

Q-learning ile Yol Planlaması

Abdullah Yaşar KISA
Mühendislik Fakültesi
Bilgisayar Mühendisliği Bölümü
Kocaeli Üniversitesi
abdullahyasarkisa@hotmail.com
180201015

İbrahim Kafkaslı
Mühendislik Fakültesi
Bilgisayar Mühendisliği Bölümü
Kocaeli Üniversitesi
ibrahim.kafkasli@hotmail.com
180201078

Özet – Bu projede temel amaç robotun Q learning algoritması kullanarak engel sütunlarından kaçması ve beyaz alanlardan geçerek doğru yol almasıdır. Proje de yol temsili olarak bir matris kullanılmıştır ve robot yolunu iki boyutlu bir ortamda bulmaktadır. Proje python programlama dili kullanılarak geliştirilmiştir. Çıktılar ise kullanıcıya bir ui ile sunulmuştur.

Anahtar Kelimeler – Q-learning, Pekiştirmeli Öğrenme, Brute-Force, Python, Pygame

Abstract – The main purpose of this project is for the robot to escape from the obstacle columns by using the Q learning algorithm and take the right path by passing through the white areas. In the project, a matrix is used as a path representation and the robot finds its path in a two-dimensional environment. The project has been developed using the python programming language. Outputs are presented to the user with a ui.

Keywords – Q-learning, reinforcement learning, Brute-Force, Python, Pygame

I. GİRİŞ

Q-Learning algoritması, pekiştirmeli öğrenmenin(reinforcement learning) en çok bilinen algoritmalarındandır. Algoritmadaki temel amaç bir sonraki hareketleri inceleyip yapacağı hareketlere göre kazanacağı ödülü görmek ve bu ödülü en çoklayıp (maximize) buna göre hareket etmektir.

Ajanın ödül haritası, gitmesini istediğimiz ya da istemediğimiz yerler daha önceden bizim tarafımızdan belirlenir ve bu değerler ödül tablosuna (Reward Table) yazılır. Ajanın tecrübeleri bu ödül tablosuna göre şekillenecektir.

Ajan ödüle giderken her iterasyonda edindiği tecrübeleri gidebildiği yerleri seçerken en çoklamak için kullanır. Bu tecrübeleri ise Q-Tablosu (Q-Table) adı verdiğimiz bir tabloda tutar.

Q-Tablomuz başlangıçta ajanımızın hiçbir tecrübesi olmadığı için sıfırlarla doludur ve bu yüzden ajanımız ilk seçimlerinde Q-Tablosundaki sıfırları en çoklayacağından rastgele hareket edecektir. Bu rastgelelik ajanın ilk ödülü bulmasına kadar sürecektir. Ajan ödülün olduğu bir yere geldiği anda ödüle gelmeden önceki yerini bilir ve bu yerin değerini kendi tecrübelerini biriktirdiği Q-Tablosuna yazar. Bu yazma işleminin belirli bir algoritması bulunmaktadır.

$$Q(s,a) = Q(s,a) + lr*(r(s,a) + Y*max(Q(s',a')) - Q(s,a))$$

“Q(s,a)” (state,action) dediğimiz değer bizim şu anda {bulunduğumuz, gideceğimiz} dizin, “lr” dediğimiz değer öğrenme katsayısı (learning rate),“r(s,a)” dediğimiz değer bizim {bulunduğumuz, gideceğimiz} ödül tablomuzdaki ödül değerimiz, “Y” değeri gamma, “max (Q (s',a'))” değeri ise gidebileceğimiz {gideceğimiz,gideceğimiz yerden gidebileceğimiz} yerlerin en yüksek Q değeridir.

Her seferinde ajan bu algoritmaya göre bir ilerisini tahmin edip hareket eder ve ödüle ulaşır. Ödüle ulaştıktan sonra ise ajan tekrardan rastgele bir yerden harekete başlar ve tekrardan ödülü bulmaya çalışır. Bu işlem devam ettikçe ajan hangi durumda nereye gideceğini bilmeye başlar.

Bu projede, ajanımız kullanıcıdan alınan başlangıç ve bitiş değerlerine göre, yine kullanıcıdan alınan matris boyutuna göre başlangıçtan, hedefe olan yolunu Brute-Force yaklaşımı ile bulacaktır.

II. KARŞILAŞILAN ZORLUKLAR

İlk olarak karşılaştığımız zorluk pandemi dolayısıyla ekip arkadaşımızdan uzak kalmamız oldu. Bu zorluğu aşmak her ne kadar zor olsa da, GIT ve GITHUB teknolojileri ile kod paylaşımını basit bir şekilde çözmüş olduk.

İkinci karşılaştığımız problem ise, python syntax'ına ekip arkadaşımızla birlikte çok fazla hakim olmayışımızdı. Daha önce python da çok az bir tecrübemiz olduğundan dolayı, dilin kendi syntax'ına alışmamız zaman aldı.

Üçüncü olarak karşılaştığımız problem, Q-learning algoritmasıyla alakalı daha önce hiçbir bilginin olmayışındı. Brute-Force yaklaşımı ile birçok kod yazmamıza rağmen algoritmaya alışmamız biraz zamanımızı aldı. Ancak bu problemi zamanla ve doğru araştırmayla aşmayı başardık.

Dördüncü karşılaştığımız problem ise, kullanıcı arayüzü kısmıydı. Daha önce birçok kullanıcı arayüzü olan uygulamalar geliştirmemize rağmen, python ile böyle bir tecrübemizin olmayışı bizi çok zorladı. Bu problemi doğru ui kütüphanelerini kullanarak aşmayı başardık.

Beşinci problemimiz ise projemizin yoğun olduğumuz bir haftaya denk gelmesiydi. Çok yoğun olduğumuzdan dolayı projeye yeterli zaman ayıramamıştık. Ancak hocalarımızın projeyi 2 gün ertelemesi sayesinde bu problemin de üstesinden gelmiş olduk.

III. YÖNTEM

```
1 from tkinter.constants import X
2 import PySimpleGUI as sg
3 import numpy as np
4 import pygame
5 from time import time, sleep
6 from random import randint as r
7 import random
8 import matplotlib.pyplot as plt
9 import math
10
11 class Form:
12     def FormPage(self):
13         sg.theme('LightGrey 6')
14         layout = [
15             [sg.Text('Please enter information')],
16             [sg.Text('Start Location', size=(15, 1)), sg.InputText()],
17             [sg.Text('Finish Location', size=(15, 1)), sg.InputText()],
18             [sg.Text('Matrix Size', size=(15, 1)), sg.InputText()],
19         ]
20         [sg.Submit(), sg.Cancel()]
21     ]
22
23     window = sg.Window('Qlearning App', layout)
24     event, value = window.read()
25     window.close()
26     return value
27
28 form = Form()
29 value = form.FormPage()
30
31 start = value[0]
32 start = start.split(",")
33 finish = value[1]
34 finish = finish.split(",")
35 size = value[2]
36 size = size.split(",")
```

Kodumuzda ilk olarak gerekli kütüphaneler dahil edilmiş ve kullanıcıdan, matris boyutu, başlangıç konumu ve bitiş konumunun alınacağı ui tasarlanıp çalıştırılmıştır. Burada alınan konumlar X,Y formatında olduğundan split edilerek belirli değişkenlerde sonradan kullanılmak üzere tutulmuştur. Kullanıcıdan bilgilerin alınacağı ui ise PySimpleGUI kütüphanesinden yararlanılarak yapılmıştır.

```
43 # Qlearning #
44 n = int(size[0]) # represents no. of side squares(n*n total squares)
45 scrx = n*30
46 scry = n*30
47 x = math.ceil(scrx/int(n))
48
49 background = (0, 0, 0) # used to clear screen while rendering
50 # creating a screen using Pygame
51
52 screen = pygame.display.set_mode((scrx, scry))
53 colors = [(255, 225, 255) for i in range(n**2)]
54
55 #create reward matrix
56 reward = np.zeros((n, n))
57 terminals = []
58 penalties = int(n*n*20/100)
59
60 startsPoint = []
61 startsPoint.append(int(start[0]))
62 startsPoint.append(int(start[1]))
63
64 finishPoint = []
65 finishPoint.append(int(finish[0]))
66 finishPoint.append(int(finish[1]))
67
```

Burada ise Q learning algoritması için başlangıç yapılıp gerekli değişkenlere atamalar yapılmıştır. Ayrıca Qlearning algoritmasının gösterileceği ekran için boyutlar ayarlanıp gerekli değişkenlere atanmıştır. Ayrıca tüm matris noktaları için beyaz renk ataması yapılmıştır. Bu renk ataması matriste bulunan labirent, başlangıç, bitiş ve yollara göre farklı renkler ile sonradan değiştirilecektir

```
68
69 while penalties != 0:
70     i = r(0, n-1)
71     j = r(0, n-1)
72     if reward[i, j] == 0 and [i, j] != [0, 0] and [i, j] != [n-1, n-1]:
73         reward[i, j] = -5
74         penalties -= 1
75         colors[n*i+j] = (255, 0, 0)
76         terminals.append(n*i+j)
77
78 reward[int(finish[0]),int(finish[1])] = 5 # finish position
79 colors[n*int(finish[0])+int(finish[1])] = (0, 255, 0)
80 colors[n*int(start[0])+int(start[1])] = (0, 137, 255)
81 terminals.append(n*(int(finish[0]) + int(finish[1]))
82 print(reward)
83
```

Burada ise, belirlenen engel oranına göre engeller rastgele olarak matriste konumlara atanmıştır. Engellerin bulunduğu konumlar kırmızı rengine boyanmıştır. Ayrıca engellerin bulunduğu konumlara -5 değeri atanırken, bitiş lokasyonunun bulunduğu konuma yine reward matrisinde 5 değeri atanmıştır. Ayrıca bitiş ve başlangıç konumları da ui için renklendirilmiştir.

```
86 Q = np.zeros((n**2, 4)) # Initializing Q-Table
87 actions = {"up": 0, "down": 1, "left": 2, "right": 3} # possible actions
88 states = {}
89 k = 0
90 for i in range(n):
91     for j in range(n):
92         states[(i, j)] = k
93         k += 1
94 alpha = 0.01
95 gamma = 0.9
96
97 #start position
98 current_pos = [int(start[0]), int(start[1])]
99
100 epsilon = 0.25
```

Burada ise Q matrisi başlangıç olarak 0 değerleriyle doldurulmuştur. Ayrıca ajan hareketi için durumlar tanımlanıp yine Qlearning algoritması için gerekli alpha, gamma, epsilon öğrenme katsayıları atanmıştır. Ayrıca başlangıç konumu olarak kullanıcının girdiği konum belirlenmiştir.

```
103 def select_action(current_state):
104     global current_pos, epsilon
105     possible_actions = []
106     if np.random.uniform() <= epsilon:
107         if current_pos[0] != 0:
108             possible_actions.append("up")
109         if current_pos[0] != n-1:
110             possible_actions.append("down")
111         if current_pos[1] != 0:
112             possible_actions.append("left")
113         if current_pos[1] != n-1:
114             possible_actions.append("right")
115         action = actions[possible_actions[r(0, len(possible_actions) - 1)]]
116     else:
117         m = np.min(Q[current_state])
118         if current_pos[0] != 0: # up
119             possible_actions.append(0(current_state, 0))
120         else:
121             possible_actions.append(m - x)
122         if current_pos[0] != n-1: # down
123             possible_actions.append(0(current_state, 1))
124         else:
125             possible_actions.append(m - x)
126         if current_pos[1] != 0: # left
127             possible_actions.append(0(current_state, 2))
128         else:
129             possible_actions.append(m - x)
130         if current_pos[1] != n-1: # right
131             possible_actions.append(0(current_state, 3))
132         else:
133             possible_actions.append(m - x)
134         action = random.choice([i for i, a in enumerate(possible_actions) if a == max(
135             possible_actions)]) # randomly selecting one of all possible actions with maximin value
136     return action
137
```

Bu kısımda ise ileride kullanılmak üzere random aksiyonlar belirlenmiştir. Ajan kazanç sağlayana kadar random hareket sağlayacaktır. Kazanç sağladığında ise kazanç sağladığı yol üzerinden hareketine devam edecektir. Burada bu tanımlama yapılmıştır.

```

139 def episode(iterasyon, startPoint, step, cost, episodeStep, episodeCost, road, probabilityRoad):
140     global current_pos, epsilon
141     current_state = states[(current_pos[0], current_pos[1])]
142     action = select_action(current_state)
143     if action == 0: # move up
144         current_pos[0] -= 1
145     elif action == 1: # move down
146         current_pos[0] += 1
147     elif action == 2: # move left
148         current_pos[1] -= 1
149     elif action == 3: # move right
150         current_pos[1] += 1
151     new_state = states[(current_pos[0], current_pos[1])]
152     if new_state not in terminals:
153         Q[current_state, action] += alpha*(reward(current_pos[0], current_pos[1]) + gamma*(
154             np.max(Q[new_state])) - Q(current_state, action))
155         step += 1
156         cost = cost+3
157         road.append([current_pos[0], current_pos[1]])
158     else:
159         Q[current_state, action] += alpha * \
160             (reward(current_pos[0], current_pos[1]) - Q(current_state, action))
161         step+=1
162         cost = cost+reward(current_pos[0], current_pos[1])
163         road.append([current_pos[0], current_pos[1]])
164         print(iterasyon, ":", iterasyon")
165         iterasyon += 1
166         # print(current_pos)
167         # step=round(startPoint[0]-current_pos[0])+round(startPoint[1]-current_pos[1])
168         # print("step: ", step)
169         # print("cost: ", cost)
170
171         episodeStep.append(step)
172         episodeCost.append(cost)
173         probabilityRoad.clear()
174         #if reward[road[len(road)-1][0],road[len(road)-1][1]]==5:
175         probabilityRoad=road[:]
176
177     #print(road)
178     road.clear()
179     cost = 0
180     step = 0
181     current_pos[0] = int(start[0])
182     current_pos[1] = int(start[1])
183
184     if epsilon > 0.05:
185         epsilon -= 3e-4 # reducing as time increases to satisfy Exploration & Exploitation Tradeoff
186
187     return iterasyon, step, cost, episodeCost, episodeStep, road, probabilityRoad
188
189

```

Bu kısımda ise Q learning hareketinde her bir adım için olan hareket tanımlanmıştır. Bu kısımda ajanın bir hareket gerçekleştirebilmesi için gerekli tüm işlemler yapılmıştır. Bu işlemler belirli değişkenlerde tutularak, kodun son çıktısında kullanılmak üzere saklanmıştır.

```

191 def layout(isFinish, probabilityRoad):
192     c = 0
193     if isFinish == False:
194         for i in range(0, scrx, 30):
195             for j in range(0, scry, 30):
196                 pygame.draw.rect(screen, (0, 0, 0),
197                     (j, i, j+30, i+30), 0)
198                 pygame.draw.rect(
199                     screen, colors[c], (j+1, i+1, j+28, i+28), 0)
200                 c += 1
201                 pygame.draw.circle(screen, (0, 0, 0), (
202                     current_pos[1]*30 + 14, current_pos[0]*30 + 14), 10, 0)
203     else:
204         for k in range(0, len(probabilityRoad)-1):
205             pygame.draw.circle(screen, (0, 0, 0), (
206                 probabilityRoad[k][1]*30 + 14, probabilityRoad[k][0]*30 + 14), 10, 0)
207             pygame.display.flip()
208

```

Burada ise ajanın yapacağı her hareket, Pygame kütüphanesi kullanılarak bir matris ekranında kullanıcıya sunulmuştur. Ayrıca tüm işlemler tamamlandıktan sonra, kullanıcıya ajanın bulduğu yolda bir ui üzerinde sunulmuştur.

```

209 def plot_results(steps, cost):
210     #
211     f, (ax1, ax2) = plt.subplots(nrows=1, ncols=2)
212     #
213     ax1.plot(np.arange(len(steps)), steps, 'g')
214     ax1.set_xlabel('Tur')
215     ax1.set_ylabel('Adım')
216     ax1.set_title('Her Tur için Adım')
217
218     #
219     ax2.plot(np.arange(len(cost)), cost, 'b')
220     ax2.set_xlabel('Tur')
221     ax2.set_ylabel('Kazanç')
222     ax2.set_title('Her Tur için Kazanç')
223
224     plt.tight_layout() # Function to make distance between figures
225
226     # Showing the plots
227     plt.show()
228
229

```

Burada ise tüm işlemlerin tamamlanmasının ardından, kullanıcıya her tur için adım ve her tur için kazanç grafikleri çıktı olarak verilmiştir.

```

230 def ciktiVer() :
231     dosya = open('engel.txt', 'w')
232     dosya.write("S= start point\n")
233     dosya.write("F= finish point\n")
234     dosya.write("R= road\n")
235     dosya.write("B= block\n")
236     dosya.write("\n\n\n")
237     for i in range(n):
238         for j in range(n):
239             if i==startPoint[0] and j==startPoint[1]:
240                 dosya.write("{}({},{},S)".format(i,j))
241             elif i==finishPoint[0] and j==finishPoint[1]:
242                 dosya.write("{}({},{},F)".format(i,j))
243             elif (reward[i][j]==0):
244                 dosya.write("{}({},{},R)".format(i,j))
245             elif (reward[i][j]==-5):
246                 dosya.write("{}({},{},B)".format(i,j))
247             dosya.write("\n")
248     dosya.close()
249

```

Burada ise, kullanıcı için bir dosyaya R matrisi kodun ilk çalışma aşamasında yazdırılmıştır. Bu matriste engeller, yollar, başlangıç ve bitiş konumları bulunmaktadır.

```

250 def ciktiVer() :
251     run = True
252     iterasyon = 0
253     step = 0
254     episodeStep = []
255     cost = 0
256     episodeCost = []
257     road=[]
258     probabilityRoad=[]
259
260     while iterasyon<1000:
261         screen.fill(background)
262         layout(False, probabilityRoad)
263         for event in pygame.event.get():
264             if event.type == pygame.QUIT:
265                 run = False
266                 pygame.display.flip()
267                 iterasyon, step, cost,episodeCost,episodeStep,road,probabilityRoad = episode(
268                     iterasyon, startPoint, step, cost, episodeStep, episodeCost,road,probabilityRoad)
269
270     print(probabilityRoad)
271     layout(True, probabilityRoad)
272     plot_results(episodeStep,episodeCost)
273
274     pygame.quit()
275

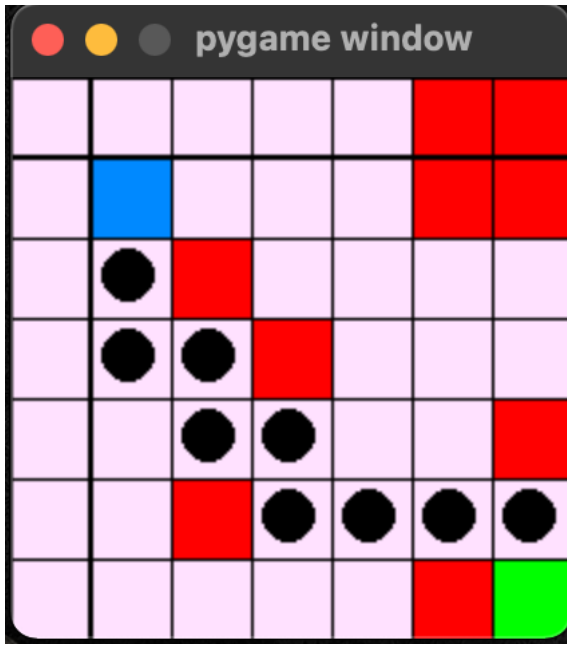
```

Burada ise, pekiştirmeli öğrenme işlemi belirli bir tekrar sayısına başlatılmış ve işlem bitince gerekli çıktılar ekrana sunulmuştur.

IV. KOD TASARIMI

Kodun akış diyagramı 5. Sayfadan sonra başlamaktadır.

V. DENEYSEL SONUÇLAR



KAYNAKÇA

WEB SİTE

<https://medium.com/deep-learning-turkiye/q-learning-giris-6742b3c5ed2b>

<https://towardsdatascience.com/simple-reinforcement-learning-q-learning-fcddc4b6fe56>

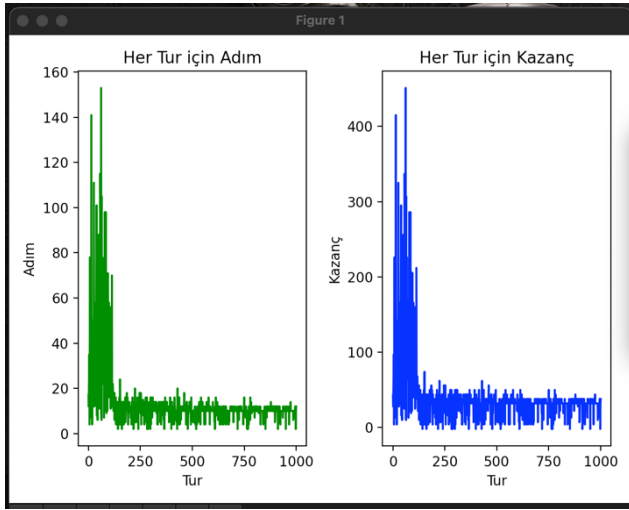
<https://medium.com/analytics-vidhya/introduction-to-reinforcement-learning-q-learning-by-maze-solving-example-c34039019317>

<https://medium.com/data-science-in-your-pocket/maze-runner-%EF%B8%8F-with-off-policy-q-learning-no-back-stepping-allowed-d01a79a6199c>

<https://github.com/erikdelange/Reinforcement-Learning-Maze>

<https://becominghuman.ai/q-learning-a-maneuver-of-mazes-885137e957e4>

https://github.com/sichkar-valentyn/Reinforcement_Learning_in_Python



VI. SONUÇ

Tüm bu işlemler ile ajan Q-learning algoritmasını kullanarak, brute-force yaklaşım ile bir labirent içerisinde başlangıçtan bitişe olan yolunu bulmuştur.

VII. KULLANILAN TEKNOLOJİLER

- Python
- Pygame
- PySimpleGUI
- Matplotlib
- GIT ve GITHUB
- Visual Studio Code tümleşik geliştirme ortamları kullanılmıştır.

VIII. KAZANIMLAR

- Q-learning algoritması hakkında bilgi sahibi olduk
- Brute-Force yaklaşım hakkında bilgi sahibi olduk
- Pygames, PySimpleGUI, Matplotlib kütüphaneleri hakkında bilgi sahibi olduk

