

# Reinforcement learning – a1rl

André Filipe Frade Guerra

Outubro/2022

## Conteúdo

Introdução.....	3
Tabuleiro de jogo .....	3
Classes e métodos.....	3
Exercício 1 .....	6
Resultados exercício 1.....	7
Exercício 2 .....	8
Resultados 2a .....	9
Resultados 2b.....	10
Exercício 3 .....	13
Resultados exercício 3.....	13
Exercício 4 .....	14
Resultados exercício 4.....	14

# Introdução

Neste relatório irá ser explicado num resumo o funcionamento do tabuleiro de jogo, as classes criadas, funcionalidade dos métodos criados e por fim a execução e discussão dos resultados dos exercícios.

## Tabuleiro de jogo

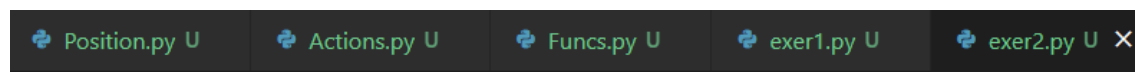
O tabuleiro de jogo é constituído por uma matriz 12x12 que contem em cada posição o objeto “Position” com os seguintes atributos:

2	<code>class Position:</code>	state é o número identificador da posição.
3	<code>state = 0</code>	parede é um booleano que indica se a posição é uma parede.
4	<code>parede = False</code>	
5	<code>reward = 0.0</code>	reward é o valor da recompensa dessa posição.

O tabuleiro contém nas laterais e rodeado, posições com a característica parede = TRUE, fazendo com que seja impossível ir para essa posição.

## Classes e métodos

Para desenvolver os exercícios da primeira prática da disciplina de Introdução à Aprendizagem Automática foi desenvolvido um código em python usando o vsCode para realizar a mesma.



As seguintes classes foram criadas, sendo as 3 primeiras auxiliares para as seguintes, que são os exercícios resolvidos.

Começando pela classe “Funcs.py”, contem métodos que são comuns nos exercícios:

```
6 def Variance(data, ddof=0):
7     n = len(data)
8     mean = sum(data) / n
9     return sum((x - mean) ** 2 for x in data) / (n - ddof)
10
11
12 def Stdev(data):
13     var = Variance(data)
14     std_dev = math.sqrt(var)
15     return std_dev
```

Métodos usados para calcular o desvio padrão.

```

17 def BuildMatrix(matrix, cols, rows, reward):
18     j = 1
19     for i in range(rows):
20         for w in range(cols):
21             if i > 0 and i < 11 and w > 0 and w < 11:
22                 p = Position()
23                 p.state = j
24                 p.reward = 0.0
25                 p.parede = False
26                 matrix[i][w] = p
27             else:
28                 p = Position()
29                 p.state = j
30                 p.reward = 0.0
31                 p.parede = True
32                 matrix[i][w] = p
33         j += 1
34     matrix[10][10].reward = reward

```

Método que cria uma matriz 12x12 em que as linhas e colunas das extremidades têm a propriedade parede = True.

Para cada posição é definido um número.

Por fim é atribuído à posição de destino o valor da recompensa.

```

50 def RandomAction():
51     r = randrange(1,5)
52     if r == 1:
53         return Action.UP
54     elif r == 2:
55         return Action.DOWN
56     elif r == 3:
57         return Action.LEFT
58     elif r == 4:
59         return Action.RIGHT

```

Método que define aleatoriamente a próxima ação.

```

61 def StateTransition(currentState, action):
62     nextState = 0
63     if action == Action.UP:
64         nextState = currentState - 12
65     elif action == Action.DOWN:
66         nextState = currentState + 12
67     elif action == Action.RIGHT:
68         nextState = currentState + 1
69     elif action == Action.LEFT:
70         nextState = currentState - 1
71     return nextState

```

Método que retorna a nova posição dependendo da ação realizada.

nota: tendo em conta que é uma matriz 12x12 as ações subir e descer subtraem e adicionam 12 à posição atual.

```

73 def IsParede(matrix, nextState):
74     return SearchStateInMatrix(matrix, nextState).parede
75
76 def SearchStateInMatrix(matrix, state):
77     for row in matrix:
78         for element in row:
79             if element.state == state:
80                 return element
81     return False

```

Método IsParede valida se a posição é uma parede.

