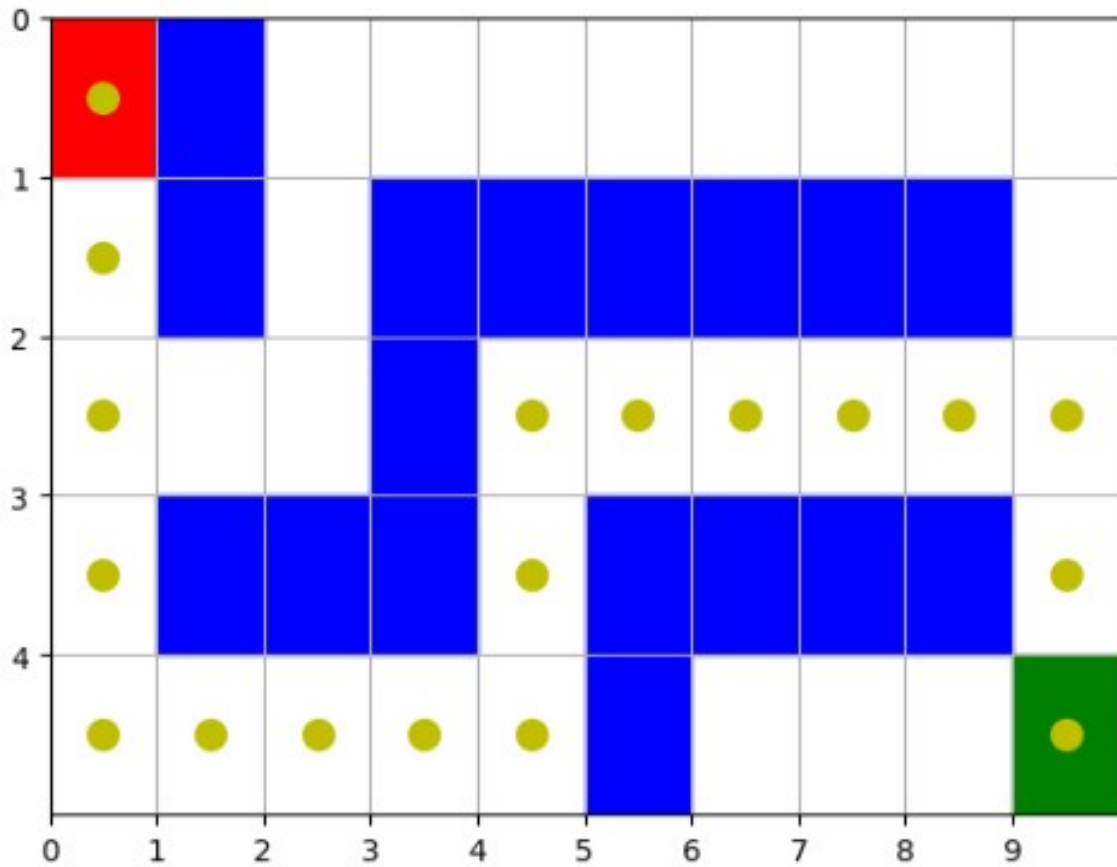


FIGURE2.1 – chemin du robot

Question3 : Faire varier ϵ , γ , (respectivement coefficient d apprentissage et facteur de depreciation) . donner à chaque fois le nombre de cycles pour trouver le chemin vers la sortie ainsi que le chemin trouvé.

Comparer les résultats obtenus à l'aide d'un graphique. Interpréter.

```
1 import numpy as np # Necessary for the operations on the
```

