

به نام خدا



دانشگاه اصفهان

دانشکده مهندسی کامپیوتر

محمد کاظم هرندی

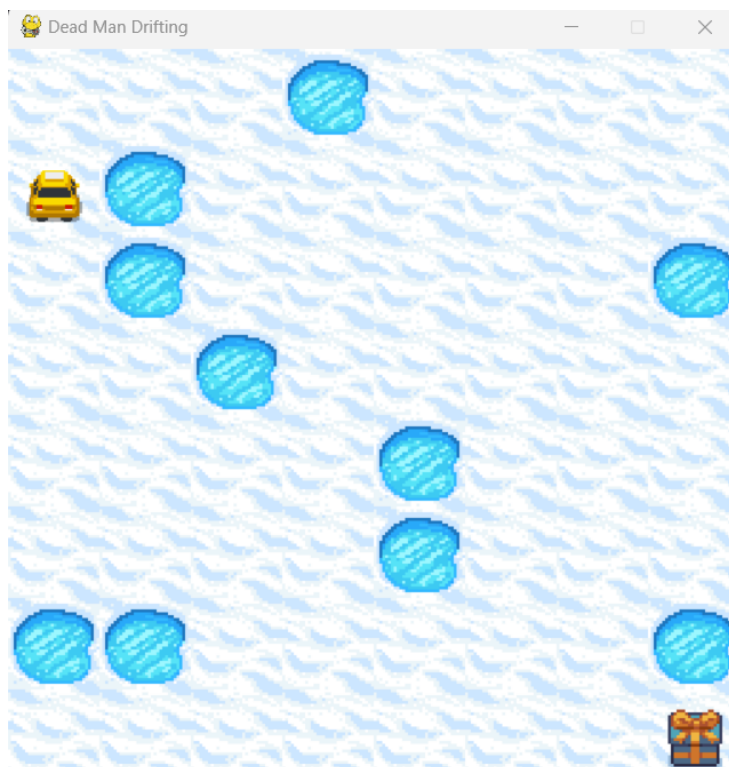
۴۰۰۳۶۲۳۰۳۹

مبانی و کاربردهای هوش مصنوعی

گزارشکار پروژه دوم

پاییز ۱۴۰۳

معرفی پروژه



در این پروژه، هدف ما پیاده‌سازی یک عامل هوشمند برای مسیردهی در محیط Frozen Lake با استفاده از الگوریتم تصمیم‌گیری مارکوف بود. محیط Frozen Lake یک شبکه‌ای ۸ خانه‌ای است که عامل (تاکسی یا الف) باید از خانه شروع (خانه ۰) به خانه هدف (خانه ۶۳) برسد. به دلیل لغزنده بودن زمین، احتمال دارد تاکسی به جهت‌های دیگر نیز برود. در این پروژه از کتابخانه‌های gymnasium و numpy استفاده شده است.

فرآیند تکرار سیاست

فرآیند تکرار سیاست یکی از الگوریتم‌های کلاسیک در حل مسائل تصمیم‌گیری مارکوف است. این فرآیند شامل دو مرحله اصلی است: ارزیابی سیاست و بهبود سیاست. در ادامه به توضیح کامل هر یک از این مراحل می‌پردازیم.

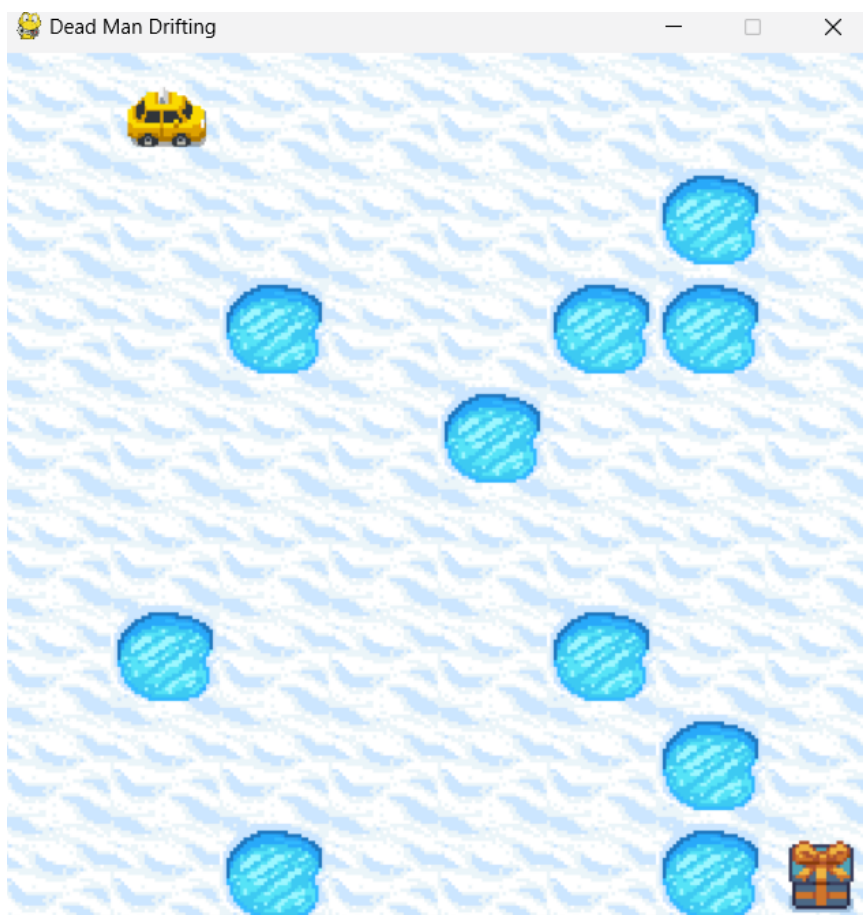
ارزیابی سیاست

هدف از ارزیابی سیاست، محاسبه مقداردهی به حالات (v) برای یک سیاست مشخص است. این کار با استفاده از معادله بلمن انجام می‌شود. معادله بلمن رابطه‌ای بین مقداردهی به یک حالت و مقداردهی به حالات بعدی آن برقرار می‌کند. در این مرحله، ما به دنبال یافتن مقداردهی پایدار (یا نزدیک به پایدار) برای هر حالت هستیم.

```
90 v = np.zeros(num_state) # Initialize the current value function to zeros
91 for i in range(1000): # Limit the number of iterations for policy iteration
92
93     # Policy Evaluation step
94     v_old = -1 * np.ones(num_state) # Reset old value function to ensure the loop condition is met
95
96     while (np.abs(v - v_old)).max() > 1e-3: # Continue until value function convergence
97         v_old = v.copy() # Update the old value function
98
99         for s in range(num_state): # For each state in the environment
100             vs = 0 # Temporary variable for storing the value of state s
101
102             # Calculate the expected value of the current policy
103             for prob, st, r, done in env.P[s][policy[s]]:
104                 vs += prob * (r + gamma * v[st]) # Bellman equation for policy evaluation
105             v[s] = vs # Update the value of state s
106
107         print(v) # Print the current value function for debugging/observation
108
109     # Policy Improvement step
110     for s in range(num_state): # For each state in the environment
111
112         q = np.zeros(num_action) # Initialize a temporary Q-value array for all actions
113
114         # For each action, calculate its Q-value
115         for a in range(num_action):
116             for prob, st, r, done in env.P[s][a]:
117                 q[a] += prob * (r + gamma * v[st]) # Bellman equation for Q-value
118
119         policy[s] = np.argmax(q) # Update the policy to choose the action with the highest Q-value
```

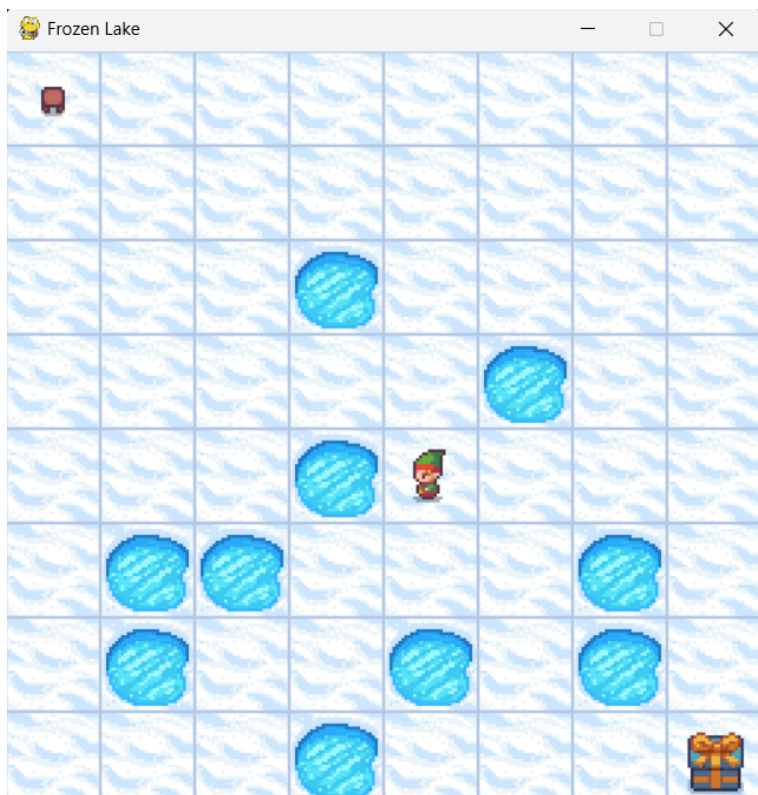
بهبود سیاست

در مرحله بهبود سیاست، هدف ما بهبود سیاست موجود با استفاده از مقداردهی به حالات به دست آمده در مرحله ارزیابی سیاست است. در این مرحله، برای هر حالت بهترین عمل را که بیشترین مقدار-Q Value را فراهم می‌کند، انتخاب می‌کنیم.



اجرای سیاست بهینه

پس از یافتن سیاست بهینه، محیط با استفاده از این سیاست اجرا شده و تاکسی از حالت اولیه تا حالت نهایی حرکت می‌کند.



گسترش پروژه

تغییر گرافیک بازی

برای بهبود تجربه کاربری و جذاب تر کردن بازی، می‌توان از کتابخانه Pygame برای تغییر گرافیک بازی استفاده کرد. با تغییر تصاویر و افزودن انیمیشن‌های جذاب، بازی برای کاربران لذت‌بخش‌تر خواهد شد.

پیاده‌سازی در گیت‌هاب

با ایجاد مخزن در گیت‌هاب برای پروژه، می‌توان کدها و مستندات پروژه را به اشتراک گذاشت. این کار امکان همکاری با دیگران و دریافت بازخورد را فراهم می‌کند. همچنین به مستندسازی پروژه کمک

می‌کند و باعث می‌شود که پروژه به صورت منظم‌تر و مستندتر باشد.

بهبود الگوریتم

علاوه بر الگوریتم تکرار سیاست، می‌توان الگوریتم‌های پیشرفته‌تری مانند یادگیری تقویتی (Reinforcement Learning) را برای بهبود عملکرد عامل هوشمند پیاده‌سازی کرد. این الگوریتم‌ها می‌توانند با استفاده از تجربه‌های گذشته، استراتژی‌های بهتری برای رسیدن به هدف پیدا کنند. از جمله الگوریتم‌های پیشرفته می‌توان به Q-Learning و SARSA اشاره کرد.

آزمایش و ارزیابی

با انجام آزمایش‌های مختلف و ارزیابی عملکرد عامل در شرایط مختلف، می‌توان نقاط قوت و ضعف الگوریتم را شناسایی کرد و بهبودهای لازم را انجام داد. این آزمایش‌ها می‌توانند شامل تغییرات در پاداش‌ها و احتمالات انتقال باشند. با استفاده از ابزارهای بصری‌سازی مانند نمودارها، می‌توان عملکرد عامل را در شرایط مختلف به صورت گرافیکی نمایش داد.