

تمرین سری چهارم: فرآیندهای تصمیم‌گیری مارکوف

لطفاً به نکات زیر توجه کنید:

- مهلت ارسال این تمرین برای هر دو گروه ۴ آذر ماه است.
- در صورتی که به اطلاعات بیشتری نیاز دارید می‌توانید به صفحه‌ی تمرین در وب‌سایت درس مراجعه کنید.
- ما همواره هم‌فکری و هم‌کاری را برای حل تمرین‌ها به دانشجویان توصیه می‌کنیم. اما هر فرد باید تمامی سوالات را به تنهایی تمام کند و پاسخ ارسالی حتماً باید توسط خود دانش‌جو نوشته‌شده باشد. لطفاً اگر با کسی هم‌فکری کردید نام او را ذکر کنید.
- لطفاً برای ارسال پاسخ‌های خود از راهنمای موجود در صفحه‌ی تمرین استفاده کنید.
- هر سؤالی درباره‌ی این تمرین را می‌توانید در گروه درس مطرح کنید و یا از دستیاران حل تمرین پرسید.

- آدرس صفحه‌ی تمرین: https://iust-courses.github.io/ai97/assignments/04_mdp

- آدرس گروه درس: <https://groups.google.com/forum/#!forum/ai97>

سؤالها

سوال یک (۲۵ نمره)

هدف یک عامل شبیه‌ساز فوتبال ۲ بعدی در زمین، کمک کردن به تیم برای رسیدن به پیروزی در بازی است. بیشترین سود برای تیم رسیدن توپ به خط حمله هست اما در این مکان احتمال از دست دادن توپ وجود دارد و برای برگرداندن توپ به خط حمله نیاز به صرف انرژی می‌باشد (در هر مرحله از دست دادن توپ، توپ به یک منطقه عقب‌تر در صورت وجود برمی‌گردد).

این عامل می‌تواند در یکی از ۳ حالت ۱- خط دفاع ۲- خط هافبک ۳- خط حمله باشد. در هر یکی از این حالت‌ها ۲ اکشن ۱- حرکت کردن ۲- حرکت نکردن را می‌تواند انجام بدهد. حرکت کردن معادل یک واحد صرف انرژی است و یک امتیاز منفی برای عامل در پی دارد. وقتی عامل در خط حمله قرار می‌گیرد ۳ امتیاز مثبت دریافت می‌کند و در حالت‌های دیگر ۱ امتیاز. برای مثال اگر عامل در خط حمله باشد و در حال حرکت (برای حفظ توپ) پاداش $2=1-3$ را دریافت می‌کند.

احتمال موقعیت بعدی عامل بعد هر عمل در هر منطقه از زمین معادل جدول زیر است:

حرکت نکردن	حرکت کردن	
دفاع ۱	دفاع ۰.۴ / هافبک ۰.۶	دفاع
دفاع ۱	دفاع ۰.۱ / هافبک ۰.۶ / حمله ۰.۳	هافبک
حمله ۰.۷ / هافبک ۰.۳	حمله ۰.۱ / هافبک ۰.۹	حمله

الف) MDP این مسئله را رسم کنید.

ب) با استفاده از ضریب تخفیف^۱ ۰.۸ و روش «تکرار ارزش^۲»، این MDP را حل کنید. با مجموعه مقادیر صفر باید این کار را آغاز کنید، و راهبرد^۳ و مقدار بهینه را نشان دهید.

پ) با استفاده از ضریب تخفیف ۰.۸ و روش «تکرار راهبرد^۴» این MDP را حل کنید. حل را با راهبرد اولیه و غیر بهینه‌ی «حرکت نکردن» آغاز کنید. راهبرد و مقدار بهینه را نشان دهید.

1 Discount Factor

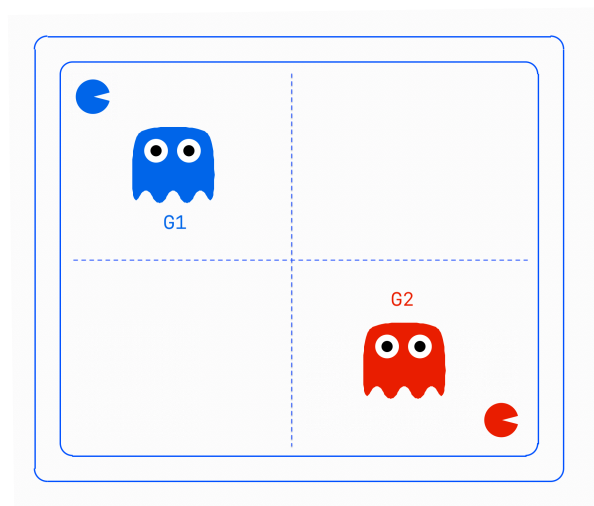
2 Value Iteration

3 Policy

4 Policy Iteration

سوال دو (۲۵ نمره)

در این سوال می‌خواهیم از MDP در شکلی از بازی پک‌من که یک بازی دو نفره‌ی نوبتی با مجموع صفر است استفاده کنیم. فرض کنید در یک صفحه‌ی ۲ در ۲، دو روح $G1$ و $G2$ به شکل زیر قرار دارند. وظیفه‌ی دو روح خوردن پک‌من طرف مقابل است. مکان پک‌من‌ها در نقشه مشخص شده‌اند و با رسیدن هر روح به پک‌من مورد نظرش آن را می‌خورد و بازی تمام می‌شود. هم‌چنین دو روح هم نمی‌توانند هم‌زمان در یک خانه از جدول باشند. هر روح فقط می‌تواند در جهت عمودی و یا افقی حرکت کند. $R(s)$ تابعی است که پاداش بودن در هر حالت را از دید روح اول ($G1$) مشخص می‌کند به این صورت که اگر روح اول زودتر پک‌من طرف مقابل را بخورد مقدار (+۱) و اگر روح دوم زودتر پک‌من طرف مقابل را بخورد مقدار (-۱) و در غیر این صورت صفر برگردانده می‌شود. شروع بازی با روح اول است.



الف) در حالت کلی اگر $V_{G1}(s)$ ارزش حالت s وقتی نوبت حرکت $G1$ باشد و هم‌چنین $V_{G2}(s)$ ارزش حالت s وقتی نوبت حرکت $G2$ باشد، معادلات بل من^۵ را برای V_{G1} و V_{G2} بنویسید (دقت کنید که تمامی محاسبات امتیاز و ارزش‌ها بر اساس دید روح اول از بازیست).

ب) توضیح دهید چگونه با استفاده از این دو معادله‌ای که به دست آوردید می‌توان الگوریتم «تکرار ارزش» را به صورت دو نفره انجام داد. ضمناً شرط پایان مناسب را هم تعیین کنید.

پ) نقشه‌ی حالات MDP را برای این بازی رسم کنید و هم‌چنین امتیاز هر حالت را برای آن بنویسید.

ت) (امتیازی / ۱۰ نمره) حال الگوریتم «تکرار ارزش»^۵ که در قسمت (ب) به دست آوردید را برای این بازی اعمال کنید و راهبرد بهینه را به دست آورید.

⁵ Bellman Equations

سوال سه (۳۰ نمره)

تا به این جا از الگوریتم «تکرار ارزش» برای حل MDPها استفاده کردیم و می دانیم با استفاده از آن حتماً به جواب بهینه می رسیم. اما در این سوال می خواهیم دلیل هم گرایی این الگوریتم را بیاییم. در ریاضی مفهومی وجود دارد به نام نگاشت انقباضی^۶. اگر بخواهیم حدودی آن را تعریف کنیم، این نگاشت یک تابع با یک ورودیست که اگر آن را به طور مجزا روی دو عدد اعمال کنیم، این تابع خروج‌هایی برای این دو تولید می کند که به هم نزدیکترند. به طور مثال «تقسیم بر دو» یک نگاشت انقباضیست. توجه کنید که این نوع عملگر یک نقطه‌ی ثابت هم دارد که در اثر اعمال تابع، بدون تغییر می ماند (مثلاً عدد صفر در نگاشت «تقسیم بر دو» یک نقطه‌ی ثابت است).

این نگاشت دو خاصیت مهم دارد: ۱- تنها یک نقطه‌ی ثابت دارد و ۲- اگر به صورت حدی آن را بر روی یک ورودی مرتباً اعمال کنیم به به نقطه‌ی ثابت آن می رسیم (چرا؟) بنابراین اگر ثابت کنیم رویه‌ی آپدیت بل من^۷ از این نوع است هم گرایی و رسیدن به جواب یکتا برای آن اثبات می شود.

اگر V_i یک بردار باشد که هر خانه‌ی آن بیان گر یک ارزش در مرحله‌ی i ام الگوریتم تکرار ارزش باشد، و همچنین اگر رویه‌ی بل من را به صورت یک عملگر ببینیم که به صورت همزمان روی تمام خانه‌ی این بردار اعمال می شود، رابطه‌ی آپدیت بل من را می توان به صورت زیر نوشت:

$$V_{i+1} \leftarrow B[V_i]$$

به علاوه برای اثبات، به روشی برای اندازه گیری «فاصله» بین دو بردار ارزش نیاز داریم، برای این منظور از MaxNorm استفاده می کنیم. اگر این عملگر روی یک بردار اعمال شود بزرگترین قدر مطلق از مقادیر آن را برمی گرداند:

$$\|V\| = \max_s |V_s|$$

با این تعریف، بیشترین تفاوت بین هر دو بردار با $\|V - V'\|$ خواهد شد، در ادامه اگر با فرض اینکه V_i و V_i' دو بردار ارزش باشد برای رسیدن به انقباضی بودن رابطه‌ی بل من شما باید رابطه زیر را اثبات کنید:

$$\|B[V_i] - B[V_i']\| \leq \gamma \|V_i - V_i'\|$$

راهنمایی: ابتدا اثبات کنید به ازای هر دو تابع g و f رابطه‌ی زیر برقرار است، سپس رابطه بل من را در نامساوی بالا قرار دهید.

$$|\max_a f(a) - \max_a g(a)| \leq \max_a |f(a) - g(a)|$$

⁶ Contraction Mapping

⁷ Bellman Update

سوال چهار (۲۰ نمره)

در یک فرآیند تصمیم‌گیری مارکوف با تابع پاداش $R(s,a)$ می‌توان راهبرد بهینه $\pi^*(s)$ را یافت. برای چه مقادیری از c ، اعمال تغییر در تابع پاداش، همان راهبرد بهینه قبلی را به ما می‌دهد؟ هر قسمت را اثبات و یا با یک مثال نقض رد کنید.

$$R'(s,a) = R(s,a) \times c \text{ (الف)}$$

$$R'(s,a) = R(s,a) + c \text{ (ب)}$$