tableaux

```
2
3 def q_learning(maze, episodes, alpha=0.1, gamma=0.9, epsilon=0.1):
4     # Initialisation de la table Q avec des zeros pour toutes les
      paires etat action
5     Q = np.zeros(maze.grid.shape + (4,))
6
7     # Liste pour enregistrer le nombre de pas necessaires pour
      atteindre le but dans chaque episode
8     steps_per_episode = []
9
10    # Boucle sur le nombre total d'episodes specifiques
11    for episode in range(episodes):
12        # Commence chaque episode a la position de depart du
          labyrinthe
13        state = maze.start
14        done = False # Indicateur de fin d'episode, True si le but
          est atteint
15        steps = 0 # Compteur de pas pour l'episode courant
16
17        # Continue tant que l'objectif n'est pas atteint
18        while not done:
19            # Selection de l'action selon la politique epsilon-
              greedy
20            if np.random.random() < epsilon:
21                # Exploration: choix d'une action aleatoire
22                action = np.random.choice(len(maze.actions))
23            else:
24                # Exploitation: choix de l'action avec la valeur Q
              maximale pour l'etat actuel
25                action = np.argmax(Q[state])
```



```
26
27     # Execution de l'action choisie et obtention du nouvel
        état et de la récompense
28     next_state, reward, done = maze.step(state, maze.actions
        [action])
29
30     # Mise a jour de Q pour l'état actuel et l'action selon
        la formule du Q-learning
31     old_value = Q[state][action]
32     future_value = np.max(Q[next_state])
33     Q[state][action] = old_value + alpha * (reward + gamma *
        future_value - old_value)
34
35     # Passage a l'état suivant
36     state = next_state
37     steps += 1 # Incrementation du compteur de pas
38
39     # Enregistrement du nombre de pas dans la liste apres la fin
        de l'episode
40     steps_per_episode.append(steps)
41
42     # Retourne la table Q mise a jour et la liste du nombre de pas
        par episode
43     return Q, steps_per_episode
```

