

Bio-inspired Intelligence and Learning for Aerospace Applications Report

Nico Voß - 4550749 - <https://github.com/NicoEVA/BioInspired>

August 31, 2021

1 Introduction

This report serves to reflect on the development and implementation of a bio-inspired algorithm for optimising control of a simplified path-finding exercise. To this extent the report will discuss the background and idea for the assignment, the theories and methods used and the creation of the resulting code.

It was decided to apply Q-Learning to a path-finding problem that represents a glider attempting to get to a target location. Firstly, a two dimensional model was created, in which the piloting agent could choose to go in either of the four cardinal directions. Later on, for the purpose of better representing the three dimensional nature of flight, it was extended to three dimensions by adding a third degree of freedom in the form of altitude.

One of the biggest issues with Q-Learning and self-learning algorithms in general is that the problem and environment need to be set up in such a way that the agent truly optimises for the desired behaviour, instead of simply cheating the system and finding valid but undesired solutions. In the light of the glider problem, this means motivating the agent to find an economic path to the target location, and thinking about the rewards and penalties to apply for e.g. violating restricted airspace or making smart use of thermals. To resemble a realistic problem it has been decided to make use of negative rewards, which will be alleviated when flying through thermals and heavily increase when violating restricted airspace. Termination statements are tied to attempted leaving of the environment and flying into restricted airspace, to prevent invalid behaviour and terminate moot runs. In the following section a short description of the Q-Learning method is given, followed by the mathematical background, the code implementation, the resulting output and the conclusion.

2 Q-Learning Background

Q-Learning is a type of reinforcement learning that consists of an agent operating in a well defined environment, being given a finite set of states, rewards and actions. By using a trial and error approach, the agent optimises its policy as it learns more about the environment. It is thus a model-free version of reinforcement learning[1].

The optimal decisions and their weights are stored in a Q-table containing the estimated sum of future rewards when choosing a specific action in the current state. This is also the source that is used for taking decisions on optimal actions, which is overruled by a random choice during training to improve especially the early exploration. Effectively, the Q-table is later the policy of the agent, which is initialised and modified during learning to create a means of extracting the best possible path from any point in the environment to the specified goal.

By having the agent complete a certain number of runs with specific settings for the learning rate (the speed with which entries in the Q-table change) and the likelihood of the random choice, the action space is explored by the agents with different behaviours, resulting in a final Q-table at the end of the learning run. An additional parameter is the discount factor, which influences the agents decision based on how much total reward can still be acquired when taking an action. All three values were kept constant in between the runs shown in Figure 1 to allow for a meaningful comparison.

In general these parameters can be adjusted to improve the agents learning process, by e.g. reducing the amount of random choices as the number of runs increases for better convergence. Similarly, the learning rate can be adjusted to either quickly or slowly react to changes, which has to be traded off against the possibility of reacting too fast to non-optimal solutions and delaying or preventing convergence. The discount factor is usually kept such that its influence is marginal, as the immediate action taken by the agent is usually of greater interest than what happens at a later point in time. For the exact values, please take a look at the supplied code.

3 Mathematical Foundation

The Q-Learning approach relies on a set of relatively simple but in combination potent formulas and constructs. The Q-table has previously been mentioned and logically takes the shape of the state space with each entry having another dimension equal to the number of possible actions at that point. Thus, assuming we have a state space with a dimension of $4 \times 4 \times 4$ and we can go in either of the 4 cardinal directions plus up/down, the Q-table will have a shape of $4 \times 4 \times 4 \times 6$.

The variables that define the learning behaviour have previously been mentioned. The randomisation variable Epsilon controls how often a random direction is taken instead of following the current best choice. The learning rate determines how fast the values in the Q-table denoting the best choice are changed per run. Finally the discount factor puts a limiter on how much influence future rewards have on the decision process.

