

بسمه تعالی



تمرین چهارم

## یادگیری تقویتی

استاد درس: دکتر آرش عبدی هیراندوست

تدریسار آموزشی: آرمین توکلی

✓ نکات تمرین :

- ❖ مهلت تعویل ۱۴۰۲/۰۲/۲۶ ساعت ۲۳:۵۵
- ❖ مهلت ارسال به هیچ وجه قابل تغییر نیست .
- ❖ مواردی که بعد از تاریخ فوق ارسال شوند قابل قبول نبوده و نمره ای نخواهند داشت .
- ❖ انجام تمرین تک نفره است. لطفا به تنهایی انجام شود، در غیر اینصورت نمره منفی در نظر گرفته خواهد شد .

❖ کل محتوای ارسالی زیپ شود و نام خایل زیپ ارسالی HW4\_StudentNumber\_YourName باشد .

❖ محتوای ارسالی دارای راهنما (read me) جهت تسهیل اجرا باشد .

❖ زبان برنامه نویسی دلخواه است. (پیشنهاد : پایتون)

❖ در صورت استفاده از زبان پایتون خایل که ترجیحا به فرمت **ipynb** بوده و خایل کد هتما به

صورت اجرا شده آپلود گردد و از وجود فروبی سلول ها اطمینان حاصل نمایید .

❖ موارد ارسال شده در تاریفی که بعدا مشفص خواهد شد و متعاقبا اعلام می گردد به صورت آنلاین

نیز تمویل گرفته خواهند شد (صرفا آنچه در LMS طبق تاریخ فوق تمویل داده شده است بعدا به

صورت حضوری تست شده و توضیح داده می شود).

❖ تنها تکالیفی که به LMS و قبل از مهلت ارسال، فرستاده می شوند بررسی خواهند شد .

❖ در صورت داشتن هر گونه سوال می توانید سوال خود را در گروه تلگرامی درس مطرح کنید .

❖ حداقل یک ساعت قبل از مهلت ارسال را احتیاطا هدف قرار دهید، تا مشکلات غیرقابل پیش

بینی مانند موارد زیر باعث عدم آپلود پاسخ ها در LMS و ارسال آنها از طریق ایمیل نشوند :

(قطعی اینترنت، تنظیم نبودن دقیق ساعت سایت با ساعت کرنویچ، کرش سیستم عامل و نیاز به فرمت، بارش زیبای

شهاب سنگ از آسمان و ...)

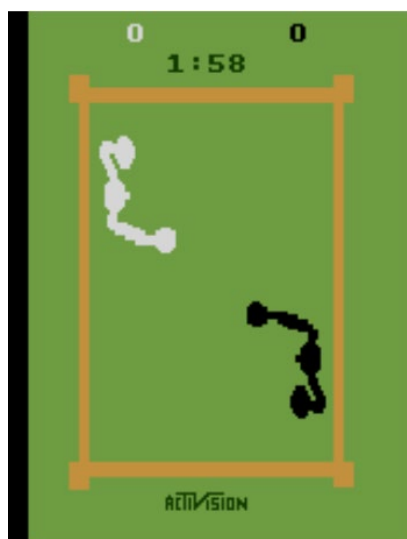
---

## شرح تمرین :

این تمرین شامل دو بخش عملی خواهد بود.

استفاده از کتابخانه های آماده مبارز نمی باشد، حتی شما دانشجوی عزیز!

بخش اول) یادگیری در محیط چندعامله



شکل (۱-۱) تصویری از محیط بوکس

بوکس یک بازی رقابتی است که در آن کنترل دقیق و واکنش های مناسب به حریف شما بسیار مهم است. بازیکنان دو دقیقه (حدود ۱۲۰۰ حرکت) برای مبارزه در رینگ زمان دارند. در هر حالت محیط، آن ها می توانند حرکت کنند و مشت بزنند. مشت های موفق امتیاز کسب می کنند، ۱ امتیاز برای یک ضربه چپ (long range jab)، ۲ امتیاز برای یک ضربه قدرتی (close power punch) نزدیک و ۱۰۰ امتیاز برای knock out (که همچنین بازی را به پایان می رساند). زمانی که شما تعدادی امتیاز کسب می کنید، متناسب با تعداد امتیاز کسب شده پاداش دریافت خواهید کرد و حریف شما نیز با آن تعداد امتیاز مبارزات می شود. بازی در یک محیط دو بعدی با اندازه کوچک اجرا می شود که در آن بوکسورها در حال حرکت بوده و سعی بر وارد کردن ضربات به حریف خود دارند و در عین حال می بایست از ضربات حریف جلوگیری کنند تا بتوانند امتیازات بیشتری را کسب کنند و در نهایت برنده بازی شوند.

