

Q-LEARNING

Karol Działowski

nr albumu: 39259

przedmiot: Uczenie ze wzmocnieniem

Szczecin, 18 stycznia 2021

Spis treści

1 Cel laboratorium	1
2 Zadanie 1 - zapałki	1
3 Zadanie 3 - dwuwymiarowa siatka	4
4 Wnioski	8

1 Cel laboratorium

Celem laboratorium nr 3 była implementacja algorytmu q-learning w celu wyznaczenia sub-optymalnej strategii przy eksploracji środowiska.

2 Zadanie 1 - zapałki

Przedmiotem rozważań jest gra w zapałki. Gracze mają do dyspozycji pewną liczbę (np. $N = 10$) zapałek. W grze bierze udział dwóch graczy, którzy na przemian zabierają 1, 2 lub 3 zapałki. Przegrywa ten gracz, który zabiera zapałki jako ostatni.

Należało wyznaczyć strategię z wykorzystaniem algorytmu Q-Learn.

Kod źródłowy 1: Zadanie 1 - zapałki

Źródło: Opracowanie własne

```
| 1 import random
```

```

2  import numpy as np
3
4  def take_action(state, action):
5      # if losing state
6      if state == 1:
7          return 0, 0, True
8
9      state_next = state - (action+1)
10     # if winning move return R and state
11     if state_next == 1:
12         return 0, 1, True
13
14     # random opponents action
15     action_opponent = random_action()
16
17     # calculate reward
18     R = 0
19     done = False
20     state_next = state_next - (action_opponent+1)
21
22     if state_next <= 0:
23         X_n = 0
24         R = 1
25         done = True
26
27     return state_next, R, done
28
29 def random_action():
30     return random.choice([0, 1])
31
32 def main():
33     N = 10 # starting state
34     q = np.random.rand(N+1, 2)
35
36     gamma = 0.01
37     epochs = 5000
38     epsilon = 0.2 # exploration
39
40     for i in range(epochs):
41         state = N
42         done = False
43
44         while not done:
45             print(f"Epoka {i}/{epochs}")
46
47             # chose action
48             if random.random() < epsilon:
49                 action = random_action()
50             else:
51                 action = q[state].argmax()
52

```

```

53         # take action
54         state_next, reward, done = take_action(state, action)
55
56         # update delta
57         q[state][action] = reward + gamma * q[state_next].max()
58         state = state_next
59
60     print()
61     print("Macierz Q")
62     print(q)
63
64     print("Wyznaczona strategia")
65     policy = np.argmax(q, axis=1)+1
66     print(policy)
67
68     print("Sprawdzenie strategii")
69     print(policy[9] == 2 and
70           policy[8] == 1 and
71           policy[6] == 2 and
72           policy[5] == 1 and
73           policy[4] == 1 and
74           policy[3] == 2 and
75           policy[2] == 1)
76
77
78 if __name__ == "__main__":
79     main()

```

Wyznaczona strategia pokrywa się z strategią optymalną wyznaczoną na poprzednich laboratoriach.

Kod źródłowy 2: Zadanie 1 - wyznaczona strategia

Źródło: Opracowanie własne

```

1  Macierz Q
2  [[9.89671802e-01  6.06568403e-01]
3   [9.89671802e-03  9.89671802e-03]
4   [1.00989672e+00  1.00861157e+00]
5   [9.89671802e-05  1.00989672e+00]
6   [1.00989672e-02  9.89671802e-05]
7   [1.00989672e-02  1.00989672e-02]
8   [9.89671802e-07  1.00989672e-02]
9   [1.00989672e-04  9.89671802e-07]
10  [1.00989672e-04  1.00989672e-04]
11  [2.06204220e-01  8.61157285e-01]
12  [1.00989672e-06  9.89671802e-09]]
13 Wyznaczona strategia
14 [1 1 1 2 1 1 2 1 1 2 1]
15 Sprawdzenie strategii
16 True

```

3 Zadanie 3 - dwuwymiarowa siatka

Zaimplementowano symulator środowiska o siatce 10×10 z przeszkodami. Ustawiono stany absorbujące. Wykorzystano algorytm q-learn do nauczania strategii.

Kod źródłowy 3: Zadanie 2 - grid

Źródło: Opracowanie własne

```
1  import random
2  import numpy as np
3
4
5  class Grid:
6      def __init__(self):
7          self.height = 10
8          self.width = 10
9          self.actions = [0, 1, 2, 3]
10         self.state_count = self.height * self.width
11         self.action_count = len(self.actions)
12
13         self.policy_print_characters = {0: "<", 1: ">", 2: "^", 3: "v"}
14
15         self.grid = np.zeros((self.height, self.width))
16         # end points
17         self.grid[0, 0] = 10
18         self.grid[-1, 0] = 50
19
20         # obstacles
21         self.grid[3, 2:8] = -1
22
23         self.grid[2:8, 6] = -1
24         self.grid[2:8, 5] = -1
25         self.grid[2:8, 4] = -1
26         self.grid[2:8, 3] = -1
27
28
29     def random_position(self):
30         """Find random starting position"""
31         while True:
32             h = random.randrange(0, self.height)
33             w = random.randrange(0, self.width)
34             if self.grid[h, w] == 0:
35                 return (h, w)
36
37     def is_terminal(self, state):
38         """Terminal states are those with reward greater than 0"""
39         if self.grid[state[0], state[1]] > 0:
40             return True
41         else:
42             return False
43
```

```

44 def to_string(self, state=None, policy=None):
45     res = " " + " " * self.width + " " + "\n"
46     for i in range(self.height):
47         res += " "
48         for j in range(self.width):
49             if self.grid[i, j] == -1:
50                 res += "#"
51             elif self.grid[i, j] > 0:
52                 res += "X"
53             elif state is not None and state[0] == i and state[1] == j:
54                 res += "a"
55             elif policy is not None:
56                 res += self.policy_print_characters[policy[i, j]]
57
58         else:
59             res += " "
60     res += " "
61     res += "\n"
62     res += " " + " " * self.width + " "
63     return res
64
65 def __str__(self):
66     return self.to_string()
67
68 def available_actions(self, state):
69     actions = []
70     if self.is_terminal(state):
71         return actions
72
73     h, w = state
74     if h > 0:
75         if self.grid[h-1, w] >= 0:
76             actions.append(2) # up
77     if w > 0:
78         if self.grid[h, w-1] >= 0:
79             actions.append(0) # left
80     if h < self.height-1:
81         if self.grid[h+1, w] >= 0:
82             actions.append(3) # down
83     if w < self.width-1:
84         if self.grid[h, w+1] >= 0:
85             actions.append(1) # right
86
87     return actions
88
89
90 def take_action(self, state, action):
91     h, w = state
92     state_next = [h, w]
93
94     if action == 0: # left

```

```

95         state_next[1] = w-1
96     if action == 1: # right
97         state_next[1] = w+1
98     if action == 2: # up
99         state_next[0] = h-1
100    if action == 3: # down
101        state_next[0] = h+1
102
103    done = self.is_terminal(state_next) or self.available_actions(state_next) ==
104    []
105    reward = self.grid[state_next[0], state_next[1]]
106    return state_next, reward, done
107
108
109 def main():
110     import numpy as np
111     import time
112     import os
113
114     env = Grid()
115
116     qtable = np.random.rand(env.height, env.width, env.action_count)
117
118     epochs = 5000
119     gamma = 0.05
120     epsilon = 0.1
121
122     for i in range(epochs):
123         state = env.random_position()
124         reward = 0
125         done = False
126         steps = 0
127
128         while not done:
129             h, w = state
130
131             # count steps to finish game
132             steps += 1
133
134             # exploration
135             if np.random.random() < epsilon:
136                 action = np.random.choice(env.available_actions(state))
137             # best action
138             else:
139                 action = qtable[h, w].argmax()
140                 available_actions = env.available_actions(state)
141                 if action not in available_actions:
142                     action = np.random.choice(env.available_actions(state))
143
144

```

```

145
146         # take action
147         state_next, reward, done = env.take_action(state, action)
148
149         # update qtable value with Bellman equation
150         qtable[h, w, action] = reward + gamma * np.max(qtable[state_next])
151
152         # update state
153         # print("Action", action)
154         # print("Step", steps)
155         state = state_next
156
157         #print(env.to_string(state=state))
158
159         #print("Epoka: ", i + 1, "kroki:", steps)
160         if i % 100 == 0:
161             print("Wyznaczona strategia")
162             policy = np.argmax(qtable, axis=2)
163             print(env.to_string(policy=policy))
164
165         print("Wyznaczona strategia")
166         policy = np.argmax(qtable, axis=2)
167         print(policy)
168         print(env.to_string(policy=policy))
169
170 if __name__ == "__main__":
171     main()

```

Wyznaczona strategia w 5000 epokach przy $\gamma = 0.9$, $\epsilon = 0.2$:

Kod źródłowy 4: Zadanie 2 - grid

Źródło: Opracowanie własne

```

1  Wyznaczona strategia
2  [[0 0 0 0 3 0 3 1 1 3]
3   [2 0 0 2 0 0 0 1 1 2]
4   [1 0 0 0 2 2 1 1 1 2]
5   [2 2 3 0 0 0 0 1 1 2]
6   [1 2 0 0 0 0 1 1 1 2]
7   [1 2 0 0 0 0 1 1 2 2]
8   [1 0 0 0 0 0 1 1 1 3]
9   [3 3 3 0 3 3 1 3 3 3]
10  [3 3 1 3 3 0 3 0 0 2]
11  [0 0 0 0 0 0 0 1 2 2]]
12
13  X<<<<v<v>>> v
14  ^<<<^<<<<>>>^
15  ><<# ###>>^
16  ^^# #####>^
17  >^<# ###>>^
18  >^<# ###>^^
19  ><<####>> v

```

```
20 vvv#### vvv
21 vv>vv<v <<^
22 X <<<<<<>^^
```

4 Wnioski

Zaimplementowano algorytm Q-Learning na dwóch środowiskach - zapalki oraz siatka z przeszkodami.

Wyznaczona strategia dla zapalek pokrywa się z wyznaczoną strategią optymalną z poprzedniego zadania.

W przypadku problemu dla siatki algorytm nie wyznaczył optymalnej strategii przy danych parametrach. W okolicy stanów końcowych strategia jest dobra, ale im dalej od punktu tym więcej błędnych akcji jest ocenianych jako najlepsze. Prawdopodobnie wynika to ze złego doboru parametrów uczących.