

# Q-Learning Applied to a Simple Board Game

João Brito, M9984  
Department of Computer Science  
University of Beira Interior  
Covilhã, Portugal  
joao.pedro.brito{at}ubi.pt

**Abstract**—In this project, a reinforcement learning technique (Q-Learning) will be employed to help an agent navigate through a coloured board. Based on this environment, several experiments will be conducted to determine the impact of specific hyper-parameters on the agent’s behaviour.

## I. INTRODUCTION

Unlike Supervised Learning, where we have access to labelled instances and try to predict new ones, or Unsupervised Learning, where we try to highlight patterns in the data, Reinforcement Learning is very much based on giving an agent a positive or negative reward, based on how it performs. That type of learning is done by repeating tasks to build experience and, hopefully, avoid the same mistakes.

The present document will give an overview of a problem that can be solved using the aforementioned techniques and, in specific, how certain parameters affect the decision-making process of our agent.

## II. GAME OVERVIEW

Reinforcement Learning algorithms usually take place within a game environment. To that end, this work consisted in guiding an agent with only four allowed actions (go left, up, right or down) from a starting position to a finishing one in a  $N \times N$  board. As figure 1 shows, there are three types of cells:

- **blue**: the agent will want to traverse the board with these cells, which give a penalty of 1 lost point;
- **red**: these are the cells to be avoided and, depending on the experiment, they will either penalise the agent with 100 or 5 lost points;
- **yellow**: these represent the starting and finishing positions, with the latter rewarding the agent with 50 points.

Based on the description above, the agent’s goal is to minimise the penalties received during the traversal of the board (i.e. take the shortest path between the top yellow cell and the bottom counterpart, while avoiding the red cells as much as possible).

To avoid backtracking to the original position, the top yellow cell equates to a penalty of 50 lost points, encouraging the agent to never choose that cell, after starting the traversal.

Finally, section IV will have examples with either 30% or 50% of the total cells in a board marked as red.

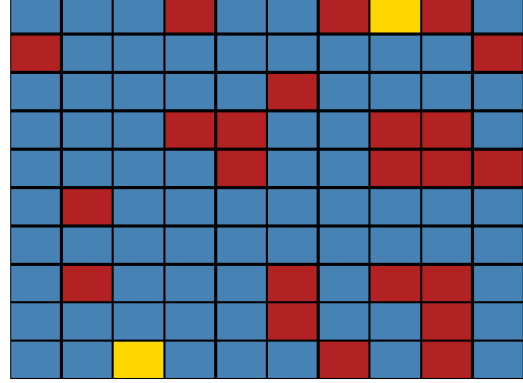


Fig. 1: Example of a 10x10 board with 30% red cells.

## III. OVERVIEW OF Q-LEARNING

To guide the agent, one of the better known algorithms is Q-Learning. Its simplicity and intuitive reasoning are key factors, especially in situations where we can store every state the agent can be in. Thus, there are 2 main data structures to keep in mind:

- **Q-Table**: this  $N \times A$  matrix (with  $N$  being the number of possible states and  $A$  the number of allowed actions in each state) represents the knowledge gathered by the agent at any time instance  $t$ . The purpose of Q-Learning is to build a Q-Table such that the agent can confidently choose actions that maximise the score at the end of the game.
- **Rewards Policy**: this structure is responsible for giving the agent the notion of immediate rewards by storing the reward/penalty of every possible state transition. It can be represented as a  $N^2 \times N^2$  matrix or, as chosen for this project, a directed graph with  $N$  nodes and bidirectional edges. The graph has the advantage of being incredibly more storage efficient when compared to a sparse matrix like the one mentioned (we would have  $N^2$  columns for each line but only 2, 3 or 4 cells of those columns would have meaningful values). The sparse nature of the matrix comes from the fact that each state only has direct access to, at most, 4 other states (figure 1).

As an example, a 10x10 board (i.e. 100 states) equates to a 100x100 matrix (i.e. 10000 cells). Using Python, this would take 80112 bytes of memory, whereas a directed graph (from the NetworkX [1] package) with 100 nodes and 2 edges linking each node pair would only amount to approximately 192 bytes.

This algorithm is mainly focused on the balance between immediate rewards (from nearby states) and future rewards (estimates of what is to come). In addition, the agent has to consider the binomial "exploration-exploitation" (i.e. roam around the environment or stick to the present knowledge).

