



دانشگاه اصفهان

دانشکده مهندسی کامپیوتر

درس مبانی و کاربردهای هوش مصنوعی
گزارش پروژه دوم (Frozen Lake Game)

نام و نام خانوادگی : مهدی رهبر

شماره دانشجویی : ۴۰۱۳۶۱۳۰۳۶

پاییز ۱۴۰۳

• فرآیند تصمیم‌گیری مارکوف (MDP)

فرآیند تصمیم‌گیری مارکوف (MDP) یک مدل ریاضی برای تصمیم‌گیری در محیط‌های نامعین است که شامل مجموعه‌ای از حالت‌ها، کنش‌ها، انتقال‌ها، پاداش‌ها و عامل تخفیف می‌باشد. دو الگوریتم رایج برای حل MDP عبارت‌اند از: Value Iteration که مستقیماً با بهینه‌سازی ارزش هر حالت، سیاست بهینه را استخراج می‌کند و Policy Iteration که با تکرار بین ارزیابی سیاست و بهبود آن، سیاست بهینه را می‌یابد. در این پروژه، این دو الگوریتم به عامل کمک می‌کنند تا بهترین مسیر را از حالت شروع به سمت هدف پیدا کند و از چاه‌ها اجتناب کند، در حالی که پاداش‌ها و احتمالات لغزش محیط نیز در نظر گرفته می‌شوند.

(۱) پیاده‌سازی الگوریتم Value Iteration:

```
3 usages
23 def create_reward_map(env):
24     reward_map = np.zeros(env.nS)
25     for s in range(env.nS):
26         row, col = np.unravel_index(s, env.shape)
27
28         if env._lake[row, col]:
29             reward_map[s] = -10
30         elif s == env.nS - 1:
31             reward_map[s] = 10
32         else:
33             reward_map[s] = -1
34     return reward_map
35
36
```

تابع `create_reward_map` نقشه‌ای از پاداش‌ها برای تمامی حالت‌های محیط ایجاد می‌کند. با پیمایش تمامی حالت‌ها، پاداش هر حالت بر اساس نوع آن تعیین می‌شود: چاه‌ها پاداش -۱۰، هدف پاداش ۱۰، و سایر خانه‌ها جریمه -۱ برای حرکت دارند. این تابع نقش کلیدی در الگوریتم‌های Value Iteration و Policy

Iteration دارد، زیرا با تعریف پاداش‌های محیط، به محاسبه ارزش حالت‌ها و استخراج سیاست بهینه کمک می‌کند.

```
38 def value_iteration(env, gamma=0.85, threshold=1e-6):
39
40     reward_map = create_reward_map(env)
41     value_table = np.zeros(env.nS)
42
43     slip_probs = {
44         0: [(0.5, 0), (0.25, 1), (0.25, 3)],
45         1: [(0.5, 1), (0.25, 0), (0.25, 2)],
46         2: [(0.5, 2), (0.25, 1), (0.25, 3)],
47         3: [(0.5, 3), (0.25, 0), (0.25, 2)]
48     }
49
50
51     iteration_count = 0
52
53     for _ in range(max_iter_number):
54         updated_value_table = np.copy(value_table)
55         iteration_count += 1
56
57         for s in range(env.nS):
58
59             Q_values = []
60             for a in range(env.nA):
61                 q_value = 0
62                 for slip_prob, slip_action in slip_probs[a]:
63                     for prob, s_, _, _ in env.P[s][slip_action]:
64                         if isinstance(s_, tuple):
65                             s_ = 8 * s_[0] + s_[1]
66                         q_value += slip_prob * (reward_map[s_] + gamma * updated_value_table[s_])
67             Q_values.append(q_value)
68             value_table[s] = max(Q_values)
69
70         if np.sum(np.abs(updated_value_table - value_table)) <= threshold:
71             break
72
73     print(f"Value Iteration converged in {iteration_count} iterations.")
74     return value_table
75
76
```

تابع **value_iteration** الگوریتم تکرار ارزش را برای یافتن سیاست بهینه در محیط اجرا می‌کند. ابتدا نقشه پاداش‌ها (**reward_map**) و جدول ارزش اولیه مقداردهی می‌شوند. در هر تکرار، با استفاده از احتمالات لغزش تعریف‌شده، مقدار Q برای هر کنش محاسبه و حداکثر آن به‌عنوان ارزش جدید حالت ثبت می‌شود. فرآیند تا زمانی ادامه می‌یابد که تغییرات جدول ارزش کمتر از آستانه همگرایی باشد یا به حداکثر تعداد تکرار برسد. این تابع ارزش نهایی هر حالت را بازمی‌گرداند.

```

78 def extract_policy(env, value_table, gamma=0.85):
79     policy = np.zeros(env.nS, dtype=int)
80
81     reward_map = create_reward_map(env)
82
83     slip_probs = {
84         0: [(0.5, 0), (0.25, 1), (0.25, 3)],
85         1: [(0.5, 1), (0.25, 0), (0.25, 2)],
86         2: [(0.5, 2), (0.25, 1), (0.25, 3)],
87         3: [(0.5, 3), (0.25, 0), (0.25, 2)]
88     }
89
90     for s in range(env.nS):
91         Q_values = []
92         for a in range(env.nA):
93             q_value = 0
94             for slip_prob, slip_action in slip_probs[a]:
95                 for prob, s_, _, _ in env.P[s][slip_action]:
96
97                     if isinstance(s_, tuple):
98                         s_ = 8 * s_[0] + s_[1]
99
100                    q_value += slip_prob * (reward_map[s] + gamma * value_table[s_])
101                Q_values.append(q_value)
102            policy[s] = np.argmax(Q_values)
103
104     return policy
105

```

تابع **extract_policy** سیاست بهینه را از جدول ارزش (value_table) استخراج می‌کند. برای هر حالت، مقادیر Q برای تمامی کنش‌ها با استفاده از احتمالات لغزش، نقشه پاداش‌ها (reward_map) و ارزش حالت‌های بعدی محاسبه می‌شود. سپس کنشی که بیشترین مقدار Q را دارد، به‌عنوان سیاست بهینه برای آن حالت انتخاب می‌شود. این تابع نقش کلیدی در هر دو الگوریتم Value Iteration و Policy Iteration دارد، زیرا بر اساس ارزش‌های محاسبه‌شده، سیاست بهینه را برای هدایت عامل به سمت هدف استخراج می‌کند.

```

108 def plot_heatmap(value_table, env_shape):
109
110     reshaped_table = value_table.reshape(env_shape)
111
112     plt.figure(figsize=(8, 8))
113     plt.imshow(reshaped_table, cmap='RdYlGn', interpolation='nearest')
114     plt.colorbar(label='Value')
115     plt.title('Value Heatmap')
116
117     for i in range(env_shape[0]):
118         for j in range(env_shape[1]):
119             plt.text(j, i, s=f'{reshaped_table[i, j]:.2f}',
120                     ha='center', va='center', color='black', fontsize=10)
121
122     plt.show()
123

```

تابع **plot_heatmap** یک نمودار Heatmap از جدول ارزش (value_table) ترسیم می‌کند تا ارزش حالت‌ها را در محیط به صورت تصویری نمایش دهد. ابتدا جدول ارزش به شکل محیط تغییر داده می‌شود و با استفاده از کتابخانه **Matplotlib**، مقادیر ارزش به صورت رنگی (سبز برای مقادیر بالا و قرمز برای مقادیر پایین) نمایش داده می‌شود. همچنین، مقدار عددی هر خانه در مرکز آن درج می‌شود. این تابع نقش مهمی در تحلیل و ارزیابی الگوریتم‌ها دارد، زیرا وضعیت ارزش حالت‌ها را به صورت گرافیکی نمایش می‌دهد و به درک بهتر نتایج کمک می‌کند.

۲) پیاده سازی الگوریتم Policy Iteration:

```
127 def compute_value_function(policy, env, gamma=0.85, threshold=1e-6):
128
129     value_table = np.zeros(env.nS)
130
131     reward_map = create_reward_map(env)
132
133     slip_probs = {
134         0: [(0.5, 0), (0.25, 1), (0.25, 3)],
135         1: [(0.5, 1), (0.25, 0), (0.25, 2)],
136         2: [(0.5, 2), (0.25, 1), (0.25, 3)],
137         3: [(0.5, 3), (0.25, 0), (0.25, 2)]
138     }
139
140     while True:
141         updated_value_table = np.copy(value_table)
142
143         for state in range(env.nS):
144
145             action = policy[state]
146             q_value = 0
147
148             for slip_prob, slip_action in slip_probs[action]:
149
150                 for trans_prob, next_state, _, _ in env.P[state][slip_action]:
151
152                     if isinstance(next_state, tuple):
153                         next_state = 8 * next_state[0] + next_state[1]
154
155                     q_value += slip_prob * (reward_map[state] + gamma * updated_value_table[next_state])
156
157             value_table[state] = q_value
158
159             if np.sum(np.abs(updated_value_table - value_table)) <= threshold:
160                 break
161
162     return value_table
163
```

تابع **compute_value_function** ارزش حالت‌ها را برای یک سیاست داده‌شده محاسبه می‌کند. ابتدا جدول ارزش با مقادیر صفر مقداردهی می‌شود و با استفاده از نقشه پاداش (**reward_map**) و احتمالات لغزش، مقدار **Q** برای هر حالت محاسبه می‌شود. در هر تکرار، ارزش حالت فعلی بر اساس کنش انتخاب‌شده در سیاست، پاداش حالت فعلی و ارزش حالت‌های بعدی به‌روزرسانی می‌شود. فرآیند تا زمانی ادامه می‌یابد که تغییرات ارزش‌ها کمتر از مقدار آستانه (**threshold**) شود. این تابع در الگوریتم **Policy Iteration** نقش ارزیابی سیاست (**Policy Evaluation**) را ایفا می‌کند و ارزش حالت‌ها را برای سیاست فعلی برآورد می‌کند.

```

167 def policy_iteration(env, gamma=0.85, threshold=1e-6):
168
169     policy = np.zeros(env.nS, dtype=int)
170
171     iteration_count = 0
172
173     for i in range(max_iter_number):
174
175         iteration_count += 1
176
177         value_table = compute_value_function(policy, env, gamma, threshold)
178
179         new_policy = extract_policy(env, value_table, gamma)
180
181         if np.array_equal(policy, new_policy):
182             print(f"Policy Iteration converged in {iteration_count} iterations.")
183             break
184
185         policy = new_policy
186
187     return policy, value_table
188
189

```

تابع **policy_iteration** الگوریتم تکرار سیاست را برای یافتن سیاست بهینه اجرا می‌کند. ابتدا سیاست اولیه مقداره‌ی می‌شود و در هر تکرار، ارزش حالت‌ها برای سیاست فعلی با استفاده از `compute_value_function` محاسبه می‌شود. سپس سیاست جدید با استفاده از `extract_policy` استخراج می‌گردد. اگر سیاست جدید با سیاست قبلی یکسان باشد (همگرایی رخ داده باشد)، الگوریتم متوقف می‌شود و تعداد تکرارها چاپ می‌شود. در غیر این صورت، سیاست به‌روزرسانی شده و تکرار ادامه می‌یابد. این تابع با ارزیابی و بهبود متوالی سیاست، به سیاست بهینه و ارزش متناظر آن دست می‌یابد.

```

192  def run_value_iteration():
193      global observation
194      optimal_value_table = value_iteration(env)
195      optimal_policy = extract_policy(env, optimal_value_table)
196      plot_heatmap(optimal_value_table, env.shape)
197
198      while True:
199          action = optimal_policy[observation]
200          next_state, reward, done, truncated, info = env.step(action)
201
202          observation = 8 * next_state[0] + next_state[1]
203
204          if observation == 63:
205              print("You reached the goal! Returning to the main menu.")
206              env.close()
207              main_menu()
208              break
209          elif env._lake[next_state]:
210              print("You fell into a hole! Restarting the game.")
211              observation, info = env.reset()

```

```

215  def run_policy_iteration():
216      global observation
217
218      optimal_policy, value_table = policy_iteration(env)
219
220      plot_heatmap(value_table, env.shape)
221
222      while True:
223          action = optimal_policy[observation]
224          next_state, reward, done, truncated, info = env.step(action)
225
226          observation = 8 * next_state[0] + next_state[1]
227
228          if observation == 63:
229              print("You reached the goal! Returning to the main menu.")
230              env.close()
231              main_menu()
232              break
233          elif env._lake[next_state]:
234              print("You fell into a hole! Restarting the game.")
235              observation, info = env.reset()

```



```

240 def main_menu():
241     global env, observation
242
243     def start_value_iteration():
244         root.destroy()
245         global env, observation
246         env = FrozenLake(render_mode="human", map_name="8x8")
247         observation, info = env.reset(seed=30)
248         run_value_iteration()
249
250     def start_policy_iteration():
251         root.destroy()
252         global env, observation
253         env = FrozenLake(render_mode="human", map_name="8x8")
254         observation, info = env.reset(seed=30)
255         run_policy_iteration()
256
257
258     root = tk.Tk()
259     root.title("Main Menu")
260     root.geometry("400x200")
261
262
263     tk.Label(root, text="Choose Algorithm", font=("Arial", 16)).pack(pady=20)
264     tk.Button(root, text="Value Iteration", font=("Arial", 14), command=start_value_iteration).pack(pady=10)
265     tk.Button(root, text="Policy Iteration", font=("Arial", 14), command=start_policy_iteration).pack(pady=10)
266
267     root.mainloop()

```

سه تابع بالا برای مدیریت اجرا الگوریتم‌ها در منو پیاده سازی شده اند.

خب بعد از بررسی کدها پروژه را اجرا می‌کنیم. وقتی پروژه اجرا می‌شود یک منو ظاهر شده و می‌توانیم از بین دو الگوریتم MDP انتخاب کنیم. با انتخاب هر کدام از الگوریتم‌ها ابتدا نقشه Heatmap نمایش داده می‌شود. سپس با بستن نقشه، عامل (تاکسی) شروع به جست‌وجو کرده و به هدف می‌رسد. البته به دلیل انتخاب های تصادفی امکان افتادن در چاه‌ها هم وجود دارد. به صورت کلی عملکرد دو الگوریتم با بررسی خطاها و پیروزی‌ها خوب ارزیابی می‌شود با این تفاوت که الگوریتم policy iteration با تعداد تکرار به مراتب کم‌تری همگرا می‌شود. در ادامه تصاویر خروجی و نقشه heatmap نشان داده شده است:

