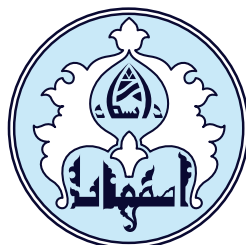


بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه اصفهان

دانشکده مهندسی کامپیوتر

فاز دوم پروژه در مبانی هوش و کاربردها

# فرآیند تصمیم مارکوف و محیط ناشناخته

استاد درس: دکتر حسین کارشناس

دستیار استاد: پوریا صامتی

دانیال شفیعی  
مهدی مهدیه  
سید امیررضا نجفی

آذر ۱۴۰۳

# فهرست مطالب

|   |       |   |
|---|-------|---|
| ۲ | ۱     | فرآیند تصمیم مارکوف                     |
| ۲ | ۱.۱   | مقدمه                                   |
| ۲ | ۲.۱   | خوشه‌بندی                               |
| ۲ | ۱.۲.۱ | انتخاب تعداد خوشه‌ها                    |
| ۳ | ۲.۲.۱ | همگرایی خوشه‌بندی و تشکیل نقشه‌ی گرمایی |
| ۳ | ۳.۱   | خروجی الگوریتم                          |
| ۵ | ۲     | محیط ناشناخته                           |
| ۵ | ۱.۲   | مقدمه                                   |
| ۵ | ۲.۲   | شیوه‌ی پیاده‌سازی الگوریتم QLearning    |
| ۵ | ۱.۲.۲ | ابزارهای آموزش                          |
| ۵ | ۲.۲.۲ | همگرایی مقادیر                          |
| ۶ | ۳.۲.۲ | استخراج سیاست                           |
| ۷ | ۴.۲.۲ | سایر ویژگی‌ها                           |
| ۸ | ۳.۲   | خروجی                                   |
| ۸ | ۴.۲   | شبکه‌ی عصبی                             |
| ۸ | ۱.۴.۲ | مقدمه                                   |
| ۸ | ۲.۴.۲ | پیاده‌سازی                              |

## فصل ۱

# فرآیند تصمیم مارکوف

### ۱.۱ مقدمه

برای بخش MDP<sup>۱</sup> در ابتدا تصمیم گرفتیم  $v\_table$  استخراج کنیم که حاوی راه رسیدن به یک خوک هستند. بعد از خوردن هر خوک به  $v\_table$  رفته و طبق سیاست‌هایی که از آن استنتاج کردیم پیش می‌رویم. این روش مشکلاتی را در پی داشت؛ مثل اینکه ممکن بود یک خوک توسط دیوارها محاصره شده باشد و یا اینکه ممکن بود در مسیر خوردن خوک مورد نظر، بعضی از خوک‌های دیگر نیز به دلیل تصادفی بودن حرکات خورده شوند و مشکلات دیگری که برای حل آنها نیاز به تعداد زیادی شرط بود. این شروط حس را می‌داد که دیگر مسئله را با فرآیند تصمیم مارکوف به صورت عمومی حل نمی‌کنیم و داریم با جزئی نگری زیادی به مسئله نگاه می‌کنیم. این جزئی نگری شاید در بعضی محیط‌ها به نفع ما باشد ولی ممکن است با تغییر محیط برایمان بسیار گران تمام شود.

### ۲.۱ خوشه‌بندی

ما تصمیم گرفتیم با استفاده از خوشه‌بندی، خوک‌های موجود در صفحه را به ۲، ۳ یا ۴ خوشه تقسیم می‌کنیم. سپس با مقایسه، خوشه‌ای که در آن خوک‌ها بهترین نماینده را درون آن خوشه دارند انتخاب کنیم. آنگاه برای هر خوشه یک  $v\_table$  می‌سازیم که به خوک‌های موجود در آن خوشه امتیاز مثبت می‌دهیم ولی خوک‌هایی که در آن خوشه نیستند را نادیده می‌گیریم. سپس عامل در هر خوشه تا جای ممکن خوک‌های موجود را خورده و به سمت خوشه‌ی بعدی حرکت می‌کند و این روند همین‌گونه ادامه پیدا می‌کند تا زمانی که به خوشه‌ی آخر (تخم مرغ) که هدف نهایی است برسیم.

#### ۱.۲.۱ انتخاب تعداد خوشه‌ها

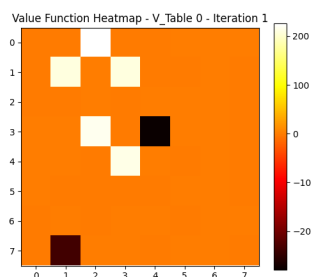
سوال کلیدی این است که الگوریتم چگونه متوجه می‌شود به چه ترتیبی  $v\_table$  ها را پیمایش کند؟ یا به زبان ساده تر بهتر است ابتدا سیاستهای متعلق به اولین خوشه را اجرا کنیم یا دومین خوشه را؟ برای مقایسه، روی تمام مراکز خوشه‌ها جایگشت را حساب می‌کنیم و بهترین جایگشت را انتخاب می‌کنیم. در بیشترین حالت اگر ۴ خوشه انتخاب کرده باشیم، ۲۴ حالت دارد پس از نظر زمانی نیز لطمه‌ای به الگوریتم نمی‌زند.

---

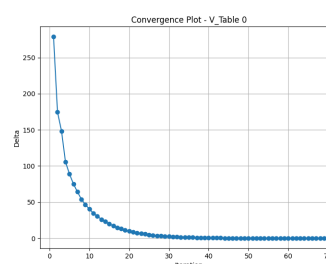
<sup>۱</sup>Markov Decision Process

## ۲.۲.۱ همگرایی خوشه‌بندی و تشکیل نقشه‌ی گرمایی

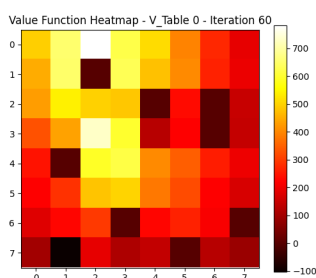
برای نمونه ما نقشه‌گرمایی محیط راه در خوشه‌بندی به ازای  $k$  برابر ۳ در چند تکرار مختلف نمایش می‌دهیم.



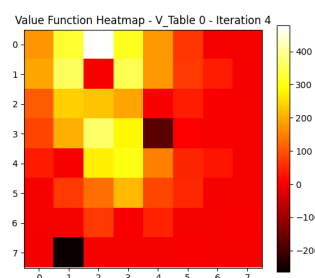
(ب) تکرار اول



(آ) همگرایی خوشه‌بندی



(د) تکرار شصتم



(ج) تکرار چهارم

شکل ۱.۱: نمونه‌ای از همگرایی خوشه‌بندی به ازای  $k$  برابر ۳

## ۳.۱ خروجی الگوریتم

در ادامه خروجی از تحلیل الگوریتم روی یک نمونه از محیط را می‌بینیم.

---

Best clustering with 4 clusters:

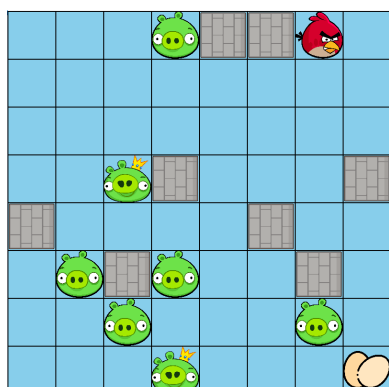
Cluster 1: [(5, 1), (5, 3), (6, 2)], Centroid: 5.33, 2.0

Cluster 2: [(0, 1), (0, 3)], Centroid: 0.0, 2.0

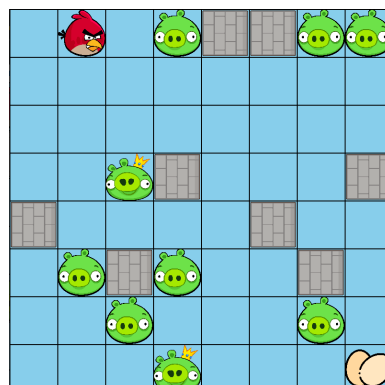
Cluster 3: [(6, 6)], Centroid: 6.0, 6.0

Cluster 4: [(0, 6), (0, 7)], Centroid: 0.0, 6.5

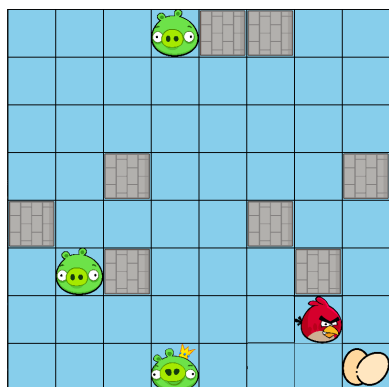
---



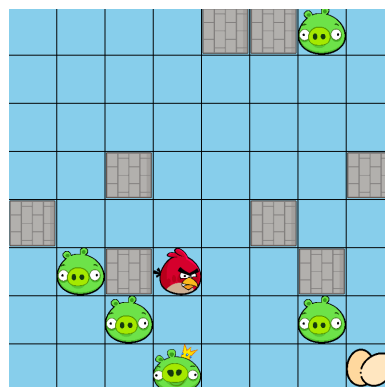
(ب) حرکت به سمت خوشه‌ی اول



(آ) شروع



(د) حرکت به سمت خوشه‌ی سوم (تخم مرغ)



(ج) حرکت به سمت خوشه‌ی دوم

شکل ۲۰۱: روند حرکت پرنده‌ی خشمگین درون یک محیط تصادفی نمونه

## فصل ۲

# محیط ناشناخته

### ۱.۲ مقدمه

در این بخش از پروژه ما قصد داشتیم با تعداد محدودی تکرار و گردش<sup>۱</sup> در محیط ناشناخته<sup>۲</sup> به یک قاعده برسیم که در آن بتوانیم با استفاده از، به نحوی در این محیط تصادفی پرنده را هدایت کنیم که بیشترین امتیاز ممکن را کسب کند. نکات قابل توجه در این پروژه این است که نباید به TNT برخورد کنیم و تعداد حرکات هم محدود است و نمی‌توانیم بی‌نهایت در محیط بگردیم تا همه‌ی حالت‌ها را بدست بیاوریم.

### ۲.۲ شیوه‌ی پیاده‌سازی الگوریتم QLearning

در مرحله‌ی اول ما یک ماتریس ۸ در ۸ در ۴ داشتیم که امتیاز هر کنش در هر حالت را مشخص می‌کرد. حالت‌ها نمایان‌گر نقطه‌ی که پرنده در آن وجود داشت بودند.

#### ۱.۲.۲ ابرپارامترهای آموزش

برای تغییرات **مقادیر جدول Q** یک سری ابر پارامتر وجود دارد.  $\gamma$  در کد ما تحت عنوان فاکتور تخفیف معرفی شده است و  $\alpha$  تحت عنوان سرعت آموزش.

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)] \quad (1.2)$$

در هر اپیزود از اجرای الگوریتم، حالت به عنوان ورودی داده می‌شود و سپس بهترین کنش از جدول انتخاب می‌شود. سپس با انجام کنش به کمک تابع step می‌دهیم و خروجی آن را دریافت می‌کنیم و توابع و حالات جدید را بدست می‌آوریم.

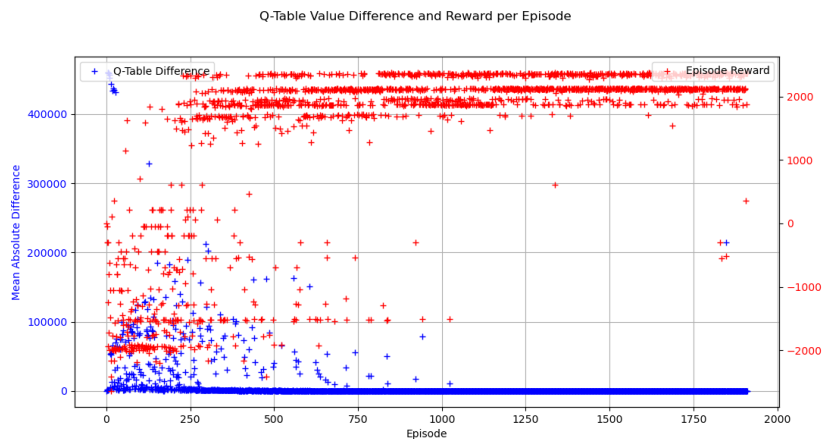
#### ۲.۲.۲ همگرایی مقادیر

برای اینکه بعد از مدتی همگرا شود، ما از تکنیک **Epsilon Greedy** استفاده کردیم. در این راهبرد. در اوایل اجرا، پرنده تا جای ممکن تصادفی عمل می‌کند. اما به مرور زمان با توجه به نتایج به دست

---

<sup>۱</sup>Explore  
<sup>۲</sup>Environment Unknown

آمده بهترین کنش را انتخاب می‌کند. نرخ کاهش نیز بدین صورت است که تا زمانی که نتیجه در هر حال بهتر شدن است، همچنان اپسیلون و نرخ یادگیری تغییری نمی‌کند اما به محض کاهش پاداش کل برای اپسیلون گریدی و افزایش value\_diff برای سرعت یادگیری، مقادیر کاهش پیدا می‌کنند.



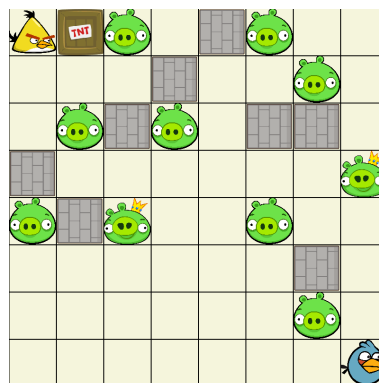
شکل ۱۰۲: همگرایی value\_diff در یک محیط نمونه که پس از ۱۰ بار کوچک‌تر شدن اپسیلون از ۱۰ الگوریتم متوقف شده است.

### ۳.۲.۲ استخراج سیاست

در ابتدای کار وقتی شرایط زنده بودن خوک‌ها را در نظر نگرفته بودیم، می‌توانستیم یک سیاست واحد استخراج کنیم اما بعد از آن، ما از بین اکشن‌ها مد می‌گیریم تا یک حس از اینکه در هر حالت به چه سمتی حرکت می‌کند مشخص شود. برای خانه‌هایی که مقادیرشان مشخص نیست، الگوریتم به صورت تصادفی سیاست استخراج می‌کند.

Policy Rule after train Q-Learning (Mode of actions)

|   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|
| ← | ↑ | → | ↑ | ↑ | ← | ↓ | ↓ |
| ↓ | ↓ | ↓ | ← | → | ↑ | ← | ← |
| ← | ↓ | → | ↓ | ← | ↑ | ← | ← |
| → | ← | ↓ | ↑ | → | → | ← | ← |
| ↑ | ↑ | ↓ | → | ↓ | ← | ← | → |
| → | → | → | → | ↑ | ← | ↓ |   |
| ↑ | → | ← | → | → | → | ← | → |
| ↑ | ↑ | → | → | → | ← | ↑ | ← |



(ب) مد کنش در حالت‌ها در محیط متناظر

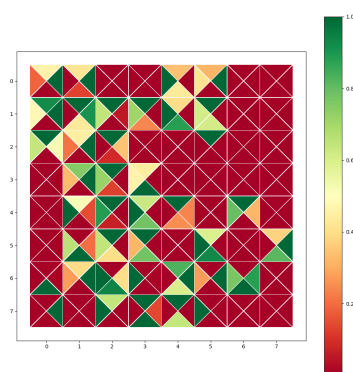
(آ) محیط نمونه

شکل ۲۰۲: استخراج سیاست به شیوه‌ی مد گرفتن از کنش‌ها در هر حالت

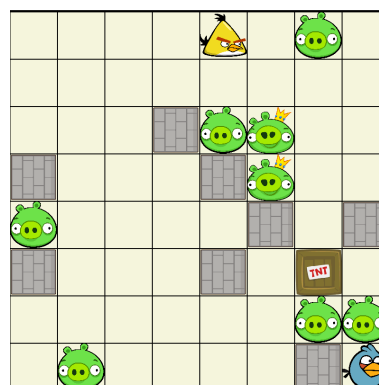
## ۴.۲.۲ سایر ویژگی‌ها

**مهندسی پاداش** وقتی الگوریتم به پایان می‌رسد، یک پاداش ۱۰۰۰- داریم که اگر این پاداش روی یک حالت و کنش برود، حتی اگر آن کنش خوب باشد، باعث می‌شود دیگر سمت این کنش نرویم! ما باید این اثر را حذف کنیم. همچنین چنین در نظر گرفته بودیم که اگر یک حالت و کنشی در یک اپیزود دفعات مکرر اجرا شد، به آن امتیاز منفی بدهیم که البته این باعث همگرایی نمی‌شود و ممکن است نتایج بدتری بدهد. همچنین برای اینکه از ملکه‌ها رد شویم و در دام آن‌ها نیفتیم، فرض کردیم اگر از آن‌ها بگذریم، امتیاز منفی کمتری دارد که البته این هم حالب نبود. در نهایت هم مقادیر جدول Q را به ازای همه‌ی کنش‌های آن حالت بروز می‌کردیم تا هیچ کنشی حتی کنش مقداره‌ی نشده هم به سمت TNT نرود.

**مصورسازی** برای نمایش اینکه در هر حالت چه کنشی بیشتر موردتوجه است و چه کنشی مضر است، از یک نقشه گرمایی استفاده کردیم که جهت‌های ممکن را به ما نشان می‌دهد. به این منظور کنش‌های هر حالت نرمالایز می‌شود و بر اساس مقادیر ۰ تا یک رنگ پیدا می‌کند. خانه‌هایی که همه‌ی آن‌ها قرمز هستند یعنی هیچ‌گاه مقداره‌ی نشده‌اند و خانه‌هایی که هر سه خانه سبز هستند یعنی یکی از خانه‌ها مقدار منفی دارد و احتمالاً قدر مطلق آن بسیار زیاد است. واضح است خانه‌ای که TNT در آن موجود است از هر ۴ جهت قرمز شده است.



(ب) نمونه نقشه گرمایی در حالی که هیچ‌کدام از خوک‌ها خورده نشده‌اند



(آ) محیط نمونه

شکل ۳.۲: نقشه گرمایی جهت حرکت بر اساس جدول Q

**اثربخشی از حالت خوک‌ها** به کمک این ویژگی متناسب با اینکه هر خوکی خورده شده است یا نه یک جدول سیاست کنش داریم. یعنی به جای ۱ جدول، ۲۵۶ جدول داریم که البته بسیاری از آن‌ها هرگز مقداره‌ی نمی‌شوند ولی همین‌هایی که پر می‌شوند به ما کمک می‌کنند تا امتیاز بسیار بالایی کسب کنیم و به‌خاطر یک خوک خورده شده گرفتار حلقه نشویم و محدودیت حرکتی نخوریم.



## ۳.۲ خروجی

نمونه‌ی خروجی این الگوریتم برای محیط نسبتاً سختی که در شکل ۳.۲ نمایش داده شده است به نحو زیر می‌باشد:

```
[False, False, False, False, False, False, False, False]
Episode finished with reward: 2347
[False, False, False, False, False, False, False, False]
Episode finished with reward: 2336
[False, False, False, False, False, False, False, False]
Episode finished with reward: 2285
[False, False, False, False, False, False, False, False]
Episode finished with reward: 2329
[False, False, False, False, False, False, False, False]
Episode finished with reward: 2303
MEAN REWARD: 2320.0
```

در هر ۵ حالت همه‌ی خوک‌ها خورده شده‌اند و امتیاز بالایی کسب شده است.

## ۴.۲ شبکه‌ی عصبی

### ۱.۴.۲ مقدمه

Q-Learning عمیق یا شبکه‌ی Q عمیق<sup>۳</sup> همان توسعه یافته‌ی الگوریتم Q-Learning است که در آن از شبکه‌های عصبی برای یادگیری عمیق و تقریب Q استفاده می‌شود. Q-Learning سنتی برای محیط‌هایی با تعداد حالت‌های کم و محدود به خوبی کار می‌کند، اما به دلیل اندازه جدول Q با فضاهای حالت بزرگ یا پیوسته مشکل دارد. Q-Learning عمیق با جایگزین کردن جدول Q با یک شبکه عصبی که می‌تواند مقادیر Q را برای هر جفت حالت-عمل تقریبی بزند و بر این محدودیت غلبه کند.

### ۲.۴.۲ پیاده‌سازی

ساختار شبکه لایه اول: ورودی با اندازه ۶۴ (که نشان‌دهنده‌ی وضعیت محیط ۸ در ۸ است) به ۱۲۸ نورون متصل می‌شود. لایه دوم: ۱۲۸ نورون به ۶۴ نورون متصل می‌شود. لایه خروجی: ۶۴ نورون به تعداد کنش‌های ممکن (بالا، پایین، چپ، راست) متصل می‌شود. تابع فعال‌سازی: از تابع ReLU برای افزودن غیرخطی بودن به شبکه استفاده می‌شود.

اپیزودها ما یک بهینه‌ساز و تابع خطا تعریف می‌شود. همچنین وضعیت محیط به بردار ۸ در ۸ تبدیل می‌شود. که این بردار ورودی شبکه عصبی است. هر سلول بر اساس پاداش خود آن مقدار خاصی به بردار اضافه می‌کند. سپس یک کنش انتخاب می‌شود و یک قدم برداشته می‌شود. بعد با استفاده از راهبر پس‌یون گزینی یک اقدام تصادفی یا اقدامی با بیشترین مقدار انتخاب می‌شود. سپس مدل بر اساس پاداش تجربیات خود را ذخیره می‌کند و وزن‌های شبکه‌ی عصبی بروزرسانی می‌شود. سیاست نهایی با انتخاب اقدام با بیشترین مقدار Q برای هر وضعیت ممکن استخراج می‌شود.

<sup>۳</sup> Deep Q Network