3.1 Temporal Difference

The Temporal Differences gives a method of calculating how much the Q-Value of a specific action taken in the previous state should be changed based on the learned knowledge of the agent from the current actions.

$$TD(s_t, a_t) = r_t + \gamma \cdot \text{Max}Q(s_{t+1}, a) - Q(s_t, a_t) \quad (1)$$

In this TD is the Temporal Difference, r_t is the reward achieved by taking an action in the previous state, γ is the discount factor in the range of zero to one, a measure of how important the future rewards are. It is usually kept relatively close to one, causing a smaller effect of future rewards[3].

MaxQ denotes the largest of the Q-values available for any action available in the current state, thus the largest predicted sum of future rewards. Q denotes the value for the action taken in the previous state.

3.2 Updating Q-values with Bellman Equation

By relying on both the old Q-values and the learned information after moving to the next state, the Q-values in the Q-table are updated from the Temporal Difference making use of a learning rate to tune the effects of how fast the ideal direction can be changed[2].

$$Q_{new}(s_t, a_t) = Q_{old}(s_t, a_t) + \alpha \cdot TD(s_t, a_t) \quad (2)$$

Here, we update the Q-Value of the Q-table of the taken action from the previous state Q_{old} by adding the product of the previously calculated Temporal Difference TD and the learning rate α . The resulting Q_{new} replaces the previous value. This also points out the necessity to initialise the Q-table with a set of random, uniform or pre-tuned values.

4 Environment Setup

To begin the development first an environment had to be generated that with its properties represents the situation of the glider. As such, it was decided to make use of negative rewards representing the sinking of the aircraft over time. Thermals were included to observe whether the agent makes use of areas that have lower penalties in favour of a direct path. Finally, restricted airspace was introduced, to add complexity and observe how the agent deals with limited decision space.

This space is defined in the code by three variables: the columns, rows and stacks. To explore the space, the initial position is randomly selected from the bounded box that is the state space. It is bounded by creating an enclosure of large negative rewards to prevent exceeding the limits of the space by tying a termination statement to said negative reward and its location. The termination on running into a set negative reward will be used to prevent flight into restricted airspace and disregard invalid routes that would otherwise continue to be evaluated. The target point is set manually and attention has to be paid that it does not lay within any of the restricted areas to be reachable by the agent.

To create the environment we make use of the following code:

```
1  ## Define Extent of Environment and create actions/q-values
2  rows = 24
3  cols = 24
4  stacks = 10 #attempt 3D
5  q_vals = np.zeros((rows,cols,stacks,6))
6  choices = np.zeros((rows,cols,stacks))
7
8  # Define Modifier areas
9  alt_rest1 = 5
10 restricted_airspace1 = [5,15,0,22]
11 thermals1 = [12,20,15,20]
12 thermals2 = [7,13,7,18]
13
14 ## Set up environment boundaries and varying negative rewards
15 actions = ['up','right','down','left',"3ddown","3dup"]
16 rewards = np.full((rows,cols,stacks),-1)
17 for z in range(stacks-1):
18     for x in range(cols-1):
19         rewards[x,0,z]=-100
20         rewards[x,rows-1,z]=-100
21     for y in range(rows-1):
22         rewards[0,y,z]=-100
23         rewards[cols-1,y,z]=-100
24     for x in range(thermals1[0],thermals1[1]):
25         for y in range(thermals1[2],thermals1[3]):
26             rewards[x,y,z] = -0.2
27     for x in range(thermals2[0],thermals2[1]):
28         for y in range(thermals2[2],thermals2[3]):
29             rewards[x,y,z] = -0.2
```

```

30 for x in range(cols-1):
31     for y in range(rows-1):
32         for z in range(0,stacks-1):
33             rewards[x,y,z] = -100
34
35 # establish restricted airspace
36 for z in range(alt_rest1,stacks):
37     for x in range(int(restricted_airspace1[0]),
38                     restricted_airspace1[1]):
39         for y in range(restricted_airspace1[2],restricted_airspace1
40                         [3]):
41             rewards[x,y,z] = -100
42 #Defining where the goal/target point is
43 rewards[2,2,3]=100