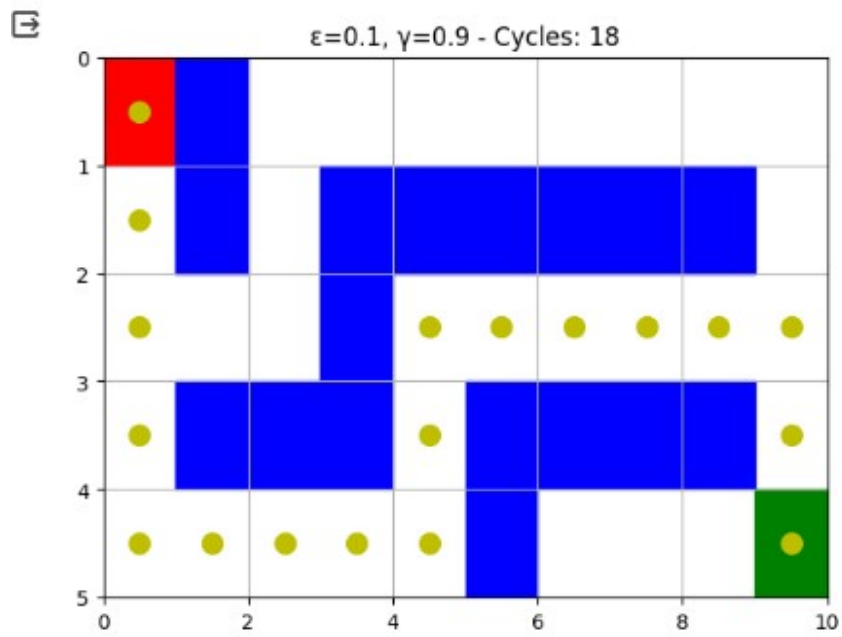


FIGURE2.2 – resultat1

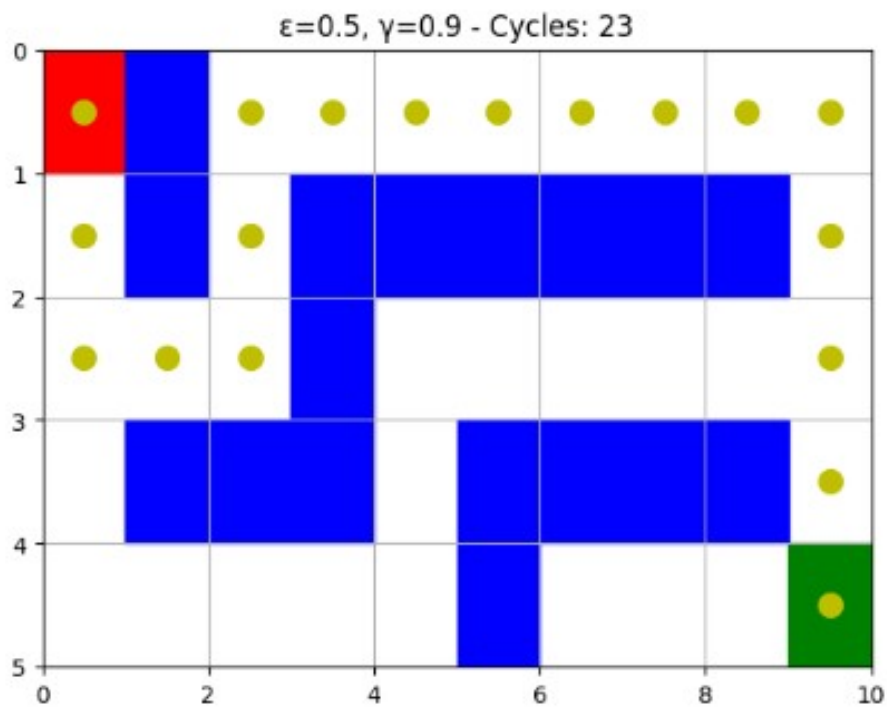


FIGURE2.3 – resultat2

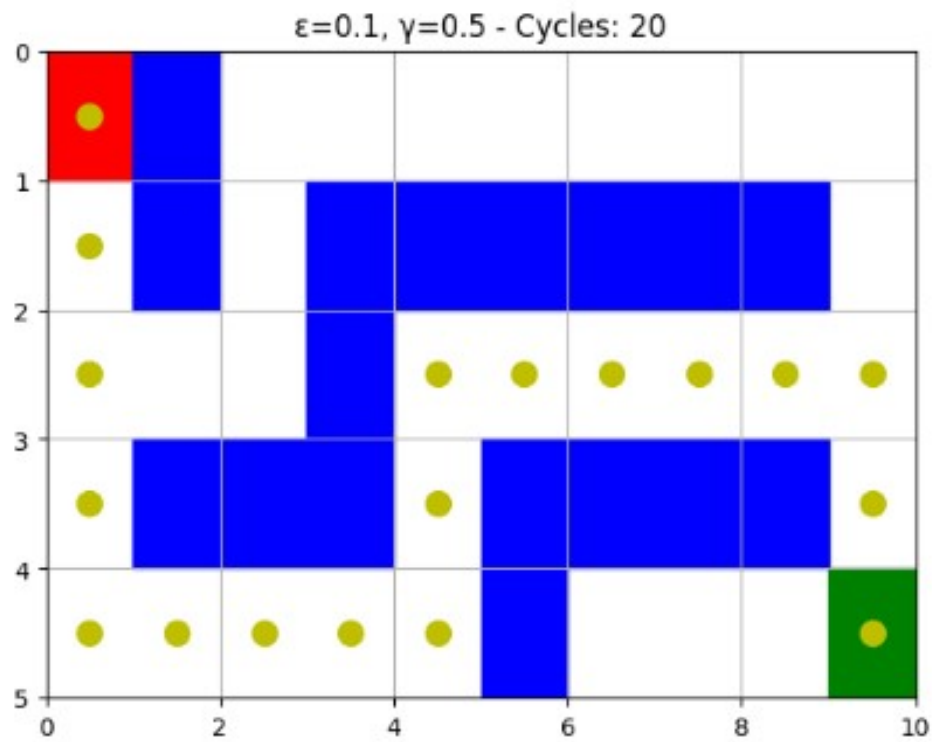


FIGURE2.4 – resultat3

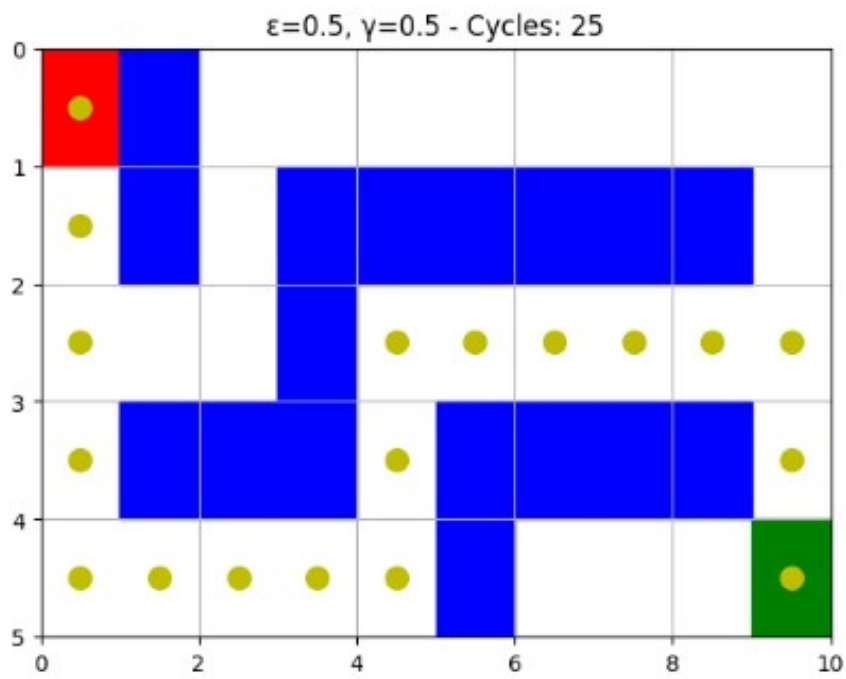


FIGURE2.5 – resultat4

2.1.2 Interpretation des Comparaison des resultats :

Première Image ($\epsilon = 0.1$, $\gamma = 0.9$, Cycles : 17) **Efficacité du chemin :**
Chemin le plus direct.

Cycles pour l'achèvement : Moins de cycles parmi tous les scénarios.

Interprétation : Un epsilon plus faible avec un gamma élevé conduit à une exploration plus directe et efficace, mettant fortement l'accent sur les récompenses à long terme et moins sur l'exploration aléatoire.

Deuxième Image ($\epsilon = 0.5$, $\gamma = 0.9$, Cycles : 31) **Efficacité du chemin :**
Plus d'exploration, chemin moins direct.

Cycles pour l'achèvement : Plus grand nombre de cycles.

Interprétation : Un epsilon plus élevé conduit à plus d'**exploration**, résultant en un chemin moins direct mais offrant plus d'apprentissage sur l'environnement. Ce scénario illustre le coût d'une exploration plus élevée avec une forte focalisation sur les récompenses futures.

Troisième Image ($\epsilon = 0.1$, $\gamma = 0.5$, Cycles : 19) **Efficacité du chemin :**
Assez direct, légèrement plus de cycles que le plus faible epsilon et le plus haut gamma.

Cycles pour l'achèvement : Modérément efficace. **Interprétation :** Un epsilon plus faible avec un gamma plus bas conduit à un chemin plus direct mais légèrement moins efficace que le paramètre de gamma le plus élevé. Cela indique une approche équilibrée entre les récompenses immédiates et futures.

Quatrième Image ($\epsilon = 0.5$, $\gamma = 0.5$ - Cycles : 30)

Efficacité du chemin : Exploration similaire à celle avec un gamma élevé mais moins efficace en termes de cycles par rapport à un gamma élevé.

Cycles pour l'achèvement : Élevé, similaire à une exploration élevée avec un gamma élevé. Interprétation : Une exploration accrue avec une focalisation équilibrée sur les récompenses immédiates et futures conduit à apprendre beaucoup sur l'environnement mais au prix de l'efficacité.

Analyse graphique et interprétation :

Pour fournir une interprétation plus claire, un graphique montrant les cycles pour chaque paramètre peut être tracé. Cependant, selon les images fournies :

*Un epsilon plus faible (0.1) conduit systématiquement à des chemins plus directs lorsque le gamma est élevé, ce qui indique qu'une exploration moins aléatoire combinée à une forte focalisation sur les récompenses futures (gamma élevé) est plus efficace pour cette configuration de labyrinthe particulière. Un epsilon plus élevé (0.5) augmente l'exploration, ce qui peut être bénéfique dans des labyrinthes plus complexes ou des environnements différents mais tend à augmenter le nombre de cycles nécessaires pour trouver le but dans ce scénario.