فرض می‌کنیم که عامل‌ها محیط را نمی‌دانند. با استفاده از الگوریتم SARSA که می‌بایست توسط خودتان پیاده شود، یکی از عامل‌ها را به گونه‌ای آموزش دهید که توانایی تعامل با محیط و برد مسابقه را داشته باشد. در نهایت نمودار متوسط مجموع پاداش دریافتی (average reward) را رسم کنید. برای اطلاعات بیشتر از محیط بازی می‌توانید به این [لینک](#) مراجعه نمایید.

## بخش دوم) عامل‌های model-free و model-based

در فصل زمستان در یک دریاچه یخی مسیرهای مختلفی برای پیمایش بین نقطه شروع و پایان وجود دارد. عامل هوشمندی قرار است طی روزهای مختلف این مسیر را بییماید. نکته مهم برای این پیمایش این است که در هر منطقه از این دریاچه یک احتمال برای شکستن یخ وجود دارد (H). نقطه‌ی شروع و پایان مسیر در شکل ۱-۲ نشان داده شده است. خانه آبی (F) نشان دهنده‌ی مکان‌های بی‌خطر برای حرکت توسط agent است.

S	F	F	F
F	H	F	H
F	F	F	H
H	F	F	G

شکل (۱-۲) محیط دریاچه یخی

عامل در هر زمان می‌تواند از بین چهار حرکت چپ (۰)، پایین (۱)، راست (۲) و بالا (۳) یکی را انتخاب کند. البته توجه کنید که در حالت‌های مرزی در صورت انتخاب حرکت غیرمجاز عامل در سر جای خود باقی می‌ماند.

ماند. عامل به ازای هر حرکت به دلیل زمان از دست رفته پاداش  $-1$  دریافت می‌کند. در صورت سقوط در هر یک از خانه ها، بازی فاتمه یافته و عامل پاداش  $-10$  دریافت می‌کند و در صورت رسیدن به هدف عامل پاداش  $+50$  دریافت می‌کند. هم چنین نکته قابل توجه دیگر آن است که حرکات عامل به دلیل سر بودن سطح دریاچه به صورت قطعی نیست به این معنا که عامل ما با احتمال  $0.9$  حرکت انتفاخ شده را انجام می‌دهد و در غیر این صورت یکی از  $ε$  حرکت را به صورت تصادفی انجام می‌دهد. محیط این مسئله در پوشه تمرین پیوست گردیده است.

در این تمرین می‌فواهیم روشهای **model-based** و **model-free** و هم چنین ترکیب آنها را بررسی کنیم:

برای تمامی بخش‌های زیر، هنگام تعریف محیط، متغیر **studentNum** را در کدهای ارائه شده برابر شماره دانشجویی خودتان قرار دهید.

#### ۱. بررسی روشهای **model-based** :

ابتدا فرض می‌کنیم که عامل مدل محیط را می‌داند. با کمک الگوریتم **value iteration** و با در نظر گرفتن مقدار **discount factor = 0.9** ارزش هر خانه را پیدا کنید و به کمک آن **Q-value** های هر حالت-عمل ( **state-action** ) را محاسبه کنید. هم چنین نهایتاً سیاست بهینه را بیابید. سپس بر روی نقشه برای هر خانه عدد عمل بهینه را نمایش دهید.

#### ۲. بررسی روشهای **model-free** :

در این بخش می‌فواهیم عامل **model-free** را بررسی کنیم سیاست مورد استفاده برای عامل را **epsilon-greedy** در نظر بگیرید. مقدار اپسیلون را به صورت کاهش و مقدار **discount factor** را  $0.9$  و همچنین مقدار نرخ یادگیری را برابر  $0.1$  در نظر بگیرید. برای تمامی روش‌های زیر مسئله را

برای تعداد اپیزود خواسته شده برای ۲۰ بار تکرار انجام دهید و متوسط پاداش دریافتی در هر اپیزود را در طول یادگیری را رسم نمایید.

الگوریتم  $q$ -learning و SARSA را به ازای نرخ یادگیری ۰.۰۵ و کاهش پیاده‌سازی نمایید و نتایج درست آمده را با استفاده از نمودار میانگین پاداش در افق ۲۰۰۰ اپیزود نمایش دهید.

آنبه تحویل داده می‌شود:

(۱) کد اجرایی برنامه با توضیحات لازم برای اجرا

(۲) گزارش کاملی از مسیر انجام کار، چالش‌هایی که مواجه شده‌اید، اجراهایی که گرفتید و نتایجی که

---

حاصل شده است. گزارش کار از اهمیت بالایی برخوردار است، مهم آن و فرمت استاندارد

---

آن اهمیت ندارد، اما باید نشان دهنده مسیر انجام پروژه، چالشها، راه‌ها و نتایج کار شما

---

باشد.