
Control of 2-degree of freedom robot using Advantage-Actor-Critic method

1
2
3
4
5
6
7
8
9

Nazerke Sandibay
Department of Robotics and Mechatronics
Nazarbayev University
nazerke.sandibay@nu.edu.kz

10

Abstract

11 The project work aims to find the optimal trajectory of 2-degree of freedom
12 robot in a space with obstacles using Advantage-Actor-Critic Algorithm.
13 The learning environment of the robot was constructed and the performance
14 of a reinforcement learning algorithm concluded to be safe and optima for a
15 robot but too cautious.

16

1 Problem formulation

18 The purpose of the project is to find the best trajectory between two positions of 2 degrees
19 of freedom robot similar to a SCARA robot in the environment with obstacles. It can be done
20 using motion planning algorithms such as A* and Reinforcement learning techniques.

21 The latter may be faster and universal compared to the former. Therefore, Advantage-Actor-
22 Critic Reinforcement learning algorithm was written and tested in a simple environment with
23 obstacles with a constant position. For the sake of convenience Open-Ai gym environment of
24 the robot was created to work with RL baselines. Later it can be applied to a changing
25 environment.



26

27

28

Figure 1. SCARA robot.

29 **1.1 Task description**

30 Our goal is to pick up an object in space at a specified position and place it to another
31 position using an end-effector electromagnet that can be activated or deactivated at will. The
32 robot's vertical movement can occur at any time without affecting obstacle collisions: as a
33 result, the vertical movement is planned independently and we only need to determine the
34 motion of the manipulator in the horizontal plane.[1] The robot can thus be modeled as a
35 simple 2R planar manipulator. The robot base (i.e., where link 1 is fixed to the ground) is at
36 $(x,y) = (0,0)$. The links have lengths $l_1 = 0.5$ m and $l_2 = 0.4$ m, respectively. The first
37 obstacle to be avoided is a wall, which runs parallel to the x-axis, keeping a distance of 0 m
38 from it. Also, there are two other obstacles: these have a fixed position and have been
39 conservatively represented by two circles, both with radius $B = 0.2$ m, and with center at
40 $(xc1,yc1) = (-0.6,0.7)$ and $(xc2,yc2) = (0.6,0.7)$ m, respectively.

41 The thickness of the links can be neglected, as the sizes of all obstacles have been already
42 augmented to account for the robot link thickness as well. The angular motion of link 1 is
43 only limited by the presence of the wall (so no additional constraints have to be inserted),
44 while link 2 can only move within a range of $\pm 90^\circ$ with respect to the configuration in which
45 it is perfectly aligned with link 1 (i.e., $\theta_2 \in [-\pi/2, \pi/2]$). Our task is to plan a motion from
46 any given initial configuration (where the object is picked) to any final configuration (where
47 the object is placed), chosen in the free space, avoiding any collision during the robot
48 motion.

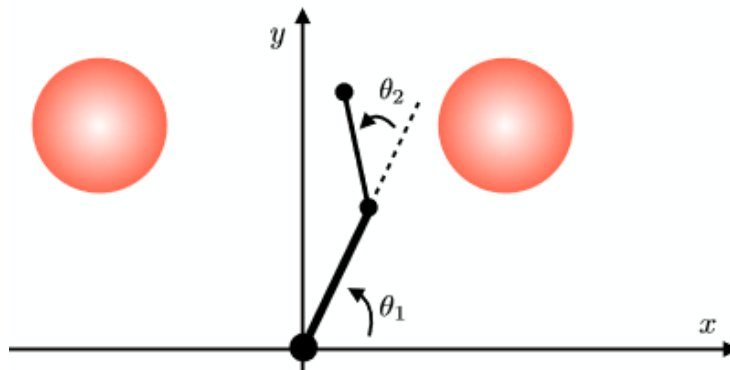
49 **2 Background**

50 **2.1 Heuristic function**

51
52
53
54 Reward functions play a crucial role in reinforcement learning. In my case reward was
55 chosen to be proportional to negative of heuristic. Heuristic function approximates the
56 distance between two objects. It was taken as 20 plus negative heuristic. It means that the
57 reward of the states that are close to the goal is higher. For example, the heuristic of the far
58 element is 15, while the heuristic of the closer element is 10. Consequently, their reward will
59 be 5 (or $20-15=5$) and 10 (or $20-10=10$). A reward of the closest element is higher, therefore
60 the algorithm will try to move closer to the goal to maximize reward.

