

TME 2 - MPD- Victor Duthoit

Imports

```
1 import gym
2 import gridworld
3 import randomAgent
4 import matplotlib
5 import matplotlib.pyplot as plt
6 # matplotlib.use("Qt5agg")
7 %matplotlib inline
8 import time
9 import numpy as np
10 import operator
11 import progressbar
12 import pandas as pd
```

Utilities

```
1 def mse(a,b):
2     #return the norm 2 (mean square error) of to functions
3     return np.sqrt(sum( (a[x]-b[x])**2 for x in a ))
4
5 def dic_argmax(d):
6     #return the argmax for dictionnaires
7     return max(d.items(), key=operator.itemgetter(1))[0]
8
9 pI, vI, rD = 'policy_iteration', 'value_iteration', 'random'
```

Policies algorithms

1. Policy iteration

```
1 def policyIteration(states, P, eps=0.01, gamma=0.5):
2     #return the policy following the policy iteration algorithm
3     final_states = set(states)-set(P)
4     pi = {s : np.random.randint(4) for s in P } #initialize policy
5     #while the policy is unstable
6     while True :
7         v = {s : 0 for s in states } #initialize the value
8         while True : #while value is not stable
9             #update of v by the Bellman equation
10            v_new = {s : sum( p*(rew + gamma*v[s_p]) for p,s_p,rew,done in P[s]
11                [pi[s]]) for s in P }
12            for f in final_states:
13                v_new[f] = 0
14            if mse(v,v_new)<eps:
```

```

14         v = v_new
15         break
16         v = v_new
17         pi_new = {}
18         for s in P:
19             #choose the action with the highest value
20             pi_action = {a : sum( p*(rew+gamma*v[s_p]) for p,s_p,rew,done in P[s]
[a] ) for a in P[s]}
21             pi_new[s] = dic_argmax(pi_action)
22             if mse(pi,pi_new)==0:
23                 break
24             pi = pi_new
25         return pi

```

2. Value iteration

```

1  def valueIteration(states, P, eps=0.01, gamma=0.5):
2      #return the policy following the value iteration algorithm
3      final_states = set(states)-set(P)
4      v = {s : 0 for s in states } #initialize the value
5      while True : #while value is not stable
6          v_new = {s: max([sum( p*(rew + gamma*v[s_p]) for p,s_p,rew,done in P[s]
[a]) for a in P[s]]) for s in P }
7          for f in final_states:
8              v_new[f] = 0
9              if mse(v,v_new)<eps:
10                 v = v_new
11                 break
12             v = v_new
13         pi = {s: dic_argmax(
14             {a : sum( p*(rew + gamma*v[s_p]) for p,s_p,rew,done in P[s][a]) for a in
P[s] }
15         ) for s in P }
16         return pi

```

Agent

Nous définissons des agents qui ont la particularité d'avoir des politiques déterministes à définir.

```

1  class Agent(randomAgent.RandomAgent):
2      #an agent pocesses a policy pi
3
4      def __init__(self, action_space):
5          super().__init__(action_space)
6          self.pi = dict()
7
8      def act(self, observation, reward, done):
9          return self.pi[observation]

```

Environnement - plan0

```

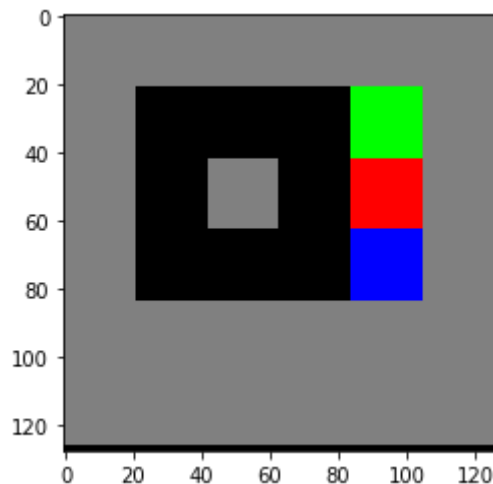
1  def createEnv(plan='0', visual=False):

```

```

2     # matplotlib.use("Qt5agg")
3     env = gym.make("gridworld-v0")
4     env.setPlan(f"gridworldPlans/plan{plan}.txt", {0: -0.001, 3: 1, 4: 1, 5: -1,
6: -1})
5     states, P = env.getMDP()
6     env.seed(0) # Initialise le seed du pseudo-random
7     # print(env.action_space) # Quelles sont les actions possibles
8     # print(env.step(1)) # faire action 1 et retourne l'observation, le reward,
    et un done un boolean (jeu fini ou pas)
9     env.render() # permet de visualiser la grille du jeu
10    #     env.render(mode="human") #visualisation sur la console
11    statedic, mdp = env.getMDP() # recupere le mdp : statedic
12    #     print("Nombre d'etats : ",len(statedic)) # nombre d'etats ,statedic :
    etat-> numero de l'etat
13    state, transitions = list(mdp.items())[0]
14    return env
15    # print(state) # un etat du mdp
16    # print(transitions) # dictionnaire des transitions pour l'etat : {action->
    [proba,etat,reward,done]}
17
18    createEnv()

```



```

1 <gridworld.gridworld_env.GridworldEnv at 0x1a30578828>

```

Test agent

On définit deux fonctions de test. L'une pour tester un agent dans un environnement connu, l'autre pour tester tous les paramètres en entrée.

```

1 def testAgent(agent, env, episode_count=1, visual=False, freq=100, FPS=0.01):
2     #test an agent for a given environment
3     if visual:
4         matplotlib.use('Qt5Agg')
5         reward = 0
6         done = False
7         rsum = 0
8         R = 0

```

```

9     actions = 0
10    for i in range(episode_count):
11        obs = env.reset()
12        if visual:
13            print(obs)
14        env.verbose = (i % freq == 0 and i > 0) # afficher 1 episode sur 100
15        if env.verbose and visual:
16            env.render(FPS)
17        j = 0
18        rsum = 0
19        while True:
20            action = agent.act(gridworld.GridworldEnv.state2str(obs), reward,
done)
21            obs, reward, done, _ = env.step(action)
22            rsum += reward
23            j += 1
24            if env.verbose and visual:
25                env.render(FPS)
26            if done:
27                if visual:
28                    print("Episode : " + str(i) + " rsum=" + str(rsum) + ", " +
str(j) + " actions")
29                    R+= rsum
30                    actions+= j
31                    break
32        avg_rsum = R/episode_count
33        avg_action = actions/episode_count
34        env.close()
35        if visual :
36            %matplotlib inline
37        return [avg_rsum, avg_action]

```

```

1    def testSeries(plan = ['0'],
2                    agents = [pI, vI, rD],
3                    epsilon = [0.1],
4                    gamma = [0.1],
5                    episode_count = [100]):
6        #tests all the agents for given parameters
7        res = []
8        N = len(plan)*len(agents)*len(epsilon)*len(gamma)*len(episode_count)
9        i=0
10
11        bar = progressbar.ProgressBar(maxval=N, \
12        widgets=[progressbar.Bar('=', '[', ']'), ' ', progressbar.Percentage()])
13        bar.start()
14        for p in plan:
15            env = gym.make("gridworld-v0")
16            env.setPlan(f"gridworldPlans/plan{p}.txt", {0: -0.001, 3: 1, 4: 1, 5: -1,
6: -1})
17            states, P = env.getMDP()
18            agent = Agent(env.action_space)
19            for eps in epsilon:
20                for g in gamma:
21                    for ec in episode_count:
22                        for ag in agents:
23                            if ag==rD:
24                                agent = randomAgent.RandomAgent(env.action_space)