Below is the pseudocode for Q-Learning, as described above [2]:

---

**Algorithm 1: Q-Learning**

---

**input :**  $B$  ( $N \times N$  matrix with  $A$  actions allowed),  $R$  (directed graph with  $N$  nodes)  
**output:**  $Q$  ( $N \times A$  matrix with the knowledge obtained)  
 $Q(s, a) \leftarrow \text{initialise\_Q\_table}()$   
**while**  $\text{episode\_count} < \text{max\_episodes}$  **do**  
   $s \leftarrow \text{start\_state}$   
  **while**  $s \neq \text{end\_state}$  **do**  
     $a \leftarrow \text{using } \epsilon\text{-greedy, choose an action from } Q(s, a)$   
     $s' \leftarrow \text{perform action } a \text{ from state } s$   
     $Q(s, a) \leftarrow Q(s, a) + \alpha[R(s, s') + \gamma * \max_a Q(s', a) - Q(s, a)]$   
     $s \leftarrow s'$   
  **end**  
**end**

---

Some notes regarding the pseudocode above:

- **Q-Table initialisation:** if an action is impossible in a given state, initialise the correspondent position with *None* or some kind of flag to indicate just that.
- **Hyper-parameters:**  $\alpha$  is a small number that determines how much the Q-Table values are updated in each time iteration  $t$  (i.e. the commonly called learning rate);  $\gamma$ , usually selected from the range  $[0,1]$ , balances the immediate ( $R(s, s')$ ) and future ( $\max_a Q(s', a)$ ) rewards;  $\epsilon$  tells the agent whether to explore or exploit.
- **$\epsilon$ -greedy policy:** given that  $\epsilon$  is constrained in the range  $[0,1]$  and represents a probability, we should choose and action at random with probability  $\epsilon$  (exploration). Otherwise, we exploit the current knowledge of the Q-Table.

#### IV. EXPERIMENTS

In this section, a few experimental exercises will take place, with respect to the way changes in hyper-parameters affect the agent's abilities and the overall complexity of the game (for example, board size and percentage of red cells).

Note that the path taken by the agent will be highlighted in white.

##### A. Base Configuration

Both 10x10 boards in this subsection were solved with parameters that, during the tests made, seem to allow the agent to reach its goal with relative ease:

- $\alpha = 0.5$  with linear decay to 0.15.
- $\gamma = 0.1$  with linear growth to 0.4.
- $\epsilon = 1.0$  with linear decay to 0.1.

Q-Learning (10x10) -  $\alpha = 0.5 \dots 0.15$  |  $\gamma = 0.1 \dots 0.4$  |  $\epsilon = 1.0 \dots 0.1$

1	2	3	4	5	6	7	8	9	10
11	12	13	14	15	16	17	18	19	20
21	22	23	24	25	26	27	28	29	30
31	32	33	34	35	36	37	38	39	40
41	42	43	44	45	46	47	48	49	50
51	52	53	54	55	56	57	58	59	60
61	62	63	64	65	66	67	68	69	70
71	72	73	74	75	76	77	78	79	80
81	82	83	84	85	86	87	88	89	90
91	92	93	94	95	96	97	98	99	100

Score: 30/50 | 1000 Ep.

Fig. 2: First example of a successful run.

Figure 2 shows a successful case where the agent had to move somewhat far to the left in order to uncover the best path.

Q-Learning (10x10) -  $\alpha = 0.5 \dots 0.15$  |  $\gamma = 0.1 \dots 0.4$  |  $\epsilon = 1.0 \dots 0.1$

1	2	3	4	5	6	7	8	9	10
11	12	13	14	15	16	17	18	19	20
21	22	23	24	25	26	27	28	29	30
31	32	33	34	35	36	37	38	39	40
41	42	43	44	45	46	47	48	49	50
51	52	53	54	55	56	57	58	59	60
61	62	63	64	65	66	67	68	69	70
71	72	73	74	75	76	77	78	79	80
81	82	83	84	85	86	87	88	89	90
91	92	93	94	95	96	97	98	99	100

Score: 27/50 | 1000 Ep.

Fig. 3: Second example of a successful run.

In figure 3 we can also observe a successful case where the agent's path has 2 distinct curves that allowed it to avoid the red walls.

From the two examples shown so far we can derive the following conclusion: on its way to the goal cell, the agent moves in straight lines, to minimise any penalties, and only makes turns or curves when needed.

## B. Changes in Hyper-parameters

### 1) $\alpha$ - learning rate:

This hyper-parameter controls how quickly we make strides towards convergence. In other words, it determines how much the Q-Table is updated with respect to  $t$ . As a test, the following 2 figures show the same board and final score but with a big difference: in figure 4, the value of  $\alpha$  was set to a really low value and that alteration lead to the need of more episodes until convergence (about 4800 episodes); on the other hand, if we linearly grow  $\alpha$  as training progresses (an increase to just 0.1 and everything else remaining the same) the number of episodes needed is reduced to just 1000 (figure 5).

Q-Learning (10x10) -  $\alpha = 0.01$  |  $\gamma = 0.1 \dots 0.4$  |  $\epsilon = 1.0 \dots 0.1$

1	2	3	4	5	6	7	8	9	10
11	12	13	14	15	16	17	18	19	20
21	22	23	24	25	26	27	28	29	30
31	32	33	34	35	36	37	38	39	40
41	42	43	44	45	46	47	48	49	50
51	52	53	54	55	56	57	58	59	60
61	62	63	64	65	66	67	68	69	70
71	72	73	74	75	76	77	78	79	80
81	82	83	84	85	86	87	88	89	90
91	92	93	94	95	96	97	98	99	100

Score: 33/50 | 4800 Ep.

Fig. 4: Example with a low value of  $\alpha$ .

Q-Learning (10x10) -  $\alpha = 0.01 \dots 0.1$  |  $\gamma = 0.1 \dots 0.4$  |  $\epsilon = 1.0 \dots 0.1$

1	2	3	4	5	6	7	8	9	10
11	12	13	14	15	16	17	18	19	20
21	22	23	24	25	26	27	28	29	30
31	32	33	34	35	36	37	38	39	40
41	42	43	44	45	46	47	48	49	50
51	52	53	54	55	56	57	58	59	60
61	62	63	64	65	66	67	68	69	70
71	72	73	74	75	76	77	78	79	80
81	82	83	84	85	86	87	88	89	90
91	92	93	94	95	96	97	98	99	100

Score: 33/50 | 1000 Ep.

Fig. 5: Example with a growing value of  $\alpha$ .

This extreme example shows how small values for  $\alpha$  delay convergence and how a simple growth policy can help the agent during training. By taking a look at the Q-Table update formula in the Q-Learning pseudocode, we can see that  $\alpha$  is affecting the sum of immediate and future rewards. A smaller  $\alpha$  will cause the values in the Q-Table to be more cluttered

together and the agent can suffer from a fuzzier sense of direction (in figure 4 the final path has unnecessary turns). On the other hand, a bigger value will give clearer directions to the agent. Obviously, the value shouldn't be that much bigger, otherwise the agent will not learn properly.

### 2) $\gamma$ - balance between immediate and future rewards:

The hyper-parameter  $\gamma$  is really important because it can give the agent a short sight or a sense of long term direction. In contrast to every other example, the following figures have red cells with a penalty of just 5 lost points (instead of the traditional 100):

Q-Learning (10x10) -  $\alpha = 0.5 \dots 0.15$  |  $\gamma = 0.1 \dots 0.4$  |  $\epsilon = 1.0 \dots 0.1$

1	2	3	4	5	6	7	8	9	10
11	12	13	14	15	16	17	18	19	20
21	22	23	24	25	26	27	28	29	30
31	32	33	34	35	36	37	38	39	40
41	42	43	44	45	46	47	48	49	50
51	52	53	54	55	56	57	58	59	60
61	62	63	64	65	66	67	68	69	70
71	72	73	74	75	76	77	78	79	80
81	82	83	84	85	86	87	88	89	90
91	92	93	94	95	96	97	98	99	100

Score: 31/50 | 1000 Ep.

Fig. 6: Case where the agent opted to follow the blue cells.

Q-Learning (10x10) -  $\alpha = 0.5 \dots 0.15$  |  $\gamma = 0.8 \dots 1.0$  |  $\epsilon = 1.0 \dots 0.1$

1	2	3	4	5	6	7	8	9	10
11	12	13	14	15	16	17	18	19	20
21	22	23	24	25	26	27	28	29	30
31	32	33	34	35	36	37	38	39	40
41	42	43	44	45	46	47	48	49	50
51	52	53	54	55	56	57	58	59	60
61	62	63	64	65	66	67	68	69	70
71	72	73	74	75	76	77	78	79	80
81	82	83	84	85	86	87	88	89	90
91	92	93	94	95	96	97	98	99	100

Score: 33/50 | 1000 Ep.

Fig. 7: Case where the agent took a shortcut (a red cell) and ended up getting a better score.

The behaviour shown in figures 6 and 7 is really interesting: with a bigger value of  $\gamma$ , the agent realised that by breaking a wall and taking a penalty of 5 lost points, it would end up benefiting it in the long run, instead of just sticking to the familiar blue cells.

### C. Board Size

With this test, the main purpose was to add more complexity to the game: the board got bigger (15x15), more red cells were distributed at random (50%) and the 2 yellow cells were placed in opposite corners, maximising the length of the ideal path:

### REFERENCES

- [1] NetworkX: Network Analysis in Python. [Online] <https://networkx.github.io/>
- [2] Q-Table update formula. [Online] <https://subscription.packtpub.com/book/data/9781789345803/4/ch04lv11sec33/the-learning-parameters-alpha-gamma-and-epsilon>

Q-Learning (15x15) -  $\alpha = 0.5 \dots 0.15$  |  $\gamma = 0.1 \dots 0.4$  |  $\epsilon = 1.0 \dots 0.1$

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
16	17	18	19	20	21	22	23	24	25	26	27	28	29	30
31	32	33	34	35	36	37	38	39	40	41	42	43	44	45
46	47	48	49	50	51	52	53	54	55	56	57	58	59	60
61	62	63	64	65	66	67	68	69	70	71	72	73	74	75
76	77	78	79	80	81	82	83	84	85	86	87	88	89	90
91	92	93	94	95	96	97	98	99	100	101	102	103	104	105
106	107	108	109	110	111	112	113	114	115	116	117	118	119	120
121	122	123	124	125	126	127	128	129	130	131	132	133	134	135
136	137	138	139	140	141	142	143	144	145	146	147	148	149	150
151	152	153	154	155	156	157	158	159	160	161	162	163	164	165
166	167	168	169	170	171	172	173	174	175	176	177	178	179	180
181	182	183	184	185	186	187	188	189	190	191	192	193	194	195
196	197	198	199	200	201	202	203	204	205	206	207	208	209	210
211	212	213	214	215	216	217	218	219	220	221	222	223	224	225

Score: 17/50 | 1000 Ep.

Fig. 8: Example of a successful path on a bigger board.

Q-Learning (15x15) -  $\alpha = 0.5 \dots 0.15$  |  $\gamma = 0.1 \dots 0.4$  |  $\epsilon = 1.0 \dots 0.1$

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
16	17	18	19	20	21	22	23	24	25	26	27	28	29	30
31	32	33	34	35	36	37	38	39	40	41	42	43	44	45
46	47	48	49	50	51	52	53	54	55	56	57	58	59	60
61	62	63	64	65	66	67	68	69	70	71	72	73	74	75
76	77	78	79	80	81	82	83	84	85	86	87	88	89	90
91	92	93	94	95	96	97	98	99	100	101	102	103	104	105
106	107	108	109	110	111	112	113	114	115	116	117	118	119	120
121	122	123	124	125	126	127	128	129	130	131	132	133	134	135
136	137	138	139	140	141	142	143	144	145	146	147	148	149	150
151	152	153	154	155	156	157	158	159	160	161	162	163	164	165
166	167	168	169	170	171	172	173	174	175	176	177	178	179	180
181	182	183	184	185	186	187	188	189	190	191	192	193	194	195
196	197	198	199	200	201	202	203	204	205	206	207	208	209	210
211	212	213	214	215	216	217	218	219	220	221	222	223	224	225

Score: 15/50 | 1000 Ep.

Fig. 9: Another example of a great run on a 15x15 board.

Once again, the agent was able to go from one side to the other, while taking the shortest route. In spite of being larger and more complex, the board was generated and solved within the time frame of the previous tests.

### V. CONCLUSION

The project described in this document showed how a problem can be solved with Reinforcement Learning strategies, like Q-Learning. Moreover, it gave examples of how hyperparameter tweaking can alter the quality of the generated solutions and/or the number of episodes needed to converge to the optimal path. Despite being a toy exercise, it brought forward transcendent concepts in this sub-area of Machine Learning.