Método SearchStateInMatrix retorna a posição através da pesquisa pelo número do estado.

```

83 def UpdateReward(matrix, currentState, nextState, alfaVar, gamaVar):
84     currentStateValue = SearchStateInMatrix(matrix, currentState).reward
85     nextStateValue = SearchStateInMatrix(matrix, nextState).reward
86     bestValueNextStepCanDo = GetBestValueAction(matrix, nextState)
87     if nextStateValue > 0.0:
88         #SearchStateInMatrix(matrix, currentState).reward = ( (1 - alfaVar) * currentStateValue ) + alfaVar * ( nextStateValue + ( gamaVar * bestValueNextStepCanDo ) )
89         #SearchStateInMatrix(matrix, currentState).reward = ( currentStateValue + ( alfaVar * ( 1 + ( gamaVar * bestValueNextStepCanDo ) ) ) - currentStateValue )
90         SearchStateInMatrix(matrix, currentState).reward = (1 - alfaVar) * currentStateValue + alfaVar * (nextStateValue + gamaVar * bestValueNextStepCanDo)

```

Método UpdateReward atualiza a recompensa da posição atual tendo em conta a posição para onde vai a seguir.

```

92 def GetBestValueAction(matrix, currentState):
93     upStateReward = SearchStateInMatrix(matrix, StateTransition(currentState, Action.UP)).reward
94     downStateReward = SearchStateInMatrix(matrix, StateTransition(currentState, Action.DOWN)).reward
95     leftStateReward = SearchStateInMatrix(matrix, StateTransition(currentState, Action.LEFT)).reward
96     rightStateReward = SearchStateInMatrix(matrix, StateTransition(currentState, Action.RIGHT)).reward
97
98     return max(upStateReward, downStateReward, leftStateReward, rightStateReward)
99
100 def GetBestStepAction(matrix, currentState):
101     upState = SearchStateInMatrix(matrix, StateTransition(currentState, Action.UP))
102     downState = SearchStateInMatrix(matrix, StateTransition(currentState, Action.DOWN))
103     leftState = SearchStateInMatrix(matrix, StateTransition(currentState, Action.LEFT))
104     rightState = SearchStateInMatrix(matrix, StateTransition(currentState, Action.RIGHT))
105
106     if upState.reward == downState.reward == leftState.reward == rightState.reward:
107         return StateTransition(currentState, RandomAction())
108
109     aux = max(upState.reward, downState.reward, leftState.reward, rightState.reward)
110     if upState.parede == False and upState.reward == aux:
111         return upState.state
112     elif downState.parede == False and downState.reward == aux:
113         return downState.state
114     elif leftState.parede == False and leftState.reward == aux:
115         return leftState.state
116     elif rightState.parede == False and rightState.reward == aux:
117         return rightState.state
118

```

Método GetBestValueAction retorna a recompensa mais alta tendo em conta a atual posição.

Método GetBestStepAction retorna o estado para o qual a atual posição deve ir tendo em conta o valor da recompensa mais alta. Caso as recompensas em volta sejam todas iguais, é retornado um estado aleatório.

```

116 def GetMatrixOnlyRewards(matrix,cols,rows):
117     matrixAux = [[0 for _ in range(cols)] for _ in range(rows)]
118     for i in range(rows):
119         for w in range(cols):
120             matrixAux[i][w] = matrix[i][w].reward
121     return matrixAux
122

```

Método que retorna uma matriz apenas com os valores das recompensas para construir o heatmap.

```

123 def Test1000(initState, finalState, steps, FAIL_FUN, matrix, mainTestStepsCount,reward):
124     currentState = initState
125     rewards = []
126     countSteps = 0
127     while currentState != finalState:
128         if countSteps >= steps:
129             countSteps = FAIL_FUN
130             break
131
132         nextState = GetBestStepAction(matrix, currentState)
133         if IsParede(matrix, nextState) == False:
134             reward = SearchStateInMatrix(matrix, currentState).reward
135             try:
136                 rewards.append(steps / reward)
137             except ZeroDivisionError:
138                 rewards.append(0)
139             currentState = nextState
140             countSteps += 1
141
142     #print("Avarage reward per "+str(steps)+" steps at " +str(mainTestStepsCount)+ " => " + str( rewards ) )
143     return rewards

```

Método que testa a alínea 2a.

# Exercício 1

Presente no ficheiro exer1.py, este ficheiro começa por declarar vários valores iniciais que definem o jogo, como cálculos, runs, steps e o tamanho da matriz.

```
10 reward = 1000
11 FAIL_FUN = "FAIL RUN"
12 alfaVar = 0.7
13 lambdaVar = 0.99
14 #final state
15 finalState = 131
16 #initial state
17 initialState = 40
18 currentState = initialState
19 #size matrix/gameboard
20 rows = 12
21 cols = 12
22
23 runs = 30
24 steps = 1000
25
26 matrix = [[0 for _ in range(cols)] for _ in range(rows)]
27
28 BuildMatrix(matrix,cols,rows, reward)
```

```
35 for r in range(runs):
36     currentState = initialState
37     countSteps = 0
38
39     st = time.time()
40     while currentState != finalState:
41         if countSteps >= steps:
42             countSteps = FAIL_FUN
43             break
44
45         nextState = StateTransition(currentState, RandomAction())
46         if IsParede(matrix, nextState) == False:
47             currentState = nextState
48             countSteps += 1
49
50     et = time.time()
51     runTimePerRun.append(et - st)
52     numStepsPerRun.append(countSteps)
```

Este pedaço de código faz correr a simulação do exercício 1 tendo em conta os runs e os steps por run.