```

In addition the previously established methods have to be implemented in the code:

```

1  ## Create required functions
2  def terminal(row,col,stack):
3      if rewards[row,col,stack]==-100:
4          return True
5      elif rewards[row,col,stack]==100:
6          return True
7      else:
8          return False
9
10 def startloc():
11     row = np.random.randint(rows)
12     col = np.random.randint(cols)
13     stack = np.random.randint(stacks)
14     while terminal(row,col,stack):
15         row = np.random.randint(rows)
16         col = np.random.randint(cols)
17         stack = np.random.randint(stacks)
18     return row,col,stack
19
20 def findact(row,col,stack,eps):
21     if np.random.random() < eps:
22         return np.argmax(q_vals[row,col,stack])
23     else:
24         return np.random.randint(6)
25
26 def newloc(row,col,stack,act):
27     newrow = row
28     newcol = col
29     newstack = stack
30     if actions[act]=="up" and row > 0:
31         newrow -= 1
32     elif actions[act] == 'right' and col < cols-1:
33         newcol += 1
34     elif actions[act] == 'down' and row < rows-1:
35         newrow += 1
36     elif actions[act] == 'left' and col > 0:
37         newcol -= 1
38     elif actions[act] == "3ddown" and stack > 0:
39         newstack -=1

```

```

40     elif actions[act] == "3dup" and stack < stacks -1:
41         newstack += 1
42     return newrow,newcol,newstack
43
44 def generatechoices(q_vals,choices):
45     for i in range(np.shape(q_vals)[0]):
46         for j in range(np.shape(q_vals)[1]):
47             for k in range(np.shape(q_vals)[2]):
48                 choice = np.argmax(q_vals[i,j,k])
49                 choices[i,j,k] = choice
50     return choices

```

In order to train the agent, a routine is set up to have it practice on the space for a limited number of episodes:

```

1  ## Set up and perform training Run
2  epsilon = 0.35
3  discount = 0.8
4  lrate    = 0.8
5  ep = 0
6  for episode in range(15000):
7      row,col,stack = startloc()
8      tries = 0
9      while not terminal(row,col,stack):
10         aind = findact(row,col,stack,epsilon)
11         orow, ocol, ostack = row,col,stack
12         row, col, stack = newloc(row,col,stack,aind)
13
14         rew = rewards[row,col,stack]
15         oq = q_vals[row,col,stack,aind]
16         tempdiff = rew + (discount*np.max(q_vals[row,col,stack]))-
17         oq
18
19         newq = oq +(lrate*tempdiff)
20         q_vals[orow,ocol,ostack,aind] = newq
21         if tries % 10 == 0:
22             print(tries)
23         tries +=1
24     print("Reached target in EP " + str(ep+1))
25     ep += 1
26 print("trained successfully")

```

The final important snippet is the function that makes the agent follow the learned actions from any chosen point in the environment:

```

1  def extrsp(srow,scol,sstack):
2      if terminal(srow,scol,sstack):
3          return []
4      else:
5          row, col, stack = srow, scol, sstack
6          sp = []
7          sp.append([row,col,stack])
8          print(sp)
9          while not terminal(row,col,stack):
10             aind = findact(row,col,stack,1.)
11

```

```

12         row, col, stack = newloc(row,col,stack,aind)
13         sp.append([row,col,stack])
14
15     return sp

