



دانشکده مهندسی کامپیوتر

هوش مصنوعی و سیستم‌های خبره

تمرین تشریحی چهارم

نام و نام خانوادگی غزل زمانی نژاد

شماره دانشجویی 97522166

مدرس محمدطاهر پیلهور - سید صالح اعتمادی

طراحی و تدوین سپهر باباپور (@Spr_Bpr)

تاریخ انتشار ۵۱ آبان ۹۹۳۱

تاریخ تحویل ۹۲ آبان ۹۹۳۱

فهرست مطالب

۱ سوالات بخش تئوري ۲

۱.۱ سوال ۱ ۲

۲.۱ سوال ۲ ۲

۳.۱ سوال ۳ ۳

۲ مسائل محاسباتي ۳

۱.۲ سوال ۱ ۳

۲.۲ سوال ۲ ۵

۳.۲ سوال ۳ ۷

۱ سوالات بخش تئوری

** در این بخش به سوالاتی که دارای * هستند پاسخ دهید **

۱.۱ سوال ۱

توضیح دهید چرا در MDP هانمی‌توان از روش planning استفاده کرد. راه حل جایگزین را توضیح دهید.

پاسخ:

.....

.....

.....

.....

.....

.....

۲.۱ سوال ۲

فرایندهای مارکوفی را تعریف کرده و بگویید کدام یک از فرایندهای زیر مارکوفی هستند.

-بازی بی‌سوالی

- بازی اسم - فامیل

- بازی سودوکو

پاسخ:

.....

.....

.....

.....

.....

۳.۱ * سوال ۳ (۰.۲ نمره)

اگر به جای ضریب تخفیف γ از توابع زیر استفاده شود:

$$e^{-t} * \quad \log(t) * \quad | \sin(t) | *$$

به سوالات زیر پاسخ دهید.

- ۱- کدام یک مشکل نامحدود شدن بازی را برطرف می‌کنند؟ توضیح دهید.
- ۲- برای تابع پاسخ قسمت اول، با فرض پاداش یک واحد در هر لحظه کوچک زمانی (dt) پاداش کل را محاسبه کنید.

پ

(۱) استفاده از تابع e^{-t} این مشکل را رفع می‌کند. چون:

(الف) $0 < |e^{-t}| \leq 1$ طبق ترین e این تابع بین صفر و یک قرار می‌گیرد.

(ب) $\lim_{t \rightarrow \infty} |e^{-t}| < \infty$ این تابع در نهایت به ۰ میل می‌کند.

(۲)

$$U_{total} = R(s_0, a_0, s_1) + \gamma R(s_1, a_1, s_2) + \gamma^2 R(s_2, a_2, s_3) + \dots$$

.....

$$\frac{U_{total}}{R_i = 1} = \sum_{t=0}^{\infty} 1 + e^{-1} + e^{-2} + \dots = \frac{1}{1-e^{-1}} = \frac{e}{e-1} \quad (t \text{ discrete})$$

.....

$$U_{total} = \int_0^{\infty} e^{-t} \times \frac{R}{1} dt = -e^{-t} \Big|_0^{\infty} = 1 - 0 = 1 \quad (t \text{ continuous})$$

.....

۲ مسائل محاسباتی

** در این بخش به سؤالاتی که دارای * هستند پاسخ دهید **

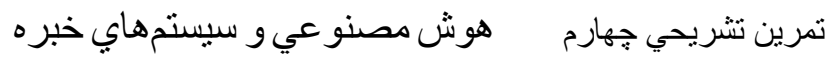
۱.۲ سوال ۱: کارت بردار!!

فرض کنید شما در یک مسابقه کارت بازی شرکت کرده‌اید که در آن ۳ نوع کارت با شماره‌های ۲، ۳، ۴ وجود دارد. شما در هر مرحله از بازی تا زمانی که به مجموع امتیاز ۶ نرسیده‌اید می‌توانید یا یک کارت بردارید یا بازی را به اتمام برسانید. احتمال آمدن هر کارت با هم برابر است. زمانی که مجموع امتیازات شما ۶ یا بیشتر شود امتیازات شما صفر می‌شود و بازی تمام می‌شود و زمانی که خودتان بازی را تمام کرده باشید امتیازتان برابر مجموع کارت‌هایی که کسب کرده‌اید می‌شود. همچنین برداشتن کارت را بدون هزینه در نظر بگیرید.

در این سوال از شما خواسته شده است که بازی فوق را به صورت یک مدل مارکوفی در نظر بگیرید و به سوالات زیر پاسخ دهید.

۱. ابتدا تابع انتقال (transition function) و تابع پاداش (reward function) را برای این مدل محاسبه کنید.

۲. سپس جدول زیر را کامل کنید.



شکل ۱: جدول سوال کارت بردار!!

[illegible]

۲.۲ * سوال ۲: تاس بریز!! (۵۴)

نمره) فرض کنید در يك بازی ریختن تاس شرکت کرده‌اید که هزینه هر بار ریختن تاس در آن ۱ سکه است و احتمال آمدن

تمام اعداد در تاس با یکدیگر برابر است. شما پس از ریختن تاس به اندازه عدد روی تاس سکه دریافت می‌کنید. قانون بازی به این شکل است که شما موظف هستید در بار اول يك تاس بریزید، اما در سایر مراحل دو انتخاب دارید:

* اتمام بازی: با این حرکت شما به اندازه عدد روی تاس سکه دریافت می‌کنید.

* تاس ریختن: يك سکه هزینه می‌کنید و بار دیگر تاس می‌ریزید.

لذا بازی را می‌توان به این صورت در نظر گرفت که بازیکن در ابتدای بازی در حالت شروع قرار دارد و در حالت شروع فقط حرکت ریختن تاس وجود دارد. در سایر حالات يك حرکت اتمام بازی وجود دارد که بازیکن را به حالت پایانی می‌برد و در حالت پایانی حرکتی وجود ندارد. هر حالت بین شروع و پایان با s_i نمایش داده می‌شود که بدین معنی است که عدد i در تاس آمده است. باتوجه به توضیحات فوق به سوالات زیر پاسخ دهید:

۱. فرض کنید π های زیر در ابتدا وجود دارد، ردیف π را کامل کنید. ($\gamma = 1$)

حالت	S_1	S_2	S_3	S_4	S_5	S_6
\square_1	تاس ریخ'	تاس ریخ'	اتمام بازی	اتمام بازی	اتمام بازی	اتمام بازی
\square''_1	3	3	3	4	5	6

شکل ۲: جدول قسمت اول سوال تاس بریز!!

۲. باتوجه به جدول فوق مقادیر π را برورسانی کنید و در جدول زیر جایگذاری کنید. این مقادیر می‌تواند سه حالت

تاس ریختن، اتمام بازی و تاس ریختن / اتمام بازی باشد. ($\gamma = 1$)

حالت	S_1	S_2	S_3	S_4	S_5	S_6
\square_1	تاس ریخ'	تاس ریخ'	اتمام بازی	اتمام بازی	اتمام بازی	اتمام بازی
\square_1, \square_2	تاس ریختن	تاس ریختن	تاس ریختن / اتمام بازی	اتمام بازی	اتمام بازی	اتمام بازی

شکل ۳: جدول قسمت دوم سوال تاس بریز!!

۳. باتوجه به مقادیر جدول فوق آیا می‌توان نتیجه گرفت که مقادیر بدست آمده بهینه هستند و دیگر نیاز به برورسانی ندارند؟ توضیح دهید.

پاسخ:

حالت	s_1	s_2	s_3	s_4	s_5	s_6
π_i	تاس ریختن	تاس ریختن	اتمام بازی	اتمام بازی	اتمام بازی	اتمام بازی
v^{π_i}	3	3	3	4	5	6

شکل ۲: جدول قسمت اول سوال تاس بریز!

۲. باتوجه به جدول فوق مقادیر π_i را بروزرسانی کنید و در جدول زیر جایگذاری کنید. این مقادیر می‌تواند سه حالت تاس ریختن، اتمام بازی و تاس ریختن / اتمام بازی باشد. ($\gamma = 1$)

حالت	s_1	s_2	s_3	s_4	s_5	s_6
π_i	تاس ریختن	تاس ریختن	اتمام بازی	اتمام بازی	اتمام بازی	اتمام بازی
π_{i+1}	تاس ریختن	تاس ریختن	تاس ریختن / اتمام بازی	اتمام بازی	اتمام بازی	اتمام بازی

$$1) \quad v_{\pi_3} = 3 \quad v_{\pi_4} = 4 \quad v_{\pi_5} = 5 \quad v_{\pi_6} = 6$$

$$v_{\pi_1} = \frac{1}{4}(-1 + v_{\pi_1}) + \frac{1}{4}(-1 + v_{\pi_2}) + \frac{1}{4}(-1 + v_{\pi_3}) + \frac{1}{4}(-1 + v_{\pi_4}) + \frac{1}{4}(-1 + v_{\pi_5}) + \frac{1}{4}(-1 + v_{\pi_6})$$

$$v_{\pi_1} = -1 + \frac{1}{4}(v_{\pi_1} + v_{\pi_2} + 18)$$

$$v_{\pi_2} = -1 + \frac{1}{4}(v_{\pi_1} + v_{\pi_2} + 18)$$

$$\Rightarrow \begin{cases} \frac{5}{6}v_{\pi_1} - \frac{1}{6}v_{\pi_2} = 2 \\ \frac{5}{6}v_{\pi_2} - \frac{1}{6}v_{\pi_1} = 2 \end{cases} \Rightarrow v_{\pi_1} = v_{\pi_2} = 3$$

$$2) \quad \pi_{i+1}(s_1) = \operatorname{argmax}\{\text{dice}: 3, \text{stop}: 1\} = \text{dice}$$

$$\pi_{i+1}(s_2) = \operatorname{argmax}\{\text{dice}: 3, \text{stop}: 2\} = \text{dice}$$

$$\pi_{i+1}(s_3) = \operatorname{argmax}\{\text{dice}: 3, \text{stop}: 3\} = \text{dice, stop}$$

$$\pi_{i+1}(s_4) = \operatorname{argmax}\{\text{dice}: 3, \text{stop}: 4\} = \text{stop}$$

$$\pi_{i+1}(s_5) = \operatorname{argmax}\{\text{dice}: 3, \text{stop}: 5\} = \text{stop}$$

$$\pi_{i+1}(s_6) = \operatorname{argmax}\{\text{dice}: 3, \text{stop}: 6\} = \text{stop}$$

۳) بله؛ اگر نه π_3 ، stop یا انتخاب کنیم، π_{i+1} با π_i یکسان می‌شود. می‌توانیم نتیجه بگیریم converge شده پس π_i همین است.

۳.۲ * سوال ۳: يك MDP ساده كه ديگر مثل قبل نيست؟ (۵۳ نمره)

يك مسئله MDP را تصور كنيد كه در آن تابع پاداش به جاي $R(s)$ ، $\eta R(s)$ باشد كه در آن η يك ثابت مثبت است. ساير خصوصيات اين مسئله MDP تغيير نكرده است. ثابت كنيد راهبرد (policy) بهينه در مسئله MDP جديد مشابه راهبرد (policy) در مسئله اوليه است.

پاسخ:

(3) بايد ثابت كنيم اگر M ، reward ها η برابر كنيم و M_{η} بهرست بايد، optimal policy تغيير نمي‌كند

مي دانيم كه M مجموع discounted reward هاي مربوط به π^* (optimal policy)، از ساير π ها برتر است يعني:

$$\sum_{t=0}^{\infty} \gamma^t R_t(s, \pi^*, s') \geq \sum_{t=0}^{\infty} \gamma^t R_t(s, \pi, s') \quad \forall \pi \in \text{policies}$$

مي توانيم در طرف نامعادله فوق را به η ضرب كنيم (چون نامعادله به دليل مثبت بودن η تغيير نمي‌كند):

$$\eta \sum_{t=0}^{\infty} \gamma^t R_t(s, \pi^*, s') \geq \eta \sum_{t=0}^{\infty} \gamma^t R_t(s, \pi, s')$$

$$\Rightarrow \sum_{t=0}^{\infty} \gamma^t \eta R_t(s, \pi^*, s') \geq \sum_{t=0}^{\infty} \gamma^t \eta R_t(s, \pi, s')$$

η جمع discounted rewards

M_{η} مشابه با M است. optimal policy M_{η} مشابه با M است.

ثابت شد كه