Confirma se a próxima ação é uma parede, caso seja, volta a tentar outra ação aleatória.

Caso chega ao máximo de steps é declarada uma FAIL RUN.

```
57 numStepsPerRunClean = []
58 numRunsReachEnd = 0
59 somaTimeOfRun = 0
60 somaStepsOfRuns = 0
61 j = 0
62 for i in numStepsPerRun:
63     if i != FAIL_FUN:
64         numStepsPerRunClean.append(i)
65         numRunsReachEnd += 1
66         somaTimeOfRun += runTimePerRun[j]
67         somaStepsOfRuns += i
68     j += 1
69
70 numberStepsPerRunMRPT = []
71 for i in numStepsPerRunClean:
72     numberStepsPerRunMRPT.append(100/i)
```

Código auxiliar para depois calcular as métricas pretendidas.

```

77 print("average reward per step in "+str(steps)+" steps => "+ str((numRunsReachEnd*100)/runs) + "%")
78 print("standard deviation steps => " + str(Stdev(numStepsPerRunClean)))
79
80 print("average time of runs "+str(steps)+" steps => "+ str(somaTimeOfRun/numRunsReachEnd))
81 print("standard deviation time => " + str(Stdev(runTimePerRun)))
82
83 print("Average number actions per run => " + str(somaStepsOfRuns/numRunsReachEnd))
84 print("standard deviation time => " + str(Stdev(runTimePerRun)))
85
86 fig, axs = plt.subplots(1,3)
87 axs[0].boxplot(numberStepsPerRunMRPT)
88 axs[0].set_title('Mean reward per step')
89 axs[1].boxplot(numStepsPerRunClean)
90 axs[1].set_title('Mean number of steps')
91 axs[2].boxplot(runTimePerRun)
92 axs[2].set_title('Mean execution time (s)')
93
94 plt.show()

```

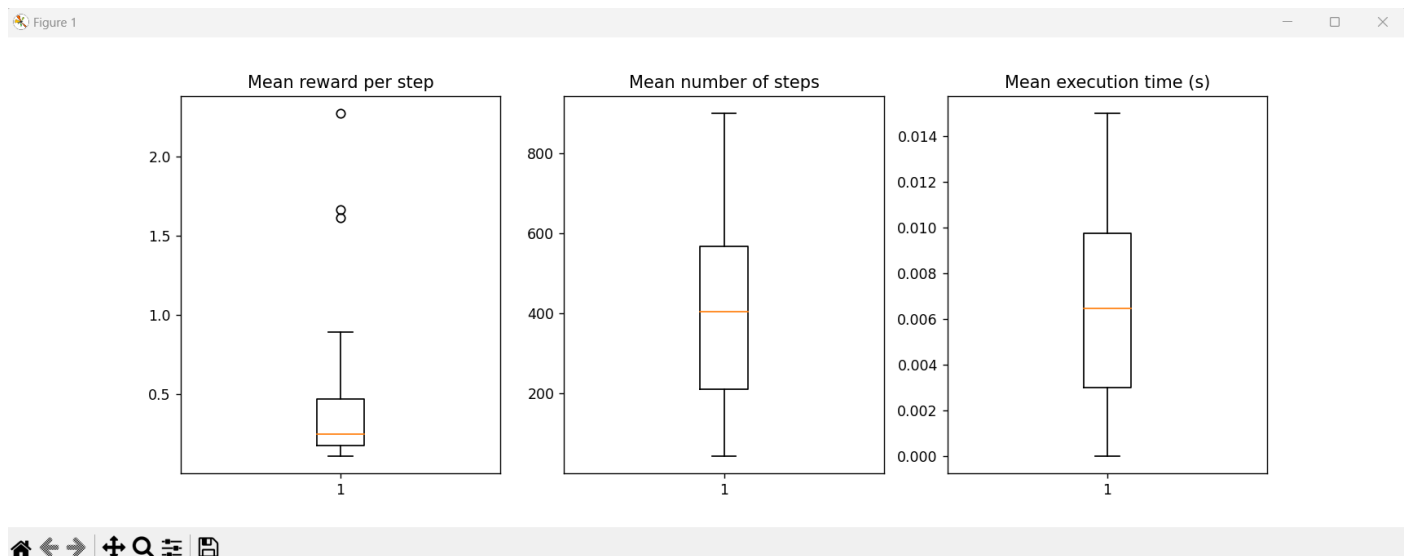
Código que mostra os cálculos pretendidos no exercício, como também no final as boxplots pretendidas.