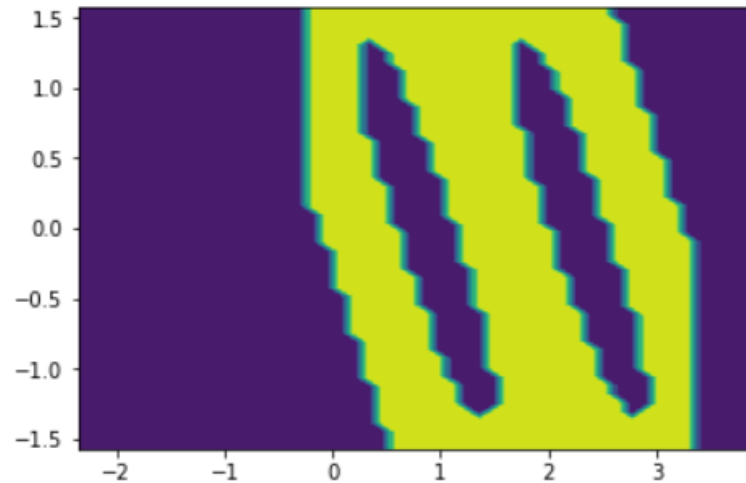
61
62
63 **3 Environment**

64 Two degrees of freedom of the robot corresponds to two agents. So, the environment is
65 multiagent with agents that depend on each other. Therefore, outputs of the algorithm should
66 be two angles (theta 1 and theta 2) corresponding to the angles of the arms with respect to
67 the neutral axis.



68
69
70
Figure 2. 2DOF robot (view from top)

71 3.1 Free space and obstacle space



73 Figure 3. Free&Obstacle space of the environment

74 First of all, we have to represent the robot configuration space, the configuration being $q =$
75 (θ_1, θ_2) . A grid of points has to be defined on both angles in a range of 2π . It is better to
76 choose the intervals for the two angles such that the free space is connected: for example,
77 rather than representing the range of both angles from 0 to 2π , one could do it between $-3\pi/4$
78 and $5\pi/4$. From a visual inspection, we notice that link 1 can collide with the wall, but not
79 with any of the circular obstacles: as a consequence, there is no need to define spheres
80 around link 1.

82 3.2 Create a gym environment

84 Gym environment with properties and functions similar to the OpenAI gym environment was
85 created.

```
86  
87 class GridEnvironment(gym.Env):  
88     metadata = { 'render.modes': [] }  
89  
90     def __init__(self, D, x, y, agent, goal):  
91         self.x=x  
92         self.y=y  
93         self.bool=D  
94         self.low = np.array([-3 * math.pi / 4, -math.pi / 2])  
95         self.high = np.array([5 * math.pi / 4, math.pi / 2])  
96  
97         self.observation_space = spaces.Box(self.low, self.high,  
98 h, dtype=np.float32)  
99         self.action_space = spaces.Discrete(4)  
100         self.max_timesteps = 25001  
101         self.agent=self.cord(agent)  
102         self.goal=self.cord(goal)  
103  
104     def reset(self):  
105         self.timestep = 0  
106         self.agent_pos = self.agent  
107         self.goal_pos = self.goal  
108         self.state = np.zeros((50,50))
```

```

109     self.state[tuple(self.agent_pos)] = 1
110     self.state[tuple(self.goal_pos)] = 0.5
111
112     return self.agent_pos
113
114     def cord(self, pos):
115         x=getcoordinates(self.x, self.y, pos)
116         return np.array(x)
117     def obs(self):
118         observation = self.state.flatten()
119         return observation
120
121     def step(self, action):
122         # 0 - down
123         # 1 - up
124         # 2 - right
125         # 3 - left
126
127         s=False
128
129         if action == 0:
130             if D[self.agent_pos[1],self.agent_pos[0]+1] and se
131 lf.agent_pos[0]<48:
132
133                 new=[self.agent_pos[0].copy() + 1,sel
134 f.agent_pos[1]]
135
136                 s=True
137
138             if action == 1:
139                 if D[self.agent_pos[1],self.agent_pos[0]-1] and se
140 lf.agent_pos[0]>1:
141
142                 new=[self.agent_pos[0].copy() - 1,sel
143 f.agent_pos[1]]
144
145                 s=True
146
147             if action == 2:
148                 if D[self.agent_pos[1]+1,self.agent_pos[0]] and se
149 lf.agent_pos[1]<48:
150
151                 new= [self.agent_pos[0],self.agent_po
152 s[1].copy() + 1]
153
154                 s=True
155
156             if action == 3:
157                 if D[self.agent_pos[1]-1,self.agent_pos[0]] and se
158 lf.agent_pos[1]>1:
159
160                 new=[self.agent_pos[0],self.agent_pos
161 [1].copy() - 1]
162
163                 s=True
164
165         if s:

```

```

163         new=np.array(new)
164     else:
165         new=self.agent_pos
166         r = heuristic(new[0],new[1],self.goal_pos[0],self.goa
167 l_pos[1])
168         if r<1:
169             if r==0:
170                 reward=np.array([1000.0])
171             else:
172                 reward=np.array([500.0+1/r])
173         else:
174             reward=np.array([20.0-r])
175
176         self.agent_pos=new.copy()
177         done = True if self.timestep >= self.max_timesteps els
178 e False
179
180
181         self.timestep += 1
182
183         info = {}
184         if not s:
185             reward= np.array([0.0])
186             self.state[tuple(new)] = 1
187
188         return new, reward, done, info
189
190     def render(self):
191         plt.imshow(self.state)
192     def pos(self,cor):
193         self.agent_pos=cor
194     def agent(self):
195         return self.agent_pos
196

```

197 3.2.1 Calculation of reward

198 Choice of reward function is very important since the performance of the algorithm will
199 depend on it.

200 The reward function is equal to twenty subtracted to the approximated distance. As the agent
201 gets closer to the goal it starts to increase.

202 In this case, when the agent reaches the goal(distance is 0) it gets 1000 reward. If the
203 distance is smaller than 1, it can get a reward between 500 and 500+1/(closest distance).
204 Performance still can be improved by changing the reward.

205

206 3.2.2 Step function

207 Our environment contains obstacles. Therefore, an agent has to make sure that the next state
208 is “safe”. It is done via the configuration space matrix described above. It makes a decision
209 based on the value of free space matrix on a given position.

210

211 4 Algorithms

212

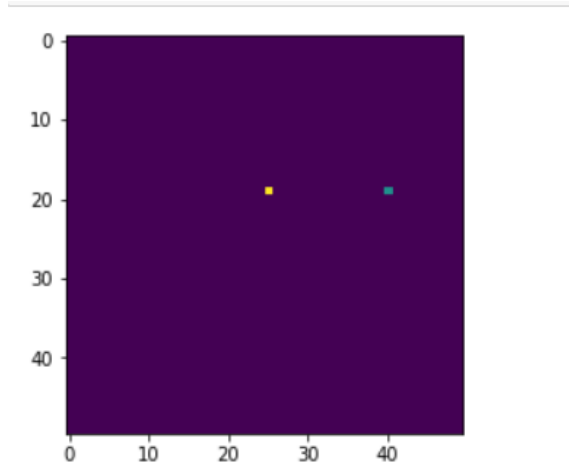
213 4.1 Actor-Critic method

214 A synchronous, deterministic variant of Asynchronous Advantage Actor-Critic (A3C)
215 algorithm from the library of Stable Baseline was used.

216

217 4.2 Results

218 Trajectory of the agent moving from start point = $([0, 0])$ to goal point = $([0, 1 \text{ rad}])$. There
219 is no obstacle between two points. Therefore, the trajectory looks like a straight line.
220

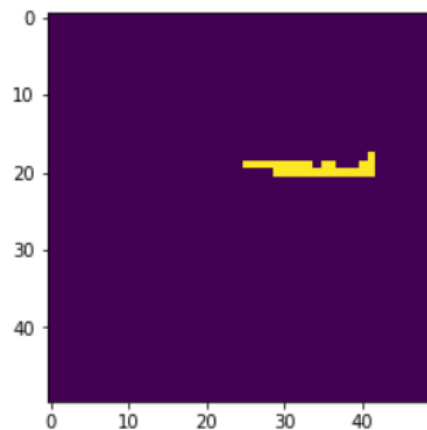


221

222 Figure 4. Grid before learning. Agent position and goal position.

223

224



225

226

227 Figure 5. Grid after learning. Trajectory

228 Trajectory of the agent moving from start point = $([2, 0.5])$ to goal point = $([3, -0.5])$ was
229 estimated. There is an obstacle between two points. Therefore, the trajectory is more
230 complex. Learning was done for the environment with one obstacle with the following
configuration space

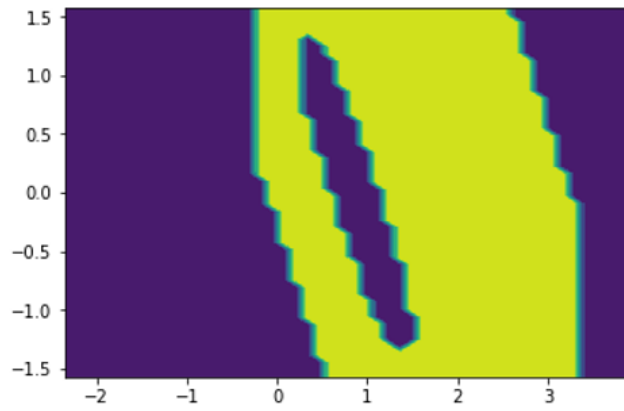


Figure 6. New free&obstacle space

231
232
233
234

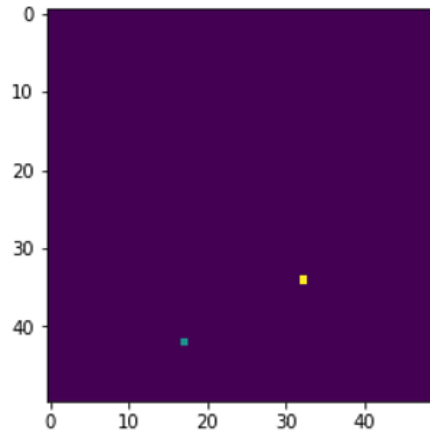


Figure 7. Grid before learning.

235
236
237
238
239
240

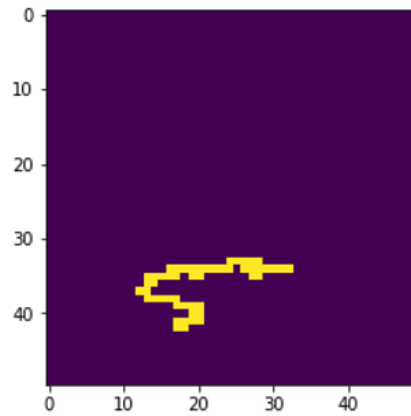


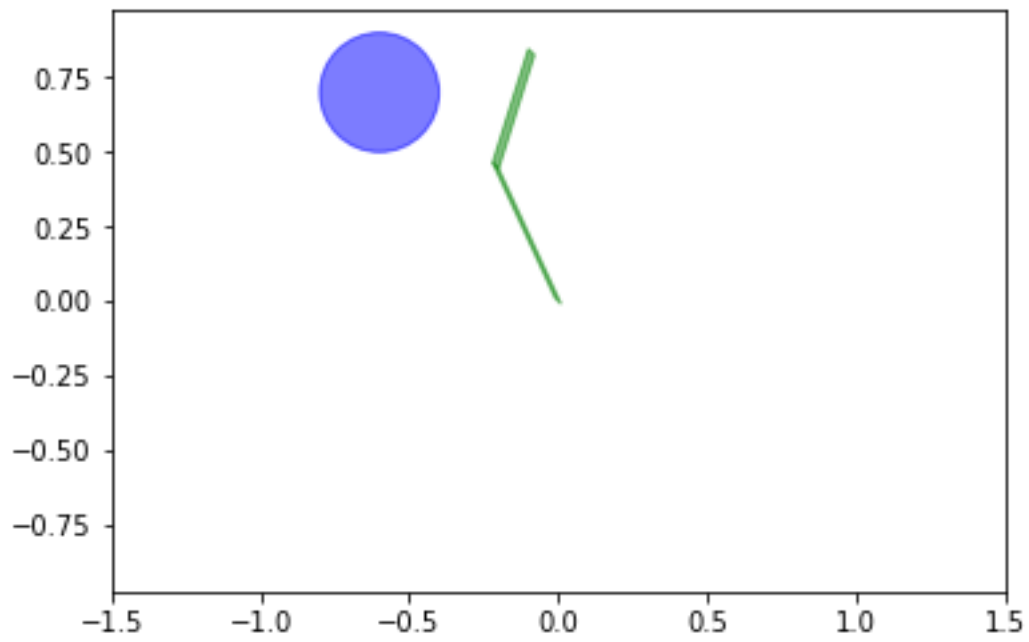
Figure 8. Grid after learning. Trajectory

241

242

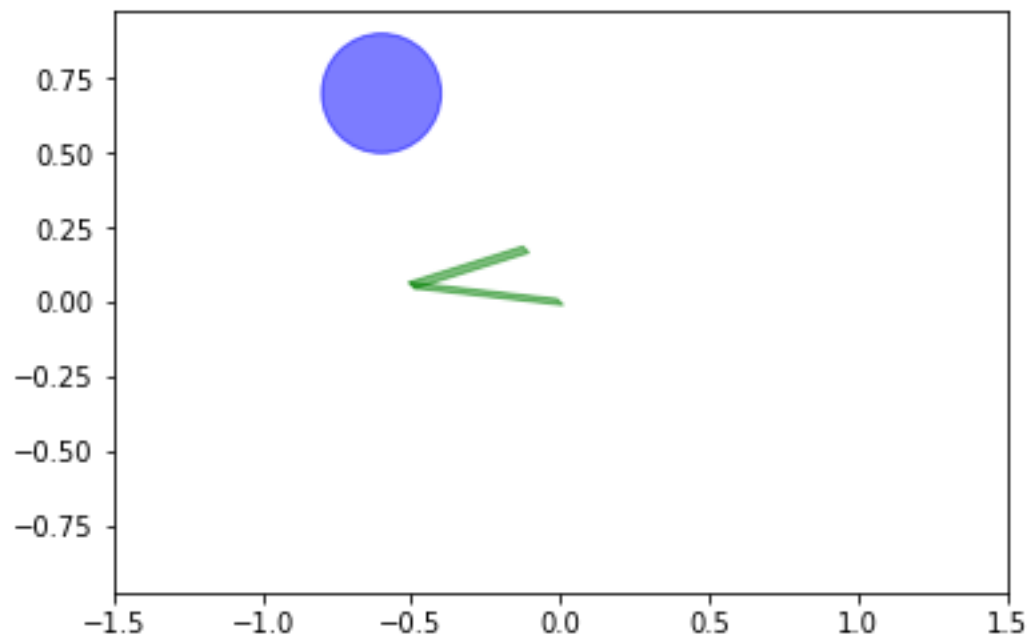
243

244 **4.3 Visualization**



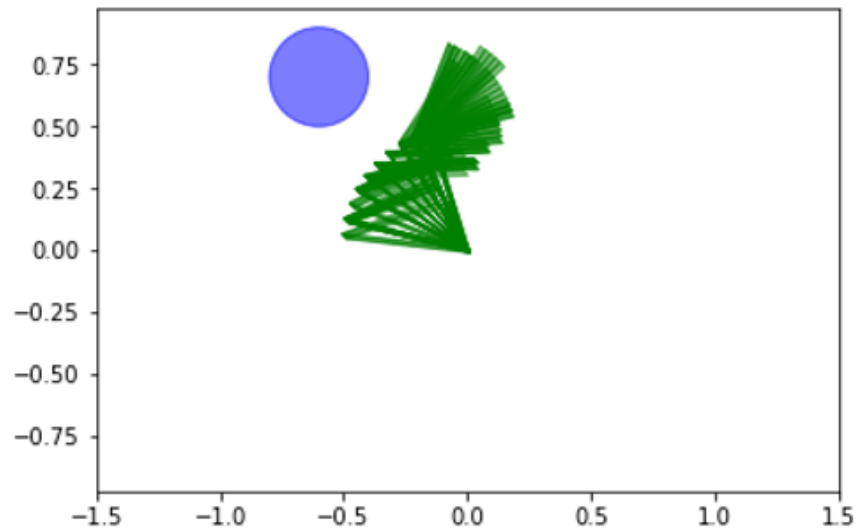
245
246

Figure 9. The initial position of a robot



247
248

Figure 10. Final Position of a robot



249

250

Figure 11. Whole trajectory

251

252 **5 Improvement**

253 The algorithm does not always generate an optimal trajectory. It might be caused by a
254 complex environment or reward values. However, the agent often tries to move toward the
255 goal. Performance of the algorithm can be improved by choosing another reward
256 function(such as another type of heuristic function)

257 **References**

258 [1] <https://www.fanuc.eu/de/en/robots/robot-filter-page/scara-series/selection-support>

259 [2] Volodymyr Mnih et.al (2016) *Asynchronous Methods for Deep Reinforcement Learning*
260 <https://stable-baselines.readthedocs.io/en/master/modules/a2c.html>

261