```

5 Results of Runs

To visualise the behaviour of the agent, a three dimensional plot has been created, showing the ideal path from the selected start point. For this, a chosen starting point is used in a function that follows the optimal actions as dictated by the trained Q-table. In addition, a figure of the ideal directions for each state is created by plotting the vector field from the Q-values and the respective optimal actions. The generated track features additional marking for restricted airspace in red and thermals in green. Within the thermal columns the negative rewards have reduced penalties.

```

1  ## Set desired Origin and create Graphs based on found Path
2  data = extrsp(20,20,7)
3
4  # create agent trace
5  x = [xi[0] for xi in data]
6  y = [yi[1] for yi in data]
7  z = [zi[2] for zi in data]
8
9  fig = plt.figure(figsize=(15,15))
10 ax = fig.add_subplot(111, projection='3d')
11
12 ax.plot(x, y, z, label='Agent Trace')
13 ax.set_xlabel("X")
14 ax.set_ylabel("Y")
15 ax.set_zlabel("Z")
16 # create modifier visualisations
17 ceilingx = np.arange(restricted_airspace1[0],restricted_airspace1
18                      [1],1)
19 ceilingy = np.arange(restricted_airspace1[2],restricted_airspace1
20                      [3],1)
21
22 wally1 = np.zeros_like(ceilingx)
23 wally1.fill(restricted_airspace1[2])
24 wally2 = np.zeros_like(ceilingx)
25 wally2.fill(restricted_airspace1[3]-1)
26 w1x, w1y = np.meshgrid(ceilingx,wally1)
27 w2x, w2y = np.meshgrid(ceilingx,wally2)
28
29 Zwallsx = np.zeros((np.shape(ceilingx)[0],np.shape(wally1)[0]))
30 for len in range(np.shape(Zwallsx)[1]):
31     Zwallsx[len][:] = alt_rest1 + len
32
33 X,Y = np.meshgrid(ceilingx,ceilingy)
34 Z = np.ones((np.shape(ceilingy)[0],np.shape(ceilingx)[0]))
35 Z.fill(alt_rest1)

```

```

35 thermalsx = np.arange(thermals1[0],thermals1[1]+1,1)
36 thermalsy = np.arange(thermals1[2],thermals1[3]+1,1)
37 Xt, Yt = np.meshgrid(thermalsx,thermalsy)
38 Zt = np.zeros((np.shape(thermalsy)[0],np.shape(thermalsx)[0]))
39
40 thermalsx2 = np.arange(thermals2[0],thermals2[1]+1,1)
41 thermalsy2 = np.arange(thermals2[2],thermals2[3]+1,1)
42 Xt2, Yt2 = np.meshgrid(thermalsx2,thermalsy2)
43 Zt2 = np.zeros((np.shape(thermalsy2)[0],np.shape(thermalsx2)[0]))
44
45 # add modifier locations to plot
46 ax.plot_surface(Xt,Yt,Zt,color="Green",alpha=0.4)
47 ax.plot_surface(Xt2,Yt2,Zt2,color="Green",alpha=0.4)
48 ax.plot_surface(X,Y,Z,color="Red",alpha=0.4)
49 ax.plot_surface(w1x,w1y,Zwallsx,color="Red",alpha=0.4)
50 ax.plot_surface(w2x,w2y,Zwallsy,color="Red",alpha=0.4)
51
52 ax.legend()
53
54 # show compount figure
55 plt.show()
56
57 Note: github code includes code to plot best-action vector field
58

```

It was tested whether the agent adjusts the best learned path when the restricted airspace and the thermals are moved with the start and end points being kept constant. As can be seen from Figure 1, the agent indeed attempts to make efficient use of the available thermals without violating restricted airspace.

During the development it was also observed, that the frequency of random choices has to be kept at a low value to prevent the agent from taking primarily random actions. This issue has caused some problems in development, as the solution did not converge to a path at all with a too high set value. That being said, the epsilon part of epsilon greedy allows for more efficient exploration of the available space and is necessary to populate the Q-table. One fact that was quickly confirmed is that the agent requires significantly more episodes to reliably find a best route from a selected starting point, owed to the fact that the space to be explored has grown significantly compared to initial 2-dimensional test runs. This coupled with the random initialisation and occasional termination of the agent can occasionally lead to problems with finding an optimal route via *extrsp*, as the agent can not decide on a sole path decision if the decision becomes ambiguous at any point. A random step in any direction, followed by a return to the original functionality could solve this problem.