## Resultados exercício 1

```

average reward per step in 1000 steps => 83.33333333333333%
standard deviation steps => 234.51210970864597
average time of runs 1000 steps => 0.0052825927734375
standard deviation time => 0.004025386729196857
Average number actions per run => 393.52
standard deviation time => 0.004025386729196857

```



## Exercício 2

O exercício dois é uma cópia do primeiro exercício, com algumas alterações no pedaço de código que corre a simulação e os resultados. Os valores iniciados na classe também mudaram de maneira que agora os steps têm um valor de 20.000.

No código que corre a simulação foi adicionado o método UpdateReward que atualiza as recompensas a cada ação feita, e o Test1000 que corre a simulação a cada X steps, neste caso nos steps 100, 200, 500, 600, 700, 800, 900, 1000, 2500, 5000, 7500, 10000, 12500, 15000, 17500 e 20000.

```
32 numStepsPerRun = []
33 runTimePerRun = []
34
35 for r in range(runs):
36     currentState = initialState
37     countSteps = 0
38
39     st = time.time()
40     while currentState != finalState:
41         if countSteps >= steps:
42             countSteps = FAIL_FUN
43             break
44
45         nextState = StateTransition(currentState, RandomAction())
46         if IsParede(matrix, nextState) == False:
47             UpdateReward(matrix, currentState, nextState, alfaVar, gamaVar)
48             currentState = nextState
49             countSteps += 1
50
51         if stepsToStop.__contains__(countSteps):
52             rewards = Test1000(initialState, finalState, 1000, FAIL_FUN, matrix, countSteps, reward)
53             #plt.plot(rewards)
54             #plt.show()
55
56     #m = GetMatrixOnlyRewards(matrix, cols, rows)
57     #ax = sns.heatmap(m, cbar=False)
58     #plt.show()
59
60     et = time.time()
61     runTimePerRun.append(et - st)
62     numStepsPerRun.append(countSteps)
63
```

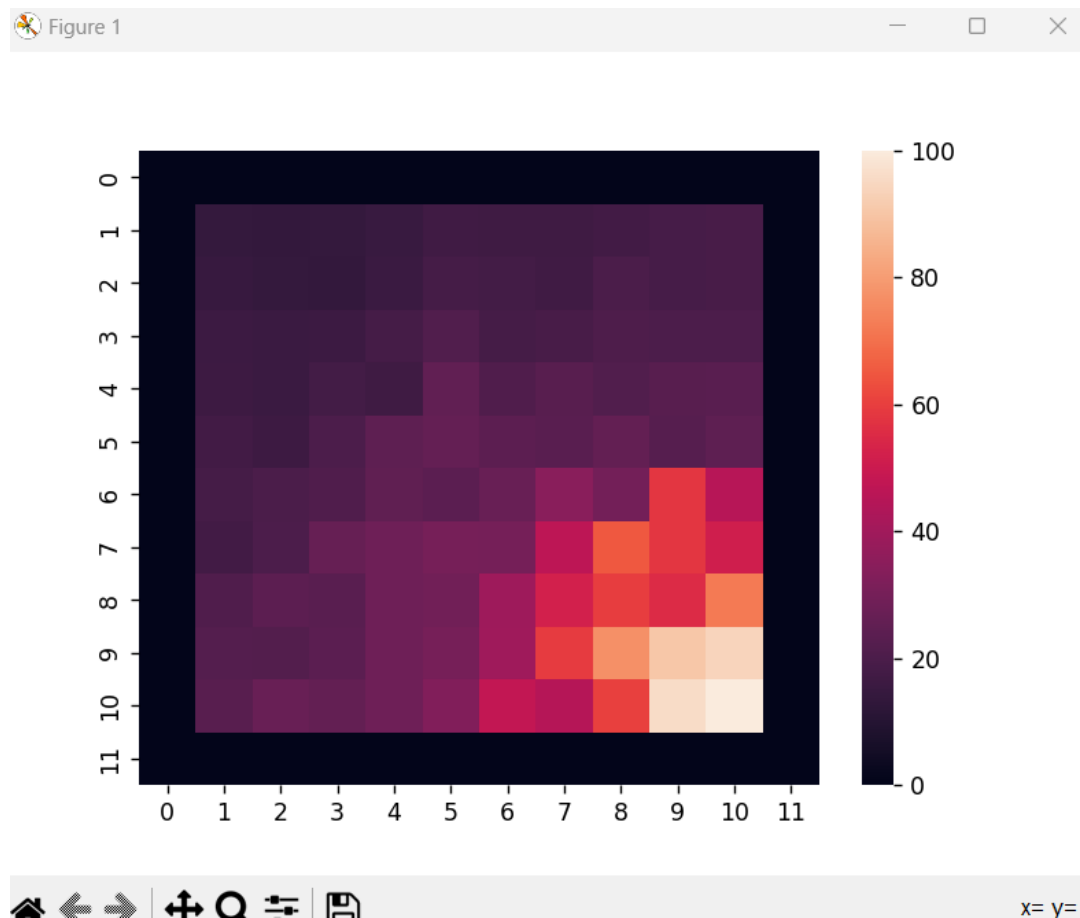
```
66 numStepsPerRunClean = []
67 numRunsReachEnd = 0
68 somaTimeOfRun = 0
69 somaStepsOfRuns = 0
70 j = 0
71 for i in numStepsPerRun:
72     if i != FAIL_FUN:
73         numStepsPerRunClean.append(i)
74         numRunsReachEnd += 1
75         somaTimeOfRun += runTimePerRun[j]
76         somaStepsOfRuns += i
77     j += 1
78
79 print("average time of runs "+str(steps)+" steps => " + str(somaTimeOfRun/numRunsReachEnd))
80 print("standard deviation time => " + str(Stdev(runTimePerRun)))
81
82 m = GetMatrixOnlyRewards(matrix, cols, rows)
83 ax = sns.heatmap(m)
84 plt.show()
```

Código auxiliar para calcular os valores pedidos e depois o gráfico do heatmap.



## Resultados 2a

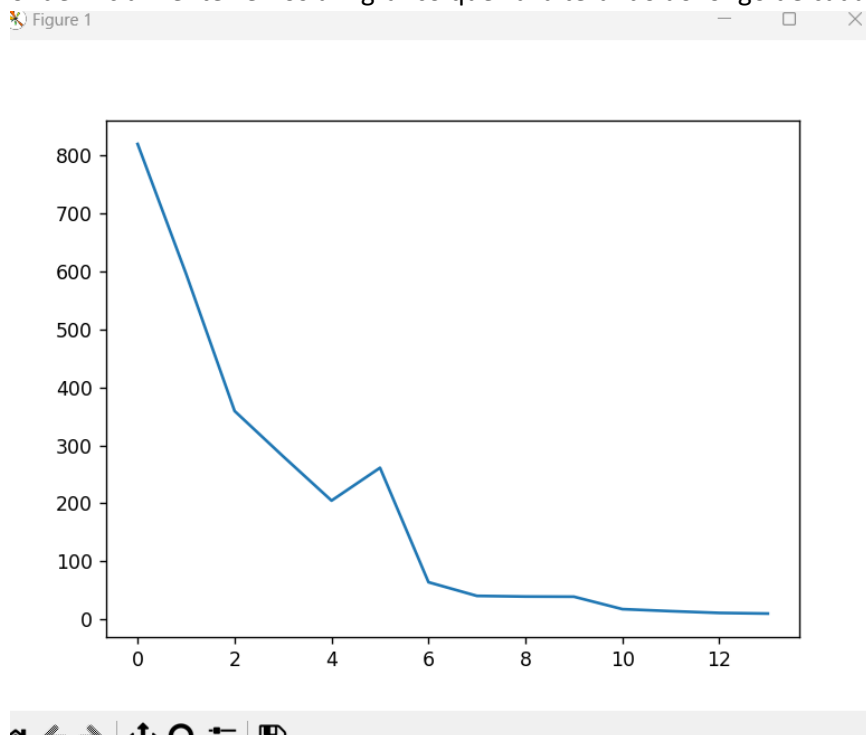
```
average time of runs 20000 steps => 0.1094772736231486  
standard deviation time => 0.11016083910338276
```

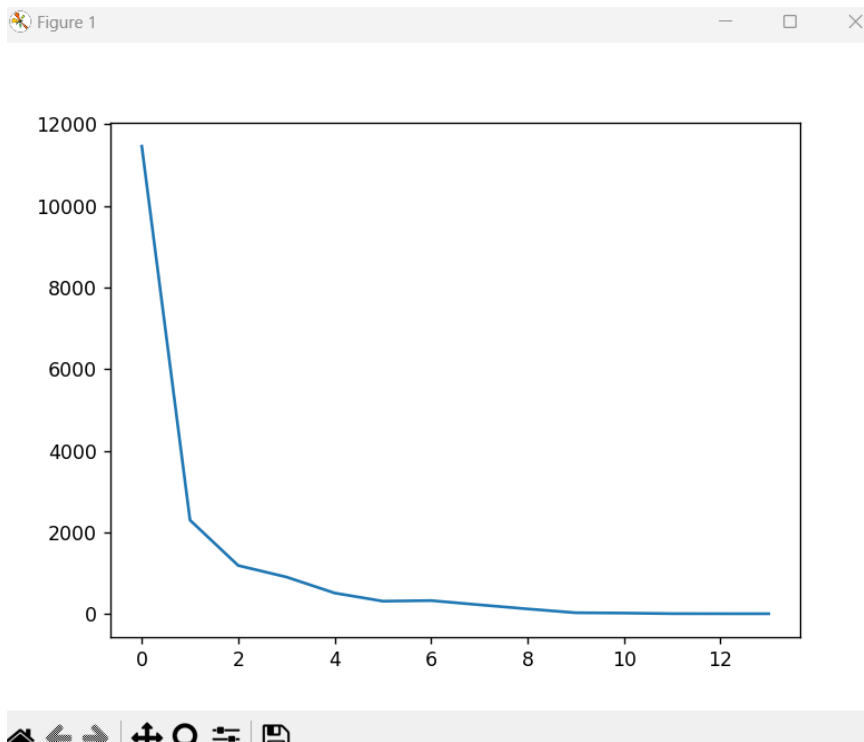


Resultado do heatmap com os valores das recompensas de cada posição.

Podemos observar este mapa a ser construído a cada run descomentando as linhas 57, 58 e 59.

Também podemos observar um gráfico para cada teste nos steps pretendidos descomentando as linhas 54 e 55, onde inicialmente vemos um gráfico que vai alterando ao longo de cada run:





## Resultados 2b

No exercício 2b a simulação foi alterada de maneira que escolha a melhor posição seguinte.

E para se obter resultados validos os runs foram alterados para 40 e os steps para 300.

```
35 for r in range(runs):
36     currentState = initialState
37     countSteps = 0
38
39     st = time.time()
40     while currentState != finalState:
41         if countSteps >= steps:
42             countSteps = FAIL_FUN
43             break
44
45         #nextState = StateTransition(currentState, RandomAction()) #Exercicio 2a
46         nextState = GetBestStepAction(matrix, currentState)#Exercicio 2b
47         if IsParede(matrix, nextState) == False:
```

Figure 1

Resultado da primeira run.

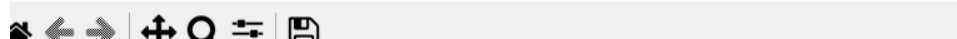
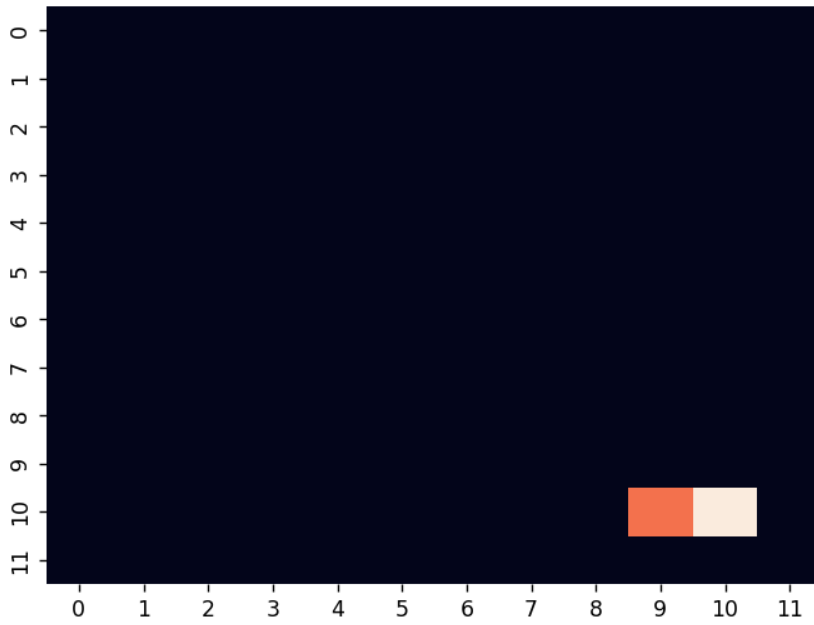


Figure 1

Resultado da segunda run.

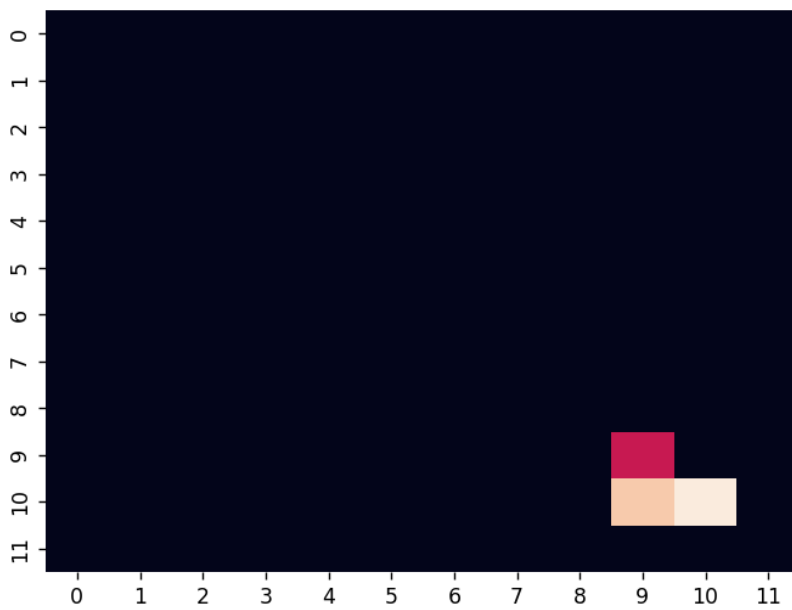
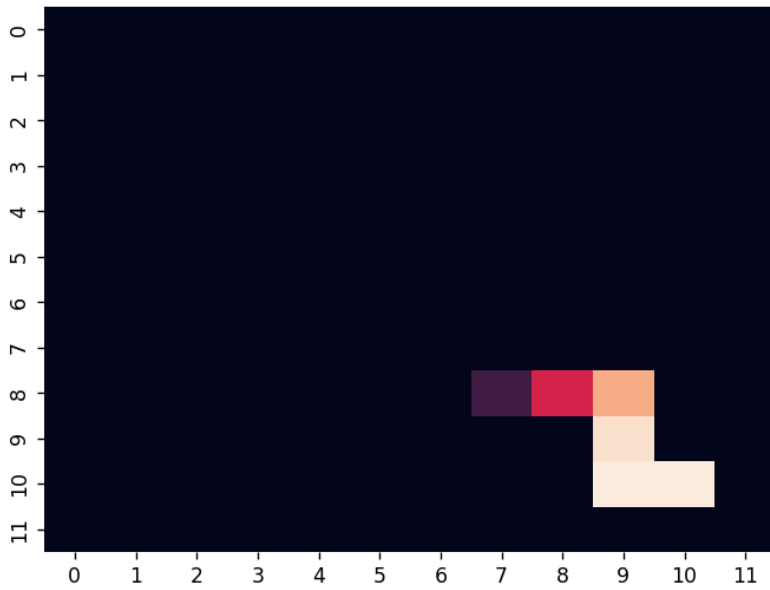


Figure 1

Resultado quinta run.



Podemos concluir que comparativamente à simulação onde não se escolhe a melhor opção seguinte, o resultado do melhor caminho possível é mais rápido a ser calculado.

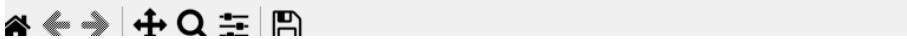
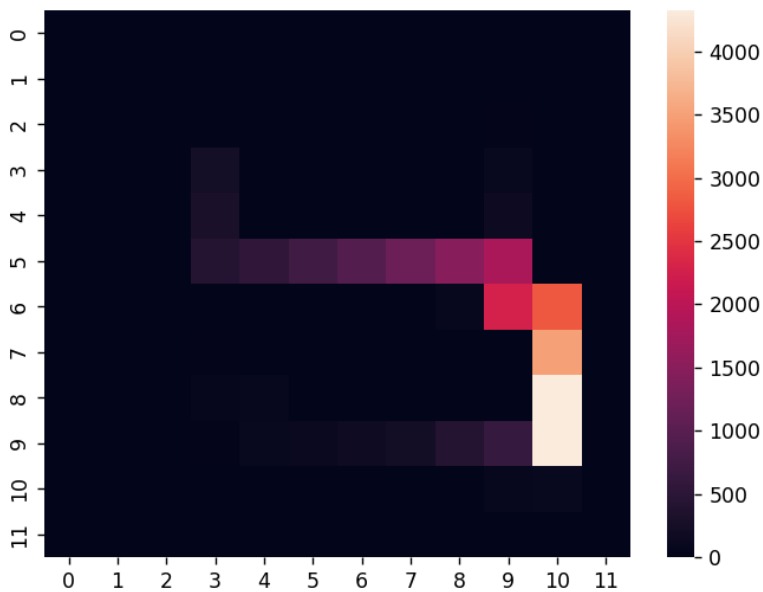


Figure 1

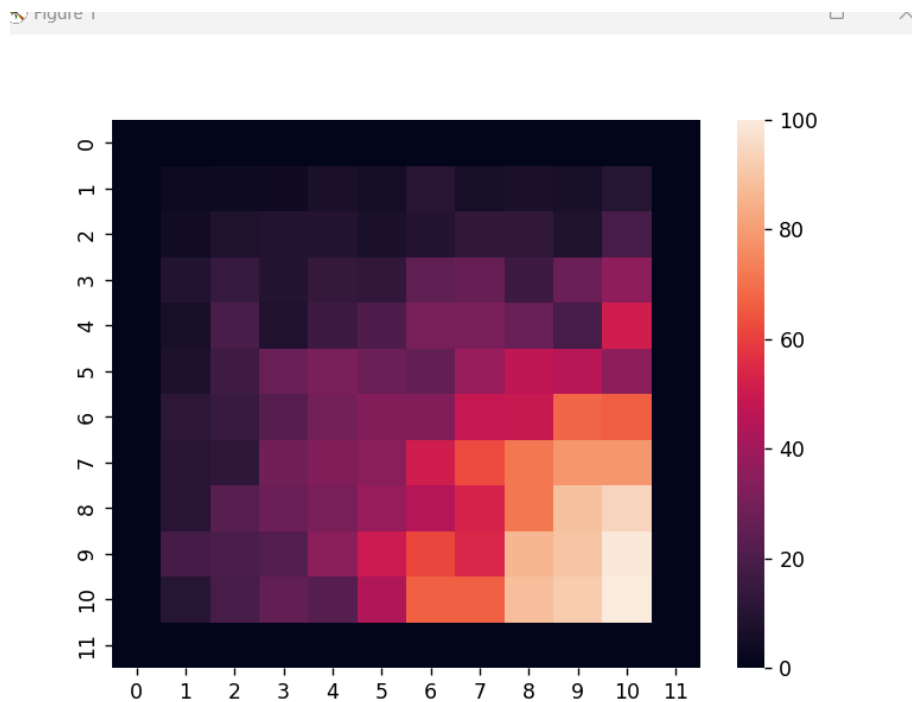
Resultado após os 40 runs com 300 steps cada uma.



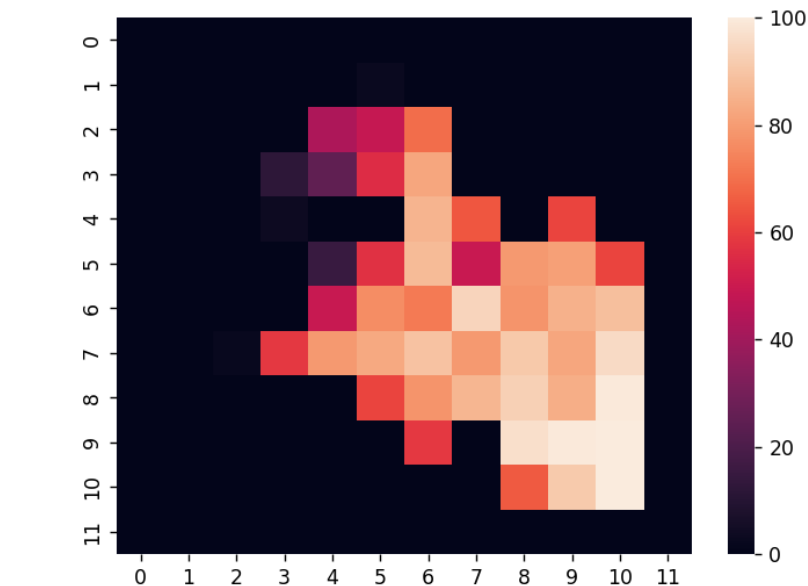
## Exercício 3

No exercício 3 foi adicionada na simulação o cálculo probabilístico que permite escolher se a ação é aleatória ou não:

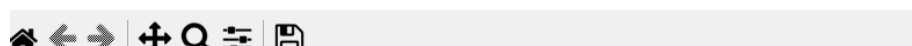
```
nextState = np.random.choice([StateTransition(currentState, RandomAction()), GetBestStepAction(matrix, currentState)], p=[0.9, 0.1])
```



Neste gráfico a probabilidade de chamarmos o método StateTransition é de 90%, enquanto o método GetBestStepAction tem apenas 10%.



Ao alterarmos as probabilidades, tendo o StateTransition 10% e o GetBestStepAction de 90% de probabilidade, podemos observar diferenças.



### Resultados exercício 3

Podemos concluir que quando o método que escolhe a melhor ação seguinte, tem mais probabilidade de acontecer, o caminho acaba por ser mais linear. Ao contrário, quando o método que escolhe uma ação aleatório, tem mais probabilidade de acontecer, o caminho fica menos linear dispersando-se.

## Exercício 4

No exercício 4 foi criado um método para construir as paredes no tabuleiro e dar-lhes recompensas negativas:

```
148 def BuildInsideWalls(matrix, rows, cols):
149     for i in range(rows-2):
150         matrix[i][4].parede = True
151         matrix[i][4].reward = -0.0000005
152     aux = cols - 1
153     for i in range(rows-2):
154         matrix[aux][7].parede = True
155         matrix[aux][7].reward = -0.0000005
156     aux -= 1
```

### Resultados exercício 4

Como podemos ver no heatmap as simulações chegaram ao caminho contornando as paredes criadas.

