

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 نوشته شده 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)

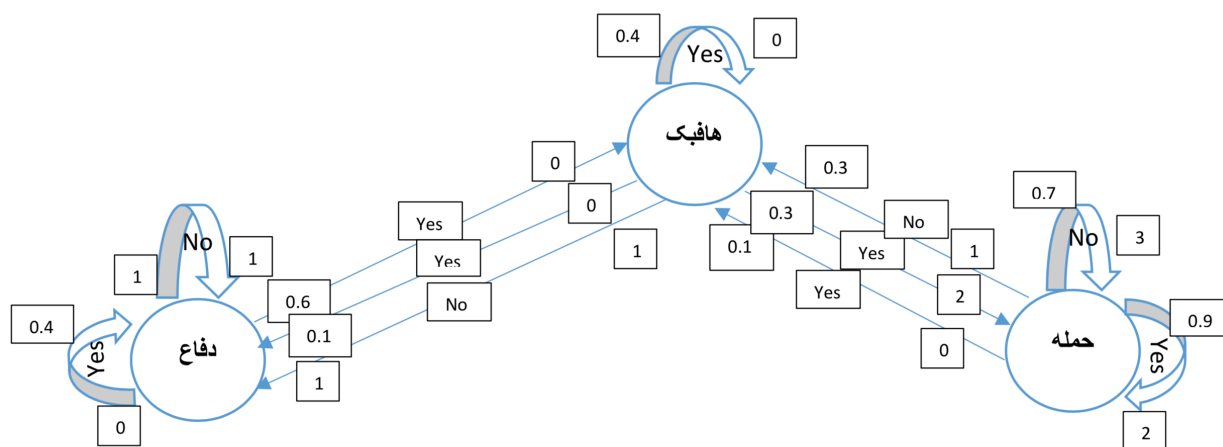
هدف یک عامل شبیه‌ساز فوتبال ۲ بعدی در زمین، کمک کردن به تیم برای رسیدن به پیروزی در بازی است. بیشترین سود برای تیم رسیدن توپ به خط حمله هست اما در این مکان احتمال از دست دادن توپ وجود دارد و برای برگرداندن توپ به خط حمله نیاز به صرف انرژی می‌باشد (در هر مرحله از دست دادن توپ، توپ به یک منطقه عقب‌تر در 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
هافبک ۰.۳ / حمله ۰.۷	هافبک ۰.۱ / حمله 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.

برای V_1 داریم:

دفاع:

$$\text{No: } 1(0+1) = 1$$

$$\text{Yes: } 0.4(0)+0.6(0)=0 \quad \text{max}=1$$

هافیک:

$$\text{No: } 1(1+0)=1$$

$$\text{Yes: } 0.1(0)+0.6(0+0.8*0)+0.3(2+0)=0.6 \quad \text{max}=1$$

حمله:

$$\text{No: } 0.7(3+0)+0.3(1+0)=2.4$$

$$\text{Yes: } 0.9(2+0)+0.1(0+0)=1.8 \quad \text{max}=2.4$$

برای V_2 داریم:

دفاع:

$$\text{Yes: } 0.4(0+0.8)+0.6(0+0.8)=0.8$$

$$\text{No: } 1(1+0.8*1)=1.8 \quad \text{max}=1.8$$

هافیک:

$$\text{Yes: } 0.1(0+0.8)+0.6(0+0.8)+0.2(2+0.8*2.4)=1.73$$

$$\text{No: } 1(1.8)=1.8 \quad \text{max}=1.8$$

حمله:

$$\text{Yes: } 0.9(2+0.8*2.4)+0.1(0+0.8)=3.528$$

$$\text{No: } 0.7(3+0.8*2.4)+0.3(1+0.8*1)=3.98 \quad \text{max}=3.98$$

¹ Discount Factor

² Value Iteration

³ Policy

به همین ترتیب مقادیر را برای V های بیشتر ادامه داده تا به عددی همگرا شویم.

همچنین در این سوال می توانستید از کد استفاده کنید.

Policy=(دفاع=Not Move, هافبک=Move, حمله=Not Move)

	دفاع	هافبک	حمله
V0	0	0	0
V1	1	1	2.4
V2	1.8	1.8	3.98
...
V25	4.98	5.91	8.67
...
V1000	~5	~6	~8.7

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.

دفاع:

$$V^*_1(\text{دفاع}) = 1(1+0) = 1 \quad \text{policy}(\text{دفاع}) = \arg \max(\text{not move}(1+0.8), \text{move}(0.6+0.8+0.6+(0.8(0))))$$

هافبک

$$V^*_1(\text{هافبک}) = 1^*(1+0) = 1 \quad \text{policy}(\text{هافبک}) = \arg \max(\text{not move}(1.8), \text{move}(0.6*0.8+0.3(2+0.8*2.4)))$$

حمله:

$$V^*_1(\text{حمله}) = 0.7*0.3+0.3(1)=2.4 \quad \text{policy}(\text{حمله}) = \arg \max(\text{not move}(0.3+(1+0.8)+0.7(3+0.8(2.4))), \text{move}(0.1(0.8+1)+0.9(2+0.8(2.4))))$$

$$\text{policy}(1) = (\text{دفاع}=\text{Not move}, \text{هافبک}=\text{Not move}, \text{حمله}=\text{Not move})$$

$$V^*_2(\text{دفاع}) = 1.8 \quad \text{policy}(\text{دفاع}) = \arg \max(\text{not move}(2.44), \text{move}(1.44))$$

$$V^*_2(\text{هافبک}) = 1.8 \quad \text{policy}(\text{هافبک}) = \arg \max(\text{not move}(2.44), \text{move}(2.564))$$

⁴ Policy Iteration

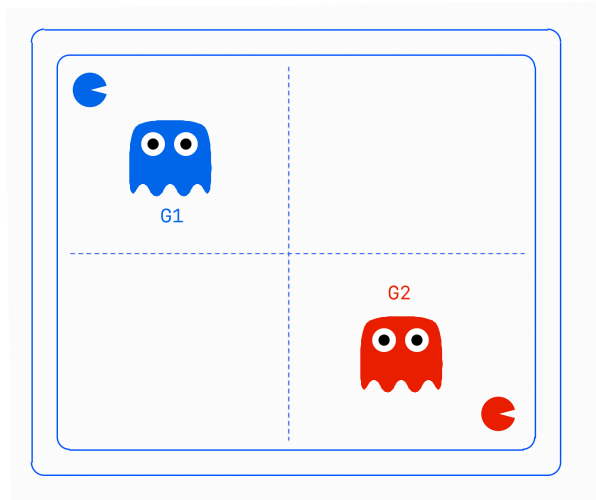
$$V_2^*(\text{حمله}) = 3.984 \quad \text{policy}(\text{حمله}) = \arg \max(\text{not move}(5.063), \text{move}(4.812))$$

$$\text{policy}(2) = (\text{دفاع} = \text{Not move}, \text{هافبک} = \text{move}, \text{حمله} = \text{Not move})$$

با محاسبه‌ی $\text{policy}(3)$ نیز به همین مقدار $\text{policy}(2)$ خواهیم رسید و مقادیر تقریباً همگرا می‌شوند و policy^* همان $\text{policy}(2)$ است

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 $(+1)$ و اگر روح دوم زودتر پک‌من طرف مقابل را بخورد مقدار (-1) 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 بل‌من^۵ را برای V_{G1} و V_{G2} بنویسید (دقت کنید که تمامی محاسبات امتیاز و ارزش‌ها بر اساس دید اول از بازیست).

⁵ Bellman Equations

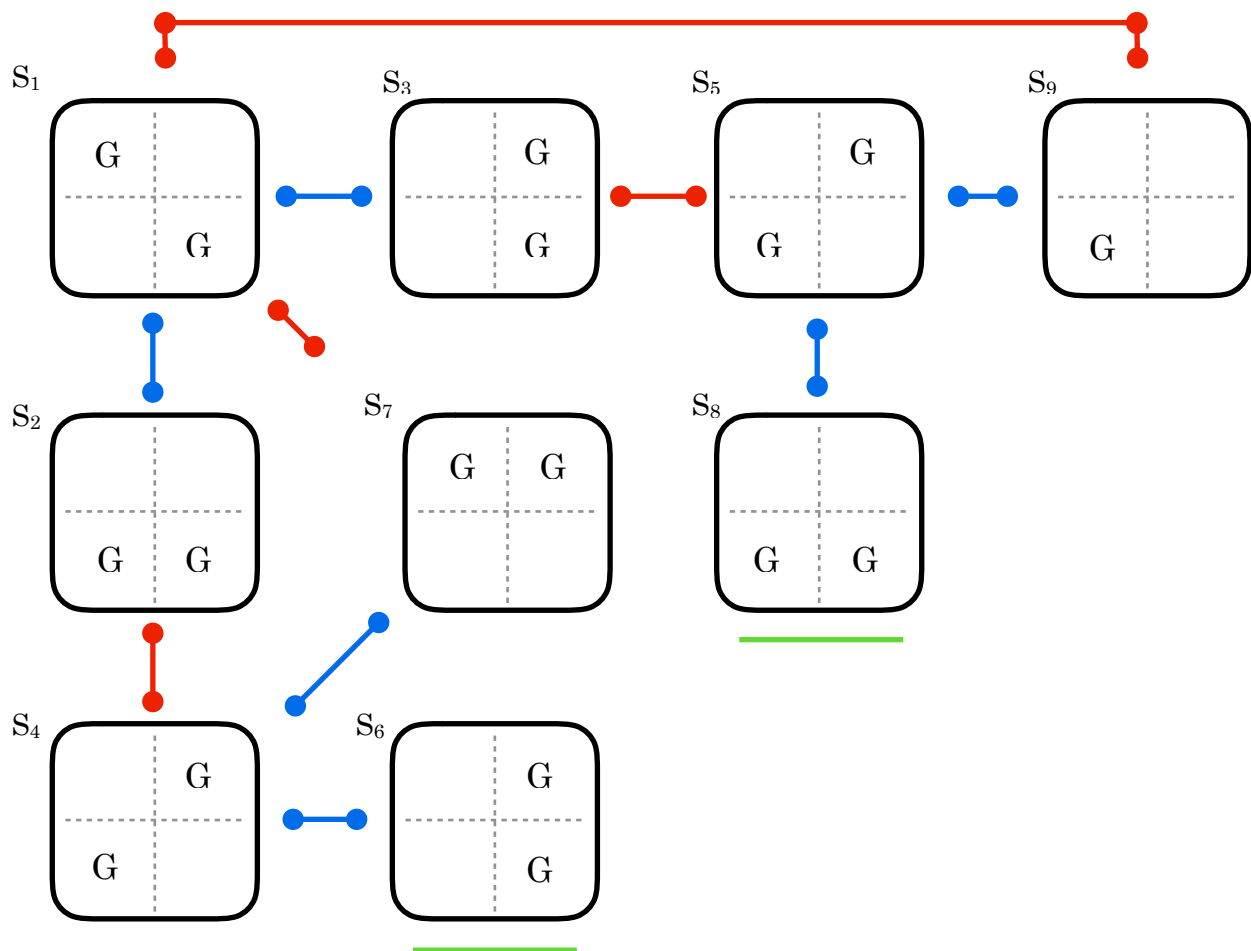
$$V_{G_1}(s) = \max_a \sum T(s, a, s') [R(s, a, s') + \lambda V_{G_2}(s')]$$

$$V_{G_2}(s) = \min_a \sum T(s, a, s') [R(s, a, s') + \lambda V_{G_1}(s')]$$

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.

برای انجام الگوریتم تکرار ارزش، هر کدام از معادلات قسمت بالا را به فرمول آپدیت بل من تبدیل می‌کنیم، سپس آن‌ها را به پشت سر هم (فرمول ۱، فرمول ۲، فرمول ۱ و ...) رو تمامی استیت‌ها اعمال می‌کنیم. شرط پایان آن است که ارزش‌های یک روح، با ارزش‌های همان روح در مرحله‌ی قبل برابر باشد

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»ی که در قسمت (ب) به دست آوردید را برای این بازی اعمال کنید و راهبرد بهینه را به دست آورید.

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 هم چنین 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 المان دو بردار برابر 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.

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

$$\begin{aligned} |\max_a f(a) - \max_a g(a)| &= \max_a f(a) - \max_a g(a) && \text{(by assumption)} \\ &= f(a^*) - \max_a g(a) \\ &\leq f(a^*) - g(a^*) \\ &\leq \max_a |f(a) - g(a)| && \text{(by definition of max)} \end{aligned}$$

$$\begin{aligned}
|(BU_i - BU'_i)(s)| &= |R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s' | s, a) U_i(s') \\
&\quad - R(s) - \gamma \max_{a \in A(s)} \sum_{s'} P(s' | s, a) U'_i(s')| \\
&= \gamma \left| \max_{a \in A(s)} \sum_{s'} P(s' | s, a) U_i(s') - \max_{a \in A(s)} \sum_{s'} P(s' | s, a) U'_i(s') \right| \\
&\leq \gamma \max_{a \in A(s)} \left| \sum_{s'} P(s' | s, a) U_i(s') - \sum_{s'} P(s' | s, a) U'_i(s') \right| \\
&= \gamma \left| \sum_{s'} P(s' | s, a^*(s)) U_i(s') - \sum_{s'} P(s' | s, a^*(s)) U'_i(s') \right| \\
&= \gamma \left| \sum_{s'} P(s' | s, a^*(s)) (U_i(s') - U'_i(s')) \right|
\end{aligned}$$

بعد از اعمال MaxNorm

$$\begin{aligned}
\|BU_i - BU'_i\| &= \max_s |(BU_i - BU'_i)(s)| \\
&\leq \gamma \max_s \left| \sum_{s'} P(s' | s, a^*(s)) (U_i(s') - U'_i(s')) \right| \\
&\leq \gamma \max_s |U_i(s) - U'_i(s)| = \gamma \|U_i - U'_i\|
\end{aligned}$$

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.

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

با ضرب مقدار c در $Q, R(s, a)$ به شکل:

$$Q^*(s) = cR(s, a) + \gamma V^*(s)$$

خواهد شد.

اگر داشته باشیم:

$$C=5$$

$$\begin{aligned} \text{Action 1: } & T=0.5 \\ & V=2 \\ & R=2 \end{aligned}$$

$$\begin{aligned} \text{Action 2: } & T=0.5 \\ & V=-2 \\ & R=5 \end{aligned}$$



$$2+2=4$$

>



$$5-2=3$$

$$C=5 \rightarrow 2(5)+2 = 12$$

<

$$5(5)-2=23$$

با مثال نقض این تساوی برقرار نیست.

$$R'(s,a) = R(s,a) + c \quad (\text{ب})$$

$$R_1(s,a) + \gamma_1 V_1 \quad >$$

$$R_2(s,a) + \gamma_2 V_2$$

$$R_1(s,a) + \gamma_1 V_1 + C \quad >$$

$$R_2(s,a) + \gamma_2 V_2 + C$$

$$R_1'(s,a) + \gamma_1 V_1 \quad >$$

$$R_2'(s,a) + \gamma_2 V_2$$

پس این تساوی برقرار است.