



آزمون پایان ترم

- زمان در نظر گرفته شده برای آزمون ۱۸۰ دقیقه است.
- لطفاً پاسخ‌های خود را خوانا و خوش خط بنویسید.

سوالات (۱۰۰ نمره)

۱. (۱۵ نمره) به سوالات زیر به طور مختصر پاسخ دهید:

- (آ) درستی یا نادرستی عبارت رو به رو را با ذکر دلیل مشخص کنید: در یک شبکه بیزی که در آن X به شرط Z از Y مستقل است؛ ممکن است فرض استقلال این دو متغیر با شرطی کردن شواهد اضافی برای متغیرهای دیگر در شبکه، برقرار نباشد.
- (ب) دو مورد از روش‌های جلوگیری از بیش‌برازش (overfitting) در درخت تصمیم را ذکر کرده و مختصراً توضیح دهید.
- (ج) آیا همواره با استفاده از مدل درخت تصمیم بدون محدودیت عمق می‌توان به دقت ۱۰۰٪ رسید؟ توضیح دهید.
- (د) برای طبقه‌بندی تصاویر به دو کلاس سگ و گربه، از یک شبکه عصبی عمیق استفاده شده است. در انتهای این شبکه، از یک تابع فعال‌ساز $ReLU$ و سپس از تابع $sigmoid$ استفاده شده است. در صورتی که خروجی نهایی شبکه $\hat{y} \geq 0.5$ باشد آن را به عنوان گربه و در غیر این صورت به عنوان سگ در نظر می‌گیریم. آیا این شبکه ایرادی در برچسب‌گذاری تصاویر دارد؟ مختصراً توضیح دهید.
- (ه) فرض کنید از رابطه‌ی

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha \left[R(s, a, s') + \gamma \max_{a'} f(Q(s', a'), N(s', a')) \right]$$

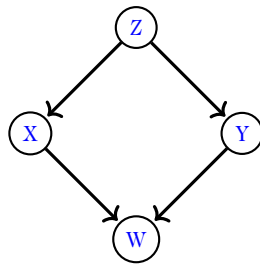
برای بروزرسانی Q -value ها استفاده کنیم که $N(s', a')$ تعداد دفعاتی است که از حالت s' کنش a' را انجام داده‌ایم و $f(x, y) = x + \frac{\beta}{y}$ که β عدد ثابت مثبتی است. توضیح دهید که این تغییر چه تاثیری بر جست‌وجوی فضای حالت خواهد گذاشت؟

حل.

(آ) درست، چون مجموعه‌ای از متغیرها به نام W می‌توانند در شبکه موجود باشند که به دلیل وجود ساختار v شکل در میان آن‌ها، گزاره $X \perp\!\!\!\perp Y | Z, W$ برقرار نباشد. (شکل ۱)

- (ب)
- جلوگیری از رشد بیش از حد درخت با استفاده از محدود کردن تعداد برگ‌ها و یا عمق درخت
 - هرس کردن درخت با استفاده از تست‌های آماری - به این صورت که در راس‌های غیر برگ با استفاده از یک تست آماری بررسی می‌کنیم که آیا افزایش دقت مدل واقعاً به علت تقسیم بر حسب معیار انتخابی بوده است یا احتمالاً به علت وجود نویز در داده است.

(ج) خیر. در صورتی می‌توان به حالت بدون خطا دست یافت که داده‌های یکسان برچسب‌های متفاوتی نداشته باشند.



شکل ۱: سوال ۱ قسمت (آ)

(د) با توجه به این که مقدار خروجی ReLU برای هر ورودی بزرگتر از صفر است، پس از اعمال تابع سیگموید روی خروجی آن، همواره عددی بزرگتر از ۰/۵ خواهیم داشت. در نتیجه برای تمامی تصاویر یک برچسب در نظر گرفته می‌شود که نامطلوب است.

(ه) تعداد دفعاتی که جفت (s, a) را دیده‌ایم در مقدار Q-value ها تاثیرگذار است و چون این تاثیر از رابطه‌ی f ناشی می‌شود که آن هم رابطه‌ی عکس با تعداد دیده شدن‌ها دارد، می‌توان نتیجه گرفت که در صورت مثبت بودن β ، با این شیوه‌ی آپدیت کردن کنش‌هایی که به حالت‌هایی که زیاد دیده نشده‌اند منجر می‌شوند برتری خواهند داشت و این کمک می‌کند تا بهتر بتوان فضای حالت را جست‌وجو کرد.