Conclusion :

Pour des labyrinthes similaires à ceux testés, un epsilon plus faible et un gamma plus élevé fournissent le parcours le plus efficace, suggérant que la minimisation des actions aléatoires tout en pondérant fortement les récompenses futures

est optimal. Cette stratégie, cependant, pourrait ne pas être aussi efficace dans des environnements où la connaissance initiale de l'environnement est faible, ou où l'environnement est très dynamique, car cela pourrait conduire à une convergence prématurée vers des chemins sous-optimaux. En revanche, une exploration plus élevée (epsilon plus élevé) peut être avantageuse dans des environnements plus stochastiques ou complexes où une exploration plus approfondie est nécessaire pour éviter les minima locaux.

Question 4 : Faire varier la taille du labyrinthe en choisissant aleatoirement des murs :

a. Grille 20x20

b. Grille 30x30

Comparer à l'aide d'un graphique, les resultats obtenus avec les 3 différentes tailles de labyrinthe.

```
1 def generate_random_maze(size, density=0.1):
2     grid = np.zeros((size, size), dtype=int) # Cree un tableau de
3         zeros pour le labyrinthe
4     start, goal = (0, 0), (size - 1, size - 1) # Definit les
5         positions de depart et d'arrivee
6     obstacles = int(size * size * density) # Calcule le nombre d'
7         obstacles a placer
8
9     # Place des obstacles sur le tableau
10    while obstacles > 0:
11        x, y = random.randint(0, size - 1), random.randint(0, size -
12            1)
13        if (x, y) != start and (x, y) != goal and grid[y, x] == 0:
14            grid[y, x] = 1
15            obstacles -= 1
16
17    return grid, start, goal
18
19 """Le but principal de la fonction simulate_maze est de tester et d'
20     evaluer l'efficacite de l'algorithme Q-learning dans un
21     environnement de labyrinthe, o l'objectif est de naviguer d'un
22     point de depart
23     a un point d'arrivee.
24     Cette fonction permet de mesurer combien de pas l'agent (le
25     simulateur du labyrinthe)
```

```
16 prend pour atteindre l'objectif a travers differents episodes, en
    utilisant divers parametres de l'algorithm Q-learning."""
17 def simulate_maze(size, density=0.1, episodes=100, alpha=0.1, gamma
    =0.9, epsilon=0.1):
18     grid, start, goal = generate_random_maze(size, density) #
        Genere le labyrinthe
19     steps_to_goal = [] # Stocke le nombre de pas pour atteindre l'
        objectif pour chaque episode
20
21     Q = {} # Initialisation du dictionnaire pour stocker les
        valeurs Q
22
23     for episode in range(episodes): # Boucle sur chaque episode
24         state = start # Commence a la position de depart
25         steps = 0 # Compteur de pas
26
27         while state != goal: # Continue tant que l'objectif n'est
            pas atteint
28             if random.random() < epsilon: # Decide aleatoirement d'
                explorer
29                 action = random.choice([(0, 1), (1, 0), (0, -1),
                    (-1, 0)])
30             else: # Exploite la meilleure action connue
31                 action = max([(0, 1), (1, 0), (0, -1), (-1, 0)], key
                    =lambda a: Q.get((state, a), 0))
32
33                 next_state = (state[0] + action[0], state[1] + action
                    [1]) # Calcule le prochain etat
34
35                 # Calcule la recompense et verifie si le prochain etat
                    est le but
```



```
36         if next_state == goal:
37             reward, done = 100, True
38         elif not (0 <= next_state[0] < size and 0 <= next_state
39                 [1] < size) or grid[next_state] == 1:
40             reward, done = -10, False
41             next_state = state # Revenir l'etat precedent
42                                 si c'est un mouvement invalide
43
44         else:
45             reward, done = -1, False
46
47             # Mise a jour de la valeur Q
48             q_key = (state, action)
49             old_value = Q.get(q_key, 0)
50             future_value = max([Q.get((next_state, a), 0) for a in
51                               [(0, 1), (1, 0), (0, -1), (-1, 0)]]))
52             Q[q_key] = old_value + alpha * (reward + gamma *
53                                             future_value - old_value)
54
55             state = next_state
56             steps += 1
57
58             if done:
59                 break
60
61             steps_to_goal.append(steps) # Ajoute le nombre de pas de
62                                         cet episode a la liste
63
64     return steps_to_goal
```

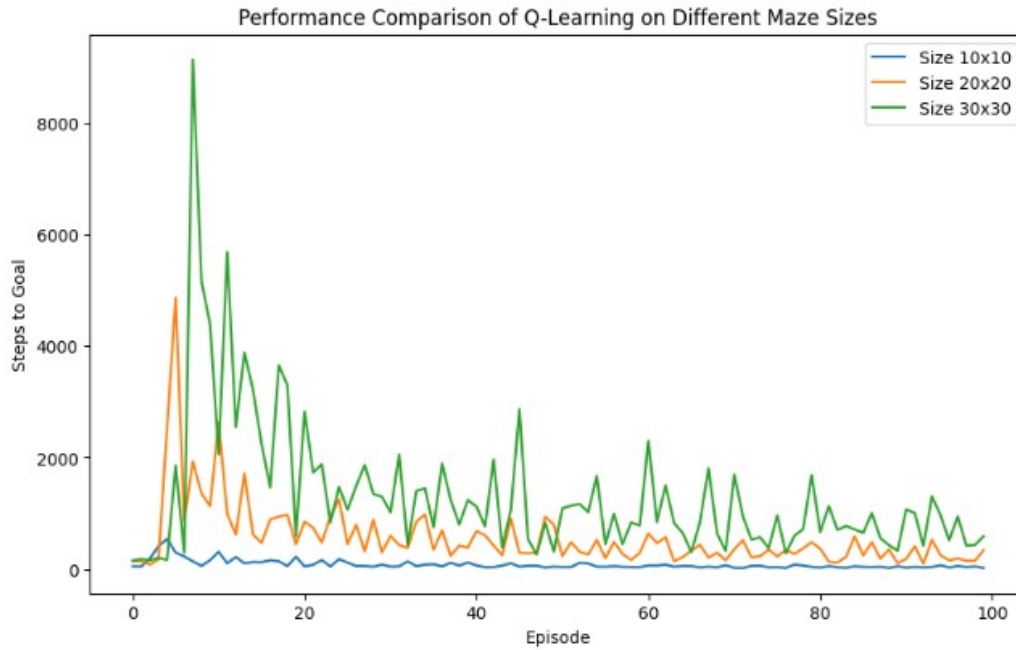


FIGURE2.6 – resultat

2.1.3 Analyse du Graphique

Décroissance du Nombre de Pas : On observe que pour toutes les tailles de labyrinthe, le nombre de pas nécessaires pour atteindre l'objectif diminue significativement au cours des premiers épisodes. Cela indique que l'algorithme apprend efficacement la structure du labyrinthe et améliore ses stratégies de navigation pour minimiser les pas.

2.1.3.1 Variabilité des Performances :

Labyrinthe 10x10 : Presente la plus rapide diminution du nombre de pas, avec une stabilisation rapide, ce qui suggere que l'algorithme parvient plus rapidement à optimiser son chemin dans un espace plus restreint.

Labyrinthe 20x20 : Montre une décroissance plus lente et plus de fluctuations dans le nombre de pas au fil des épisodes, ce qui peut refléter une complexité accrue et une difficulté plus grande pour l'algorithme à optimiser son approche dans un espace plus grand.

Labyrinthe 30x30 : Bien que les performances s'améliorent, les pics de variabilité restent élevés tout au long des simulations, indiquant des défis continus en termes d'optimisation de la navigation dans un espace encore plus vaste.

2.1.3.2 Comparaison des Tailles de Labyrinthe :

Les performances tendent à être moins stables dans les labyrinthes plus grands. Cela pourrait être dû à un plus grand nombre de configurations possibles de murs et de chemins, rendant l'apprentissage et la généralisation plus difficiles pour l'algorithme.

Conclusion et Implications Efficacité de l'Apprentissage : L'algorithme Q-learning montre une capacité notable à apprendre et à s'adapter à l'environnement du labyrinthe, comme en témoigne la réduction globale du nombre de pas nécessaires pour atteindre l'objectif au fil du temps.

Influence de la Taille du Labyrinthe : La taille du labyrinthe a un impact significatif sur la rapidité et la stabilité de l'apprentissage. Dans les environnements plus grands, il pourrait être nécessaire d'ajuster les paramètres de l'algorithme (comme epsilon, alpha, ou gamma) ou d'augmenter le nombre d'épisodes de formation pour obtenir des performances optimales.

Utilisation Pratique : Ces résultats peuvent guider la configuration des paramètres de l'algorithme lors de l'application de Q-learning à des problèmes réels de navigation ou de planification de trajet, en montrant l'importance de considérer la taille et la complexité de l'environnement lors de la conception et du réglage des algorithmes d'apprentissage par renforcement.

2.1.4 Visualisation du parcours du robot dans labyrinthe

Ce module définit une classe Maze pour simuler un labyrinthe où un agent peut naviguer en utilisant l'algorithme de Q-learning. La classe Maze est conçue pour générer dynamiquement un labyrinthe avec une densité spécifiée de murs et de positions aléatoires pour le départ et l'arrivée. Elle permet également de définir les récompenses et les pénalités pour différentes actions dans le labyrinthe.

Fonctionnalités : - Initialiser un labyrinthe avec des paramètres spécifiques pour la taille, la densité des murs, et les valeurs de récompenses/pénalités.

- Simuler le déplacement d'un agent dans le labyrinthe avec gestion des collisions et récompenses.
- Visualiser graphiquement le chemin parcouru par l'agent en temps réel, ce qui

aide à comprendre comment l'agent apprend et navigue dans le labyrinthe.

- Implémenter l'apprentissage par renforcement avec une stratégie d'action epsilon-greedy pour optimiser le chemin pris par l'agent en maximisant les récompenses cumulées.
- Extraire et afficher le chemin optimal à partir des valeurs Q apprises après une série d'épisodes d'apprentissage.

L'objectif est d'utiliser ce modèle pour étudier et démontrer les capacités de l'apprentissage par renforcement, en particulier l'algorithme de Q-learning, pour résoudre des problèmes de navigation complexes dans des environnements avec obstacles.

```
1 class Maze:
2     def __init__(self, grid_size, wall_density=0.2, wall_penalty
3         =-10, step_penalty=-1, goal_reward=100):
4         # Initialisation du labyrinthe avec une taille specifiee,
5         # densite des murs, et des valeurs de penalite et de
6         # recompense.
7         self.grid = self.initialize_grid(grid_size, wall_density) #
8         # Creation du grille avec des murs
9         self.start, self.goal = self.random_positions() #
10        # Definition aleatoire des positions de depart et d'arrivee
11        self.actions = [(0, 1), (1, 0), (0, -1), (-1, 0)] #
12        # Definition des actions possibles : droite, bas, gauche,
13        # haut
14
15    def initialize_grid(self, grid_size, wall_density):
16        # Cree un grille ou les murs sont places aleatoirement selon
```

```
        la densite specifiee.
10     return np.random.choice([0, 1], size=grid_size, p=[1-
        wall_density, wall_density])
11
12     def step(self, state, action):
13         # Determine le nouvel etat apres avoir pris une action, gere
        les collisions avec les murs et les recompenses.
14         next_state = tuple(np.array(state) + np.array(action))
15         if self.grid[next_state] == 1 or not (0 <= next_state[0] <
        self.grid_size[0] and 0 <= next_state[1] < self.grid_size
        [1]):
16             return state, self.wall_penalty, False # Reste en place
        si collision avec un mur, retourne une penalite
17         elif next_state == self.goal:
18             return next_state, self.goal_reward, True # Retourne
        une grande recompense si l'objectif est atteint
19         else:
20             return next_state, self.step_penalty, False # Mouvement
        normal avec un petit cout
21
22     def render(self, path, interactive=False):
23         # Affichage graphique du labyrinthe, des positions de depart
        et d'arrivee, et du chemin parcouru.
24         for (y, x) in path:
25             ax.plot(x + 0.5, y + 0.5, 'yo') # Dessine le chemin en
        jaune
26         if interactive:
27             clear_output(wait=True)
28             time.sleep(0.1) # Attend un peu pour une visualisation
        en temps reel
29         plt.show()
```

```
30 """Dans cette fonction, chaque fois que l'agent prend une action et
    met e jour son etat, la methode render est appelee pour
    visualiser l'etat actuel
31     du labyrinthe et le chemin parcouru par l'agent jusqu'e ce
        point. Cela permet de voir comment l'agent se deplace
        dans le labyrinthe, quels chemins il prend,
32     et comment il reagit aux recompenses et aux penalites."""
33 def real_time_q_learning(maze, episodes, alpha=0.1, gamma=0.9,
    epsilon=0.1):
34     # Implementation de l'apprentissage par renforcement Q-learning
        pour naviguer dans le labyrinthe.
35     while not done:
36         if np.random.random() < epsilon:
37             action = np.random.choice(len(maze.actions)) #
                Exploration aleatoire
38         else:
39             action = np.argmax(Q[state]) # Exploitation : choisit
                la meilleure action connue
40         next_state, reward, done = maze.step(state, maze.actions[
            action])
41         # Mise e jour de Q-value pour l'action prise
42         Q[state][action] = old_value + alpha * (reward + gamma *
            future_value - old_value)
43         path.append(state) # Ajoute l'etat au chemin
44         maze.render(path, interactive=True) # Affichage du chemin
            en temps reel
45
46 def extract_optimal_path(maze, Q):
47     # Extraction du chemin optimal base sur la table Q apres l'
        apprentissage.
48     while state != maze.goal:
```

```
49     action = np.argmax(Q[state])
50     state = tuple(np.array(state) + np.array(maze.actions[action
51         ])) # Avance vers le meilleur etat suivant
    optimal_path.append(state)
```

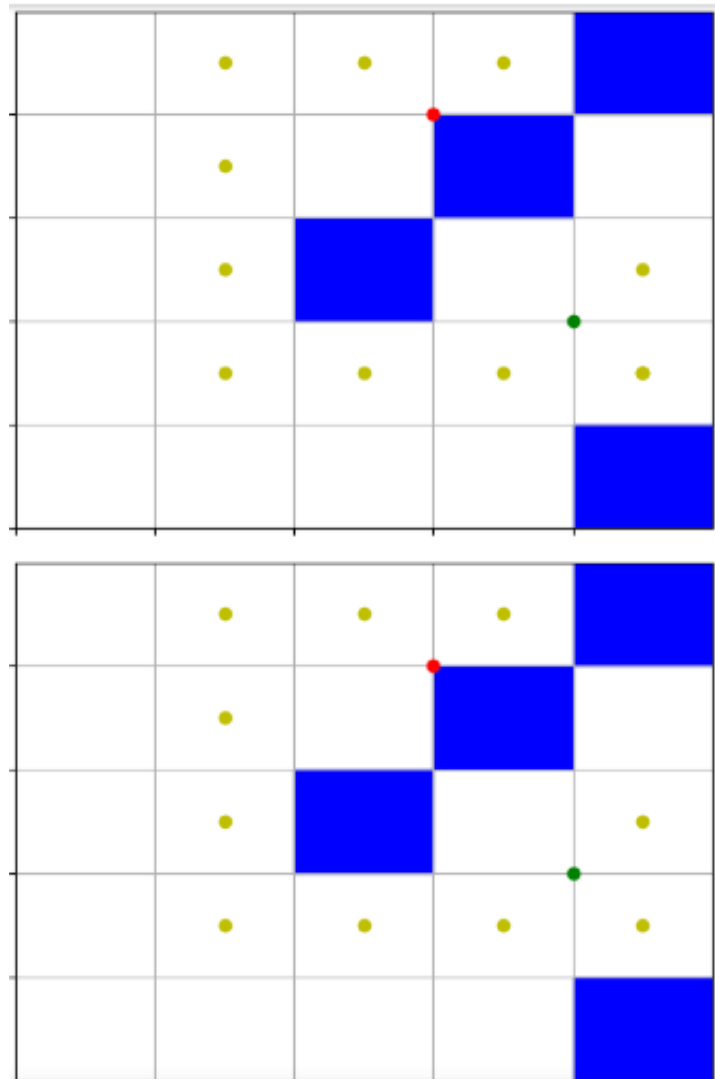


FIGURE2.7 – le output

=> L’algorithme de Q-learning démontre son efficacité à naviguer et trouver le chemin optimal dans le labyrinthe en s’adaptant intelligemment aux récompenses et pénalités rencontrées.