```

```

25         t_pi = 0.
26     else:
27         t1_pi = time.time()
28         if ag==pI:
29             policy = policyIteration(states, P, eps=eps,
gamma=g)
30         if ag==vI:
31             policy = valueIteration(states, P, eps=eps,
gamma=g)
32         t_pi = time.time()-t1_pi
33         agent.pi = policy
34         t1_test = time.time()
35         test = testAgent(agent,env, episode_count = ec)
36         t_test = time.time()-t1_test
37         res.append([p,eps,g,ec,ag]+test+[t_pi,t_test])
38         i+=1
39         bar.update(i)
40     bar.finish()
41     columns =
['plan','eps','gamma','episode_count','agent','rsum','action','t_pi','t_test']
42     return pd.DataFrame(res,columns=columns)

```

```

1  def plotSeries(df, xaxis='gamma',yaxis='rsum',legend='agent',legendValue=
[pI,vI],bar=False):
2      plot = ()
3      for legendV in legendValue:
4          subdf = df[df[legend]==legendV]
5          x = subdf[xaxis].to_numpy()
6          y = subdf[yaxis].to_numpy()
7          plot+=(x,y)
8      if bar:
9          plt.bar(*plot)
10     else:
11         plt.plot(*plot)
12         plt.ylabel(yaxis)
13         plt.xlabel(xaxis)
14         plt.legend(legendValue)

```

```

1  # Execution avec un Agent
2
3  # agent = randomAgent.RandomAgent(env.action_space)
4  eps = 0.1
5  gamma = 0.1
6  policy = policyIteration(states, P, eps=eps, gamma=gamma)
7  # policy = valueIteration(states, P,eps=eps, gamma=gamma)
8  agent = Agent(env.action_space)
9  agent.pi = policy
10 testAgent(agent, env)

```

```

1  [0.979, 22.0]

```

Premier tests

Test for plan0, seed(0), 10000 episode, esp = 0.0001, gamma = 0.9 :

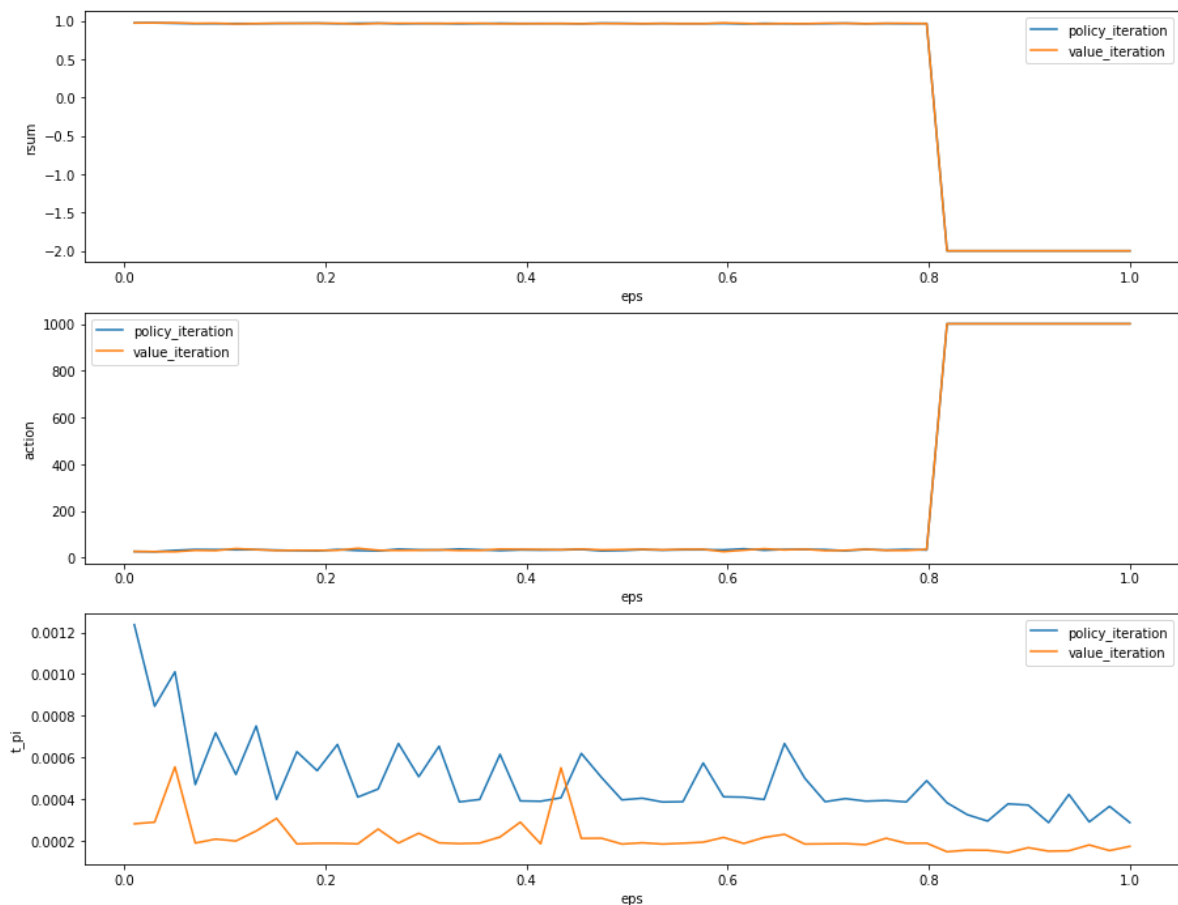
- **policy iteration agent:** average rsum : 0.56, average action : 4.79
- **value iteration:** average rsum : 0.55, average number of action : 4.80
- **random agent:** average rsum : -0.81, average action : 11.33

Les agents "politisés" sont bien meilleurs.

Choix d'epsilon

```
1 epsilon = np.linspace(0.01,1,num=50)
2 df_eps = testSeries(epsilon=epsilon,agents=[pI,vI])

1 #plots
2 plt.figure(figsize=(15,12))
3 plt.subplot(311)
4 plotSeries(df_eps, xaxis='eps', yaxis='rsum')
5 plt.subplot(312)
6 plotSeries(df_eps, xaxis='eps', yaxis='action')
7 plt.subplot(313)
8 plotSeries(df_eps, xaxis='eps', yaxis='t_pi') #t_pi is the time to determine the
policy
```



On remarque qu'un epsilon relativement élevé est suffisant pour une bonne détermination de la politique. On choisira ainsi un $\epsilon = 0.1$ pour la suite des calculs. On remarque de plus que l'itération de la politique est plus longue à déterminer.

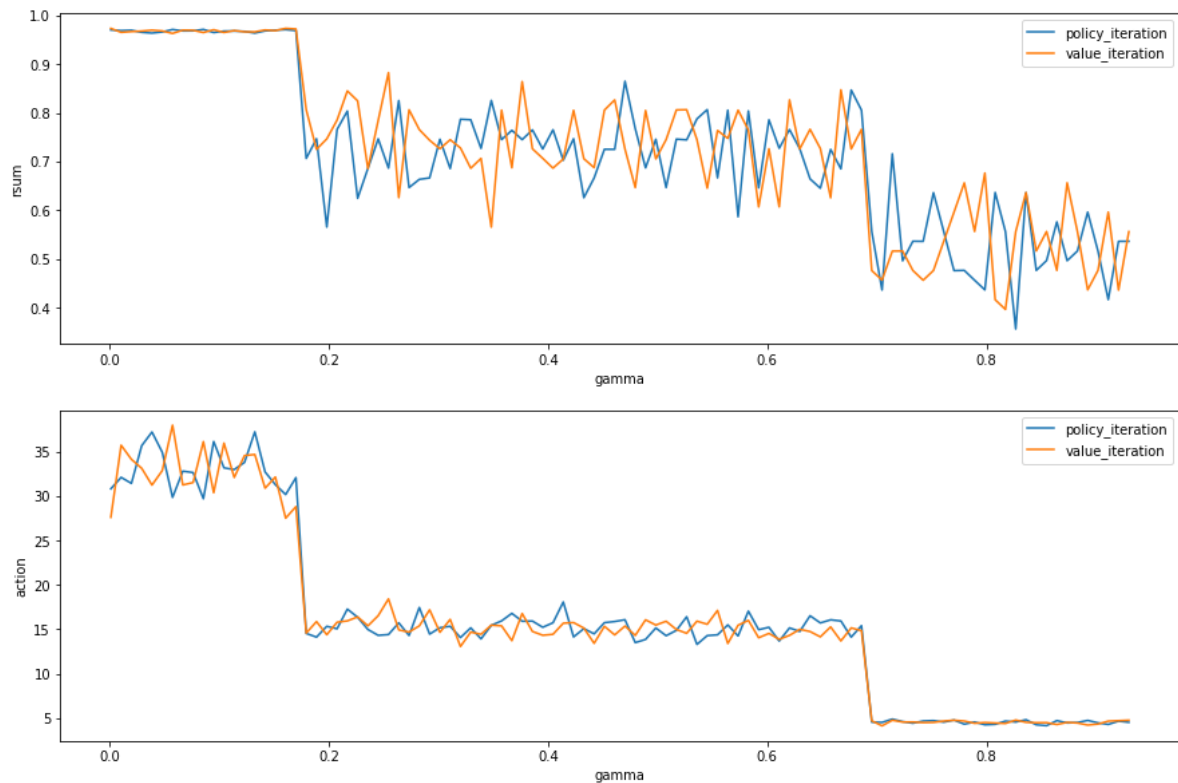
Rôle de gamma

```
1 gamma = np.linspace(0.001,0.93,num=100)
2 df_gamma = testSeries(gamma=gamma,agents=[pI,vI])
```

```
1 [= ] 2%
2
3 0
```

```
1 [======] 100%
```

```
1 #plots
2 plt.figure(figsize=(15,10))
3 plt.subplot(211)
4 plotSeries(df_gamma, xaxis='gamma', yaxis='rsum')
5 plt.subplot(212)
6 plotSeries(df_gamma, xaxis='gamma', yaxis='action')
```

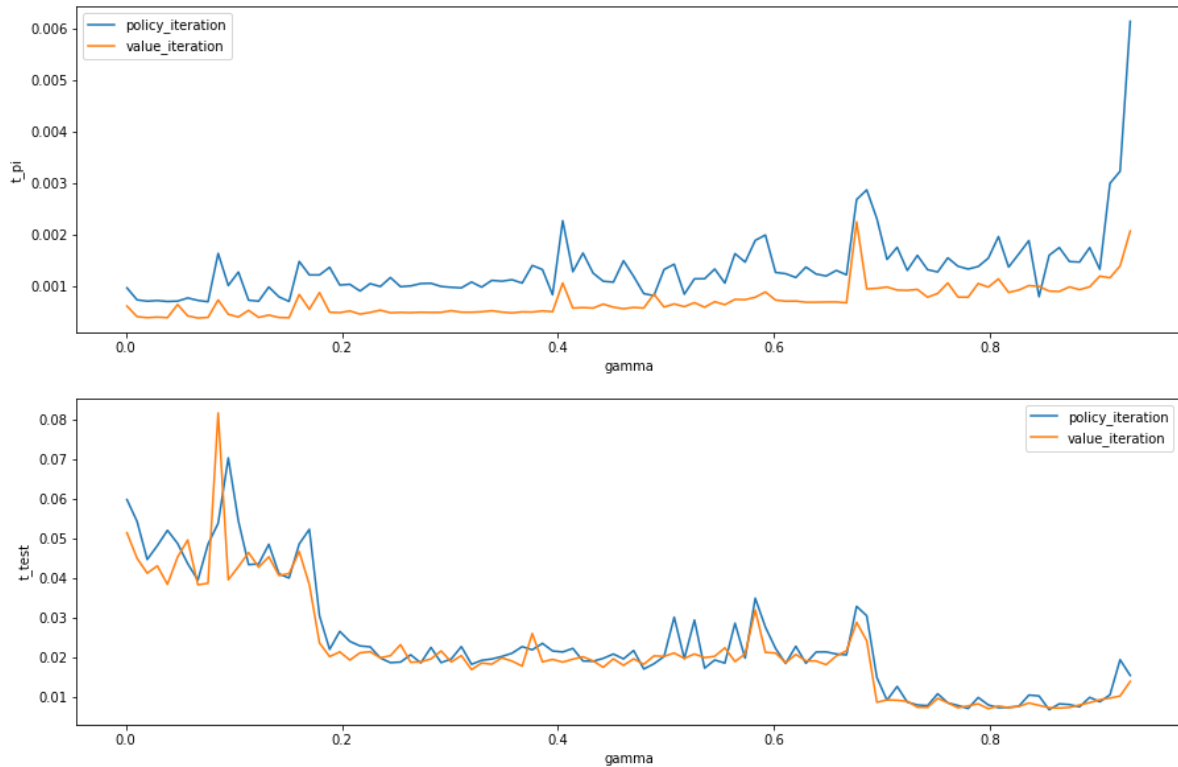


On peut noter deux points importants :

1. Les deux algorithmes semblent procurer des résultats similaires. On pourra ainsi se contenter de l'algorithme d'itération de la valeur.
2. Gamma joue un rôle essentiel dans le choix des stratégies de jeux. Apparemment, un gamma élevé privilégiera des gains accessibles rapidement, la partie s'écourtera ainsi plus rapidement. En revanche un gamma moindre a tendance à ne pas trop dévaloriser les gains accessibles en plus de coûts. Les paliers incitent à croire qu'un gain est possible en effectuant plus d'actions. Ainsi, cela donne à penser que le gamma dépend de la topologie du terrain.
3. L'optimal de reward est atteint pour des valeurs proches de 0 ou de 1. On retiendra la valeur 0.1 pour gamma (proche de 0) car le palier proche de 1 semble limité.

Gestion du temps

```
1 #plots
2 plt.figure(figsize=(15,10))
3 plt.subplot(211)
4 plotSeries(df_gamma, xaxis='gamma', yaxis='t_pi')
5 plt.subplot(212)
6 plotSeries(df_gamma, xaxis='gamma', yaxis='t_test')
```

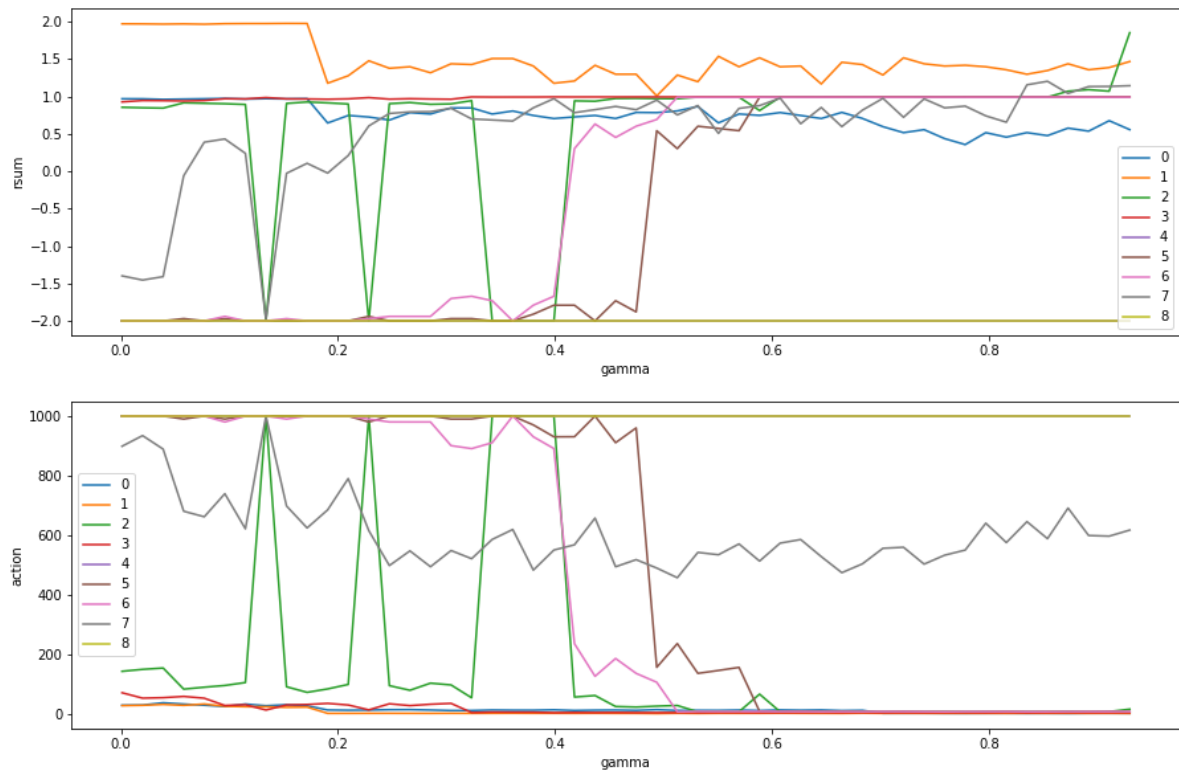


1. La première observation que l'on peut faire est que l'ordre de grandeur de calcul de la politique est très faible devant les temps de tests.
2. De plus, on remarque que l'algorithme d'itération de la valeur est plus rapide à terminer que son homologue d'itération de la politique.
3. Les temps de tests sont similaires pour les deux algos et ressemblent grossièrement à la distribution du nombre d'actions.

Sur de nouveaux terrains

```
1 gamma = np.linspace(0.001,0.93,num=50)
2 plan = list(range(9))
3 df_plan = testSeries(plan=plan, gamma=gamma,agents=[vI])
```

```
1 #plots
2 plt.figure(figsize=(15,10))
3 plt.subplot(211)
4 plotSeries(df_plan, xaxis='gamma', yaxis='rsum', legend='plan', legendValue=plan)
5 plt.subplot(212)
6 plotSeries(df_plan, xaxis='gamma', yaxis='action', legend='plan',
7 legendValue=plan)
```

On remarque que les plans n'ont pas les mêmes ordres de grandeur concernant le nombre d'actions et les gains. Néanmoins, nous pouvons voir que la dépendance en fonction de gamma change avec les plans. Voyons de plus près.

```

1  def plotPlan(df, p):
2      sub_df = df[df['plan'] == p]
3      plt.figure(figsize=(15,8))
4      plt.subplot(211)
5      plotSeries(sub_df, xaxis='gamma', yaxis='rsum', legendValue=[vI])
6      plt.subplot(212)
7      plotSeries(sub_df, xaxis='gamma', yaxis='action', legendValue=[vI])

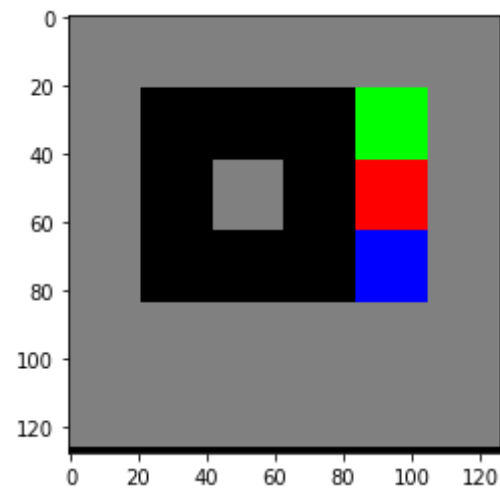
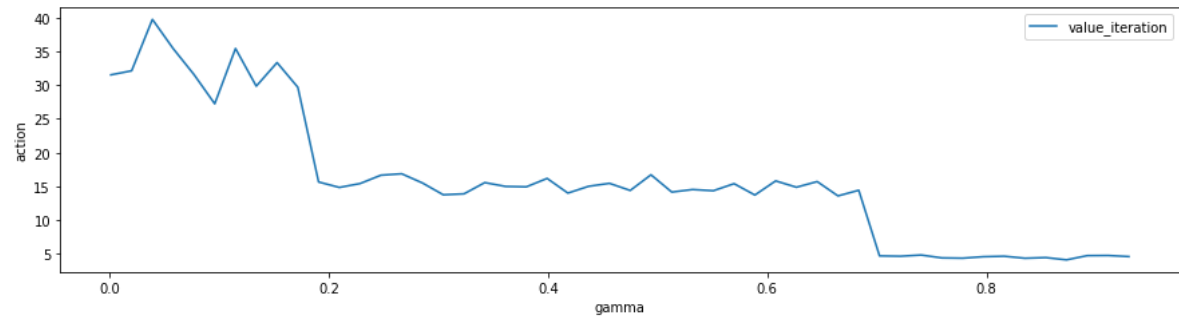
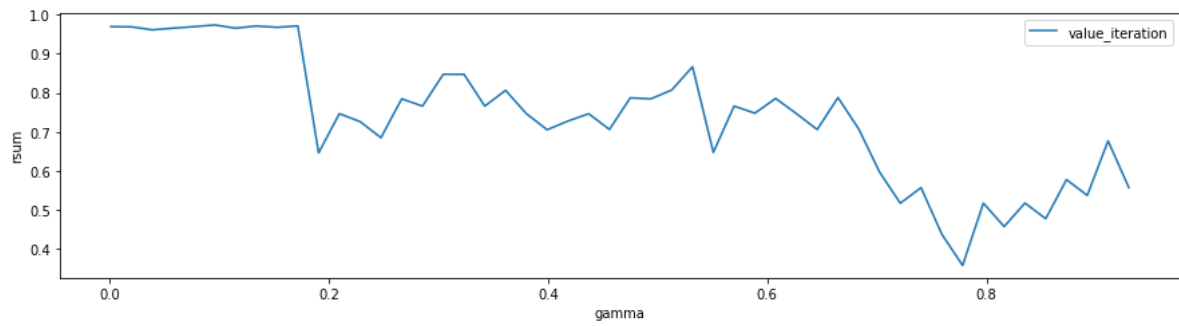
```

Plan 0

```

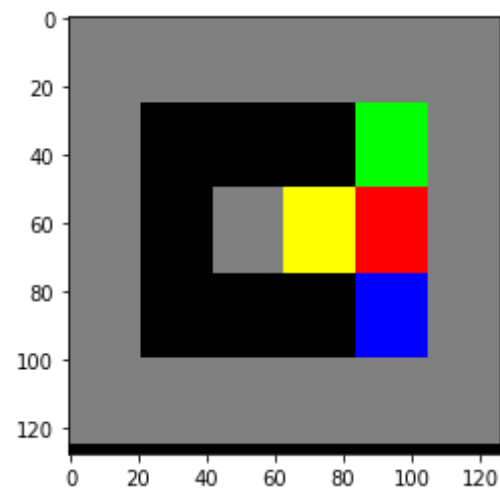
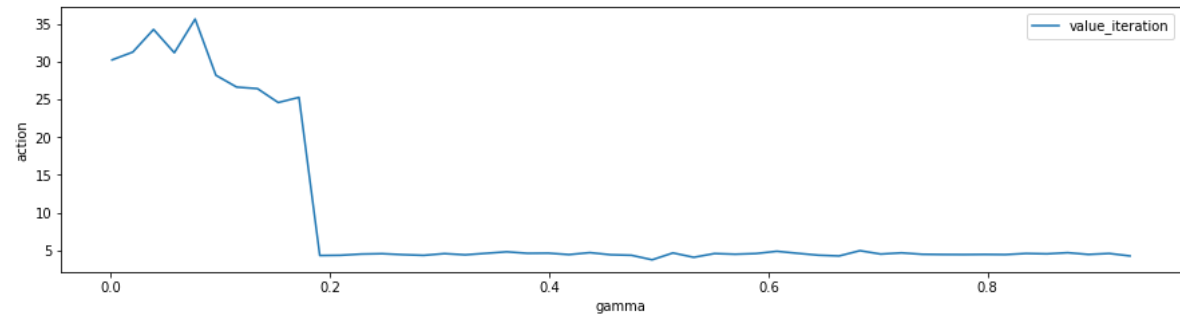
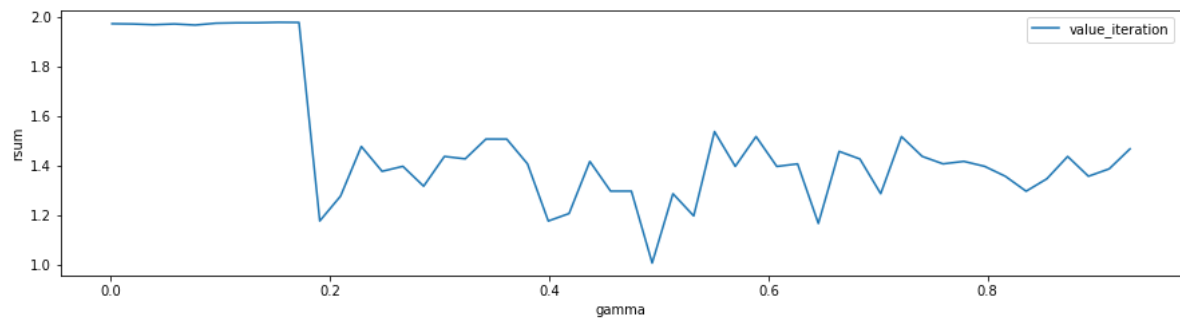
1  plotPlan(df_plan, 0)
2  createEnv(0)

```



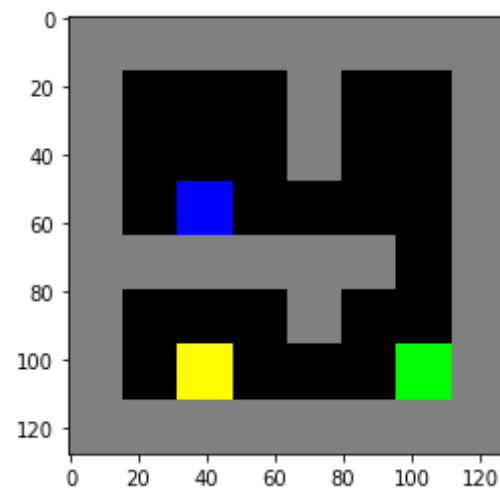
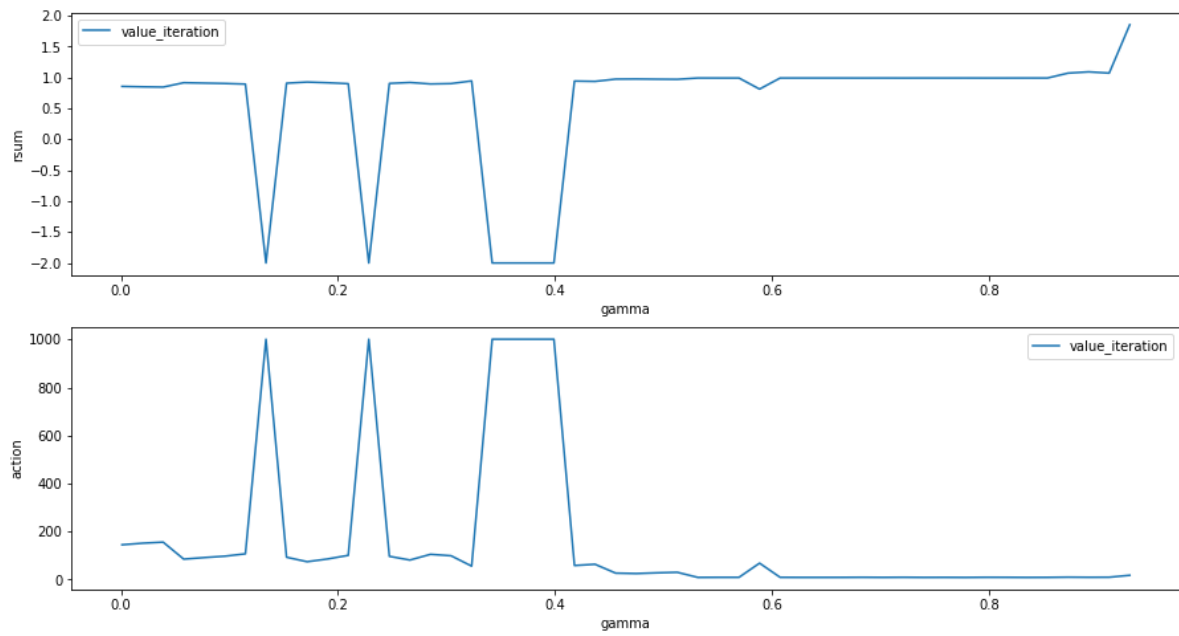
Plan 1

```
1 plotPlan(df_plan,1)
2 createEnv(1)
```



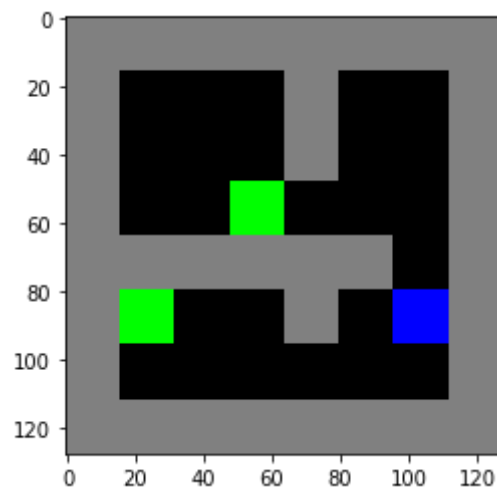
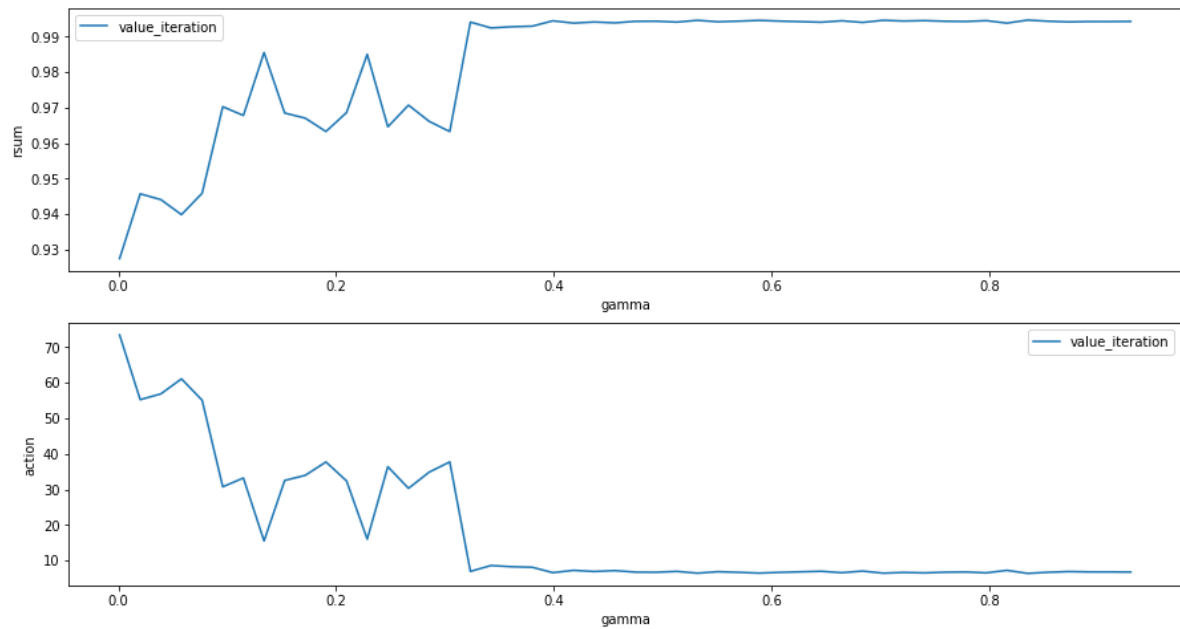
Plan 2

```
1 plotPlan(df_plan,2)
2 createEnv(2)
```



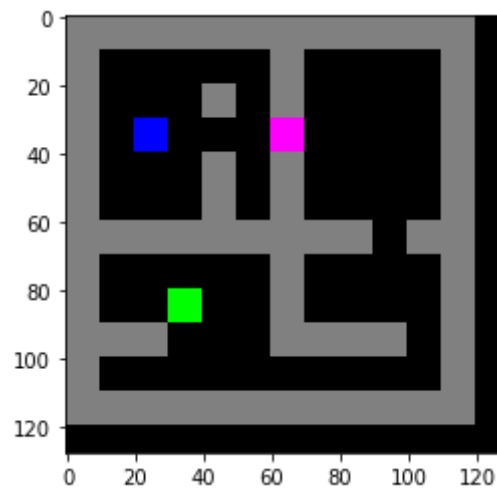
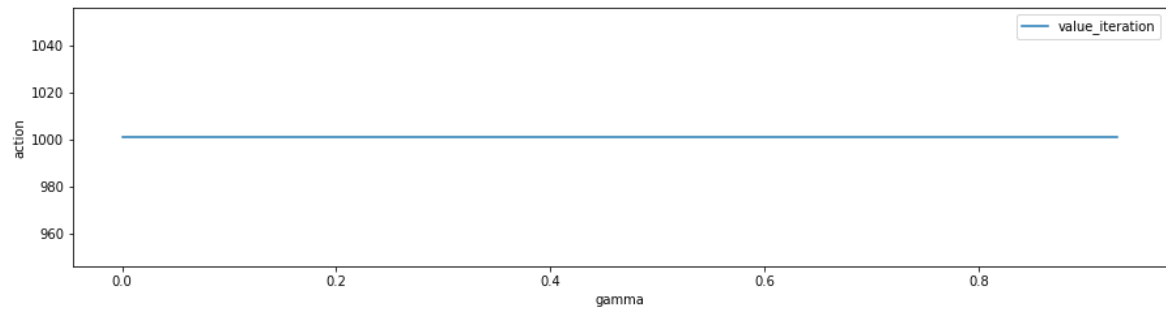
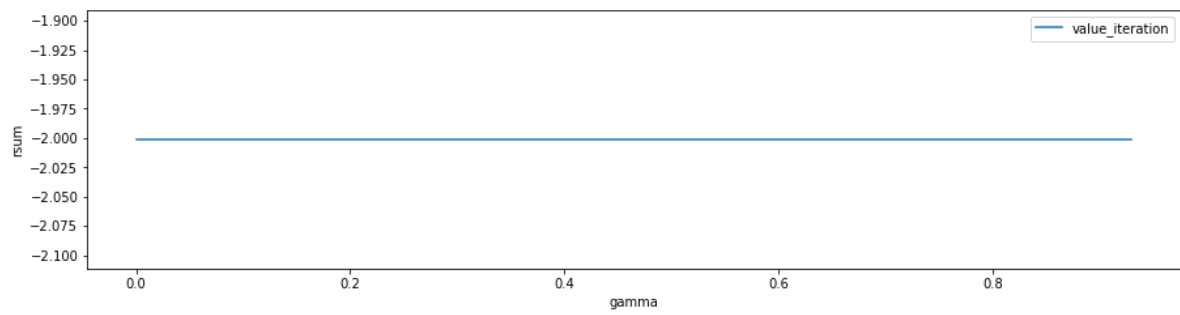
Plan 3

```
1 plotPlan(df_plan,3)
2 createEnv(3)
```



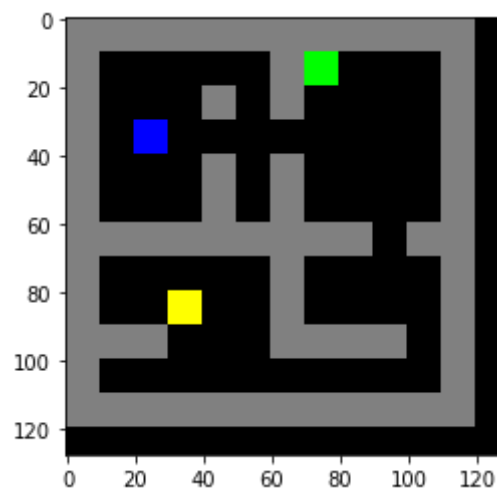
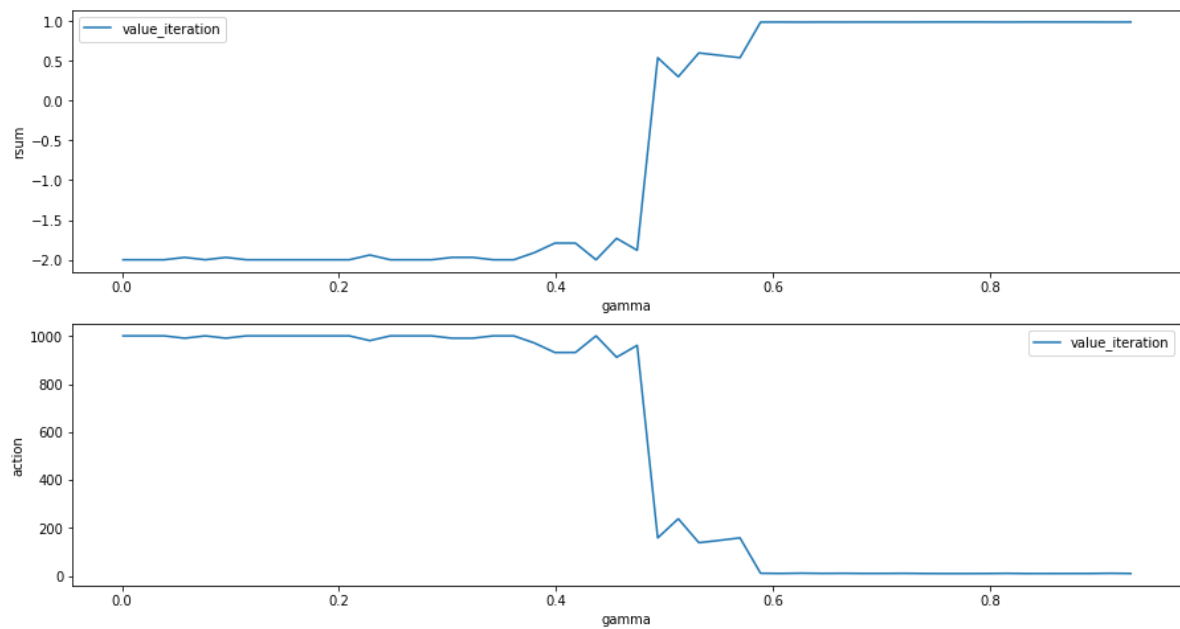
Plan 4

```
1 plotPlan(df_plan,4)
2 createEnv(4)
```



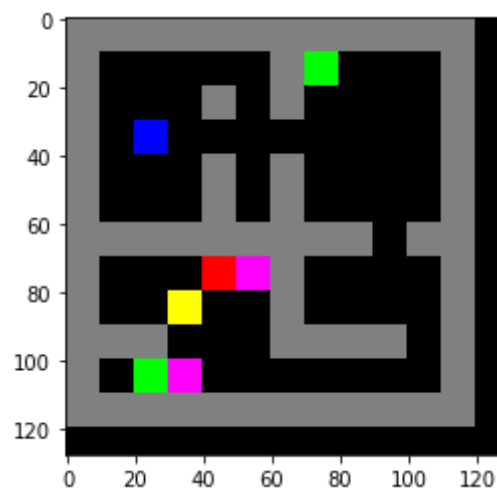
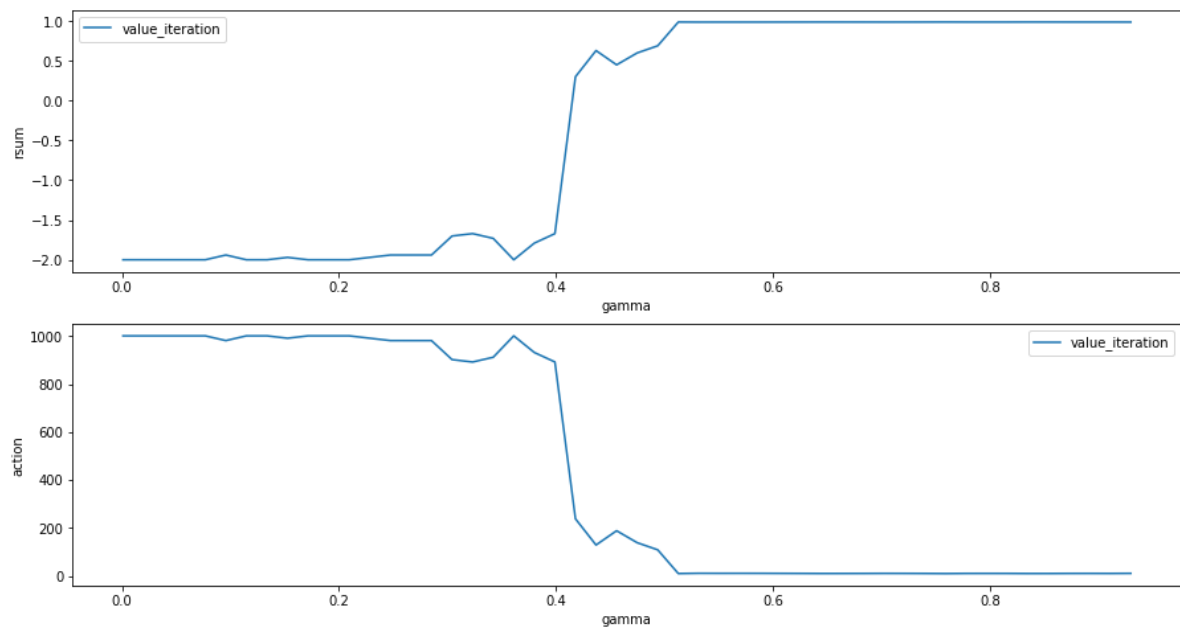
Plan 5

```
1 plotPlan(df_plan,5)
2 createEnv(5)
```



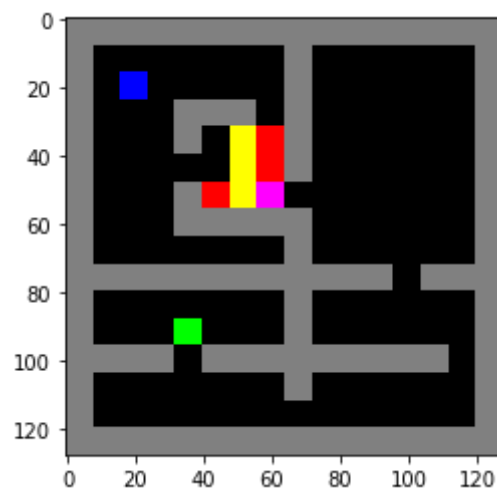
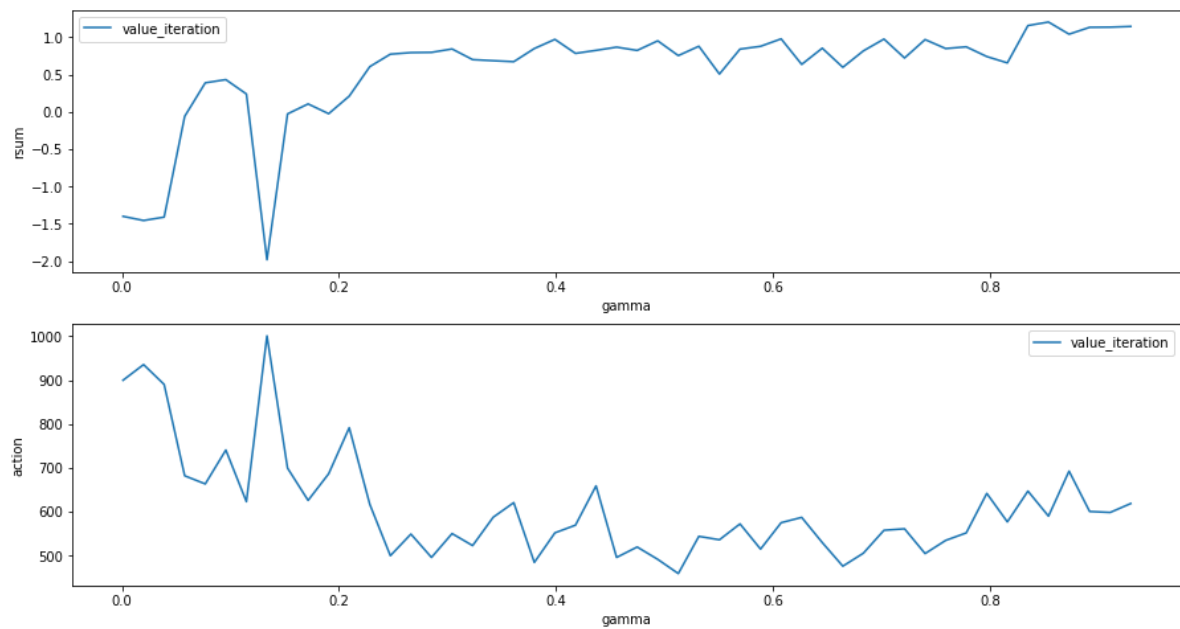
Plan 6

```
1 plotPlan(df_plan,6)
2 createEnv(6)
```



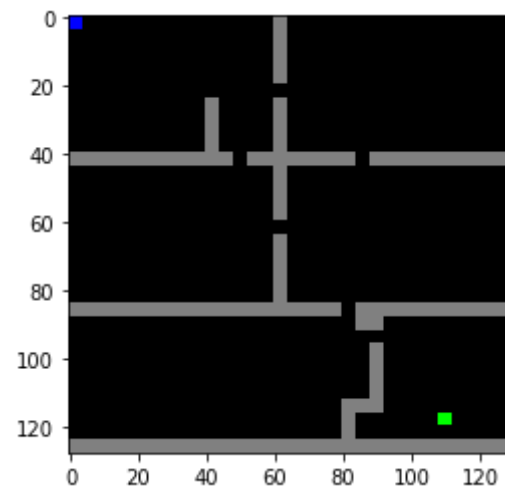
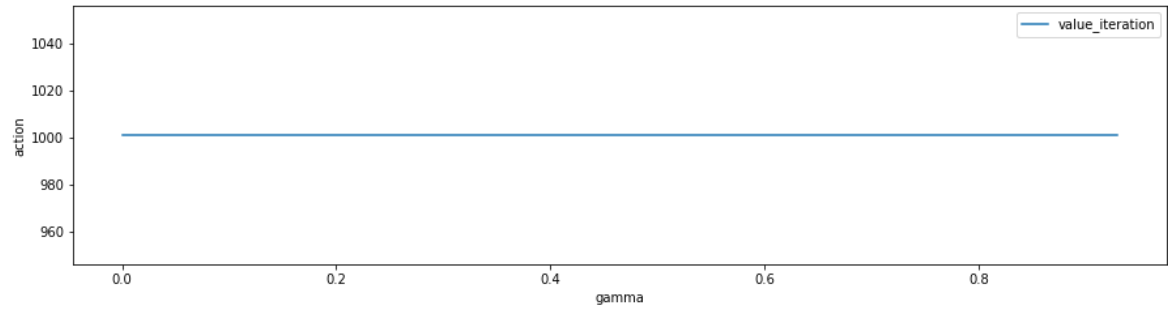
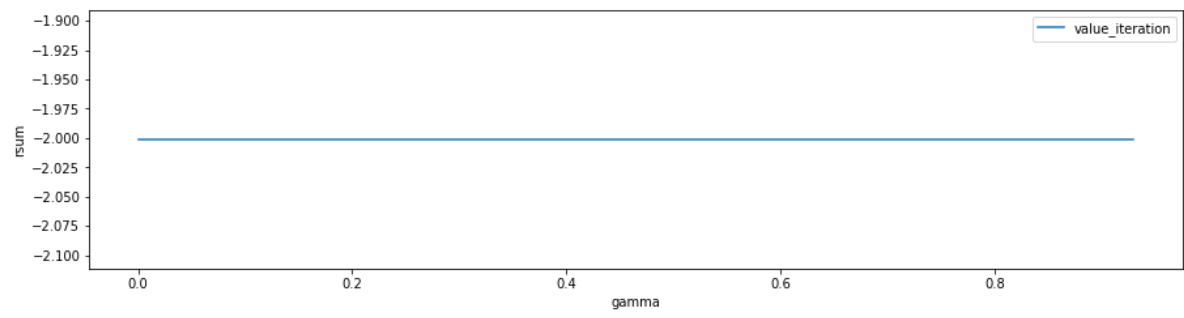
Plan 7

```
1 plotPlan(df_plan,7)
2 createEnv(7)
```

Plan 8

```
1 plotPlan(df_plan,8)
2 createEnv(8)
```



```
1 createEnv(10)
```