۲. (۸ نمره) مدل مارکوف زیر را در نظر بگیرید که سه استیت ۱ و ۲ و ۳ دارد. این مدل دارای Transition Matrix به صورت زیر است:

$$P = \begin{bmatrix} \frac{1}{2} & \frac{1}{4} & \frac{1}{4} \\ \frac{1}{3} & 0 & \frac{2}{3} \\ \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix}$$

(آ) برای این مدل یک نمودار حالت (state diagram) رسم کنید.

(ب) اگر بدانیم که $\frac{1}{4} = P(X_1 = 1) = P(X_1 = 2)$ مقدار $P(X_1 = 1, X_2 = 2, X_3 = 1)$ را محاسبه کنید.

حل.

(آ) نمودار حالت این مدل در شکل ۲ قابل مشاهده است.

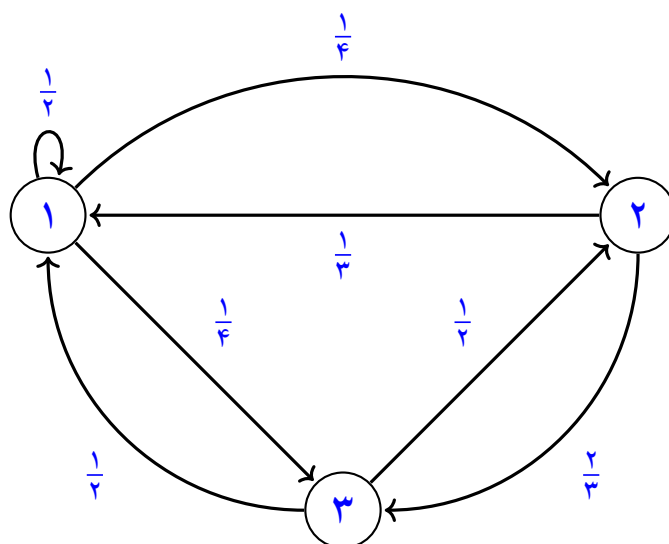
(ب) در ابتدا دقت کنیم که داریم:

$$P(X_1 = 3) = 1 - P(X_1 = 1) - P(X_1 = 2) = 1 - \frac{1}{4} - \frac{1}{4} = \frac{1}{2}$$

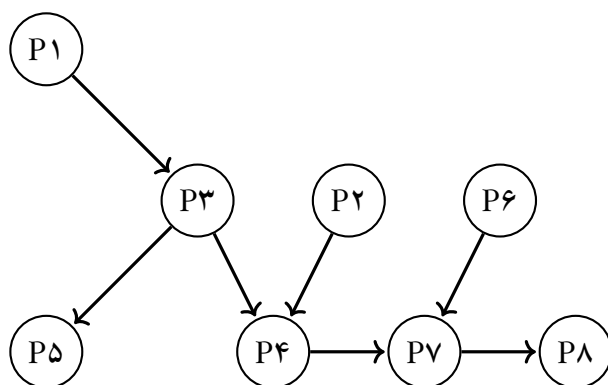
حال می‌توان نوشت:

$$P(X_1 = 3, X_2 = 2, X_3 = 1) = P(X_1 = 3) \times p_{32} \times p_{21} = \frac{1}{2} \times \frac{1}{2} \times \frac{1}{3} = \frac{1}{12}$$

۳. (۱۲ نمره) فرض کنید گراف زیر را به عنوان یک شبکه‌ی بیزین داریم:



شکل ۲: سوال ۲ قسمت (آ)



درستی یا نادرستی هر کدام از عبارتهای زیر را با ذکر دلیل مشخص کنید.

- (آ) $P_5 \perp\!\!\!\perp P_6 | P_7$
- (ب) $P_1 \perp\!\!\!\perp P_2 | P_5$
- (ج) $P_1 \perp\!\!\!\perp P_2 | P_8$
- (د) $P_5 \perp\!\!\!\perp P_8 | P_7$

حل.

- (آ) نادرست. همه‌ی سه‌تایی‌ها در مسیر P_5, P_3, P_4, P_7, P_6 فعال هستند در نتیجه P_6 و P_5 به شرط P_7 مستقل نیستند.
- (