It is also noticeable that the limited actions of the agent restrict its motion significantly. For further exploration of the topic, one should either think about adding more available actions, such as Even.Up-Right, Up.Up-Right, Down.Up-Right or consider changing over to a continuous system, which would also allow for a more detailed description of the aircraft dynamics.

6 Conclusions

Q-Learning, albeit mostly used for two dimensional path optimisation problems has proven to be a feasible, if not optimal method to have an agent reach a goal within a 3D environment. The agent is capable of taking into account different types of environment modifiers, in this case: Restricted Airspace, Boundaries of the space to be explored and Thermals. These are defined by modifying the reward table to represent them via appropriate rewards. The results are shown in Figure 1, where the agent trace can be seen as a blue line, with the Thermals being shown as green ground-trace and the restricted airspace as simplified red box. During development of the code it was often noticeable that the agent would move in ways, that were not immediately obvious as an ideal path. However, this is to be expected when simplifying the dynamics of the aircraft to the extend that was decided upon. As such, the aircraft can be observed to climb without moving in either of the cardinal directions. A better idea might be to allow only slanted climbs, but it would likely be a better idea to develop a new continuous state system for the aircraft.

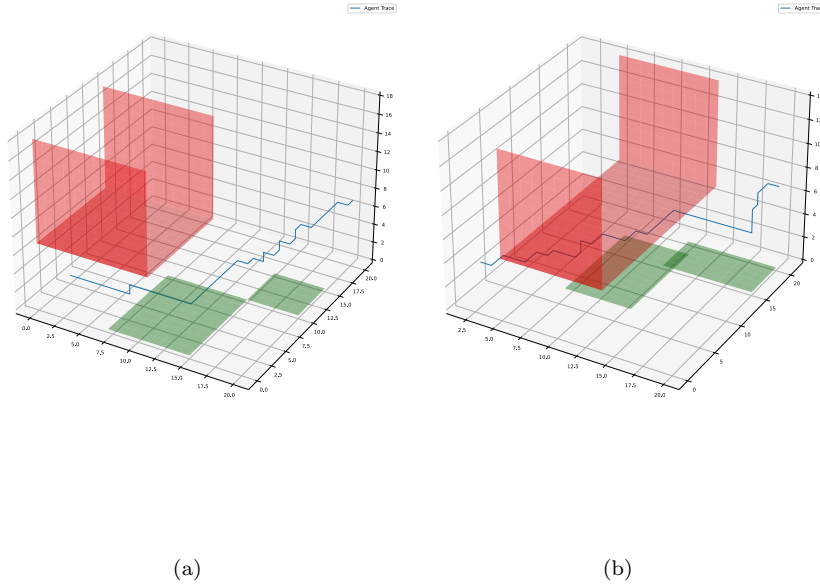


Figure 1: (a) Initial Run (b) Modified (Green shows Thermals, Red Restricted Airspace)

Note that the figures are available in higher resolution and better scale in the github repository <https://github.com/NicoEVA/BioInspired!>

References

- [1] ANDREW, A.: Reinforcement Learning: An Introduction by Richard S. Sutton and Andrew G. Barto, Adaptive Computation and Machine Learning series, MIT Press (Bradford Book), Cambridge, Mass., 1998, xviii + 322 pp, ISBN 0-262-19398-1. In: *Robotica* 17 (1999), S. 229–235
- [2] KAMPEN, Erik-Jan van: *Reinforcement Learning for Flight Control*. – URL <https://bit.ly/3mIrDQe>
- [3] SOPER, Dr. D.: *Foundations of Q-Learning*. – URL <https://bit.ly/3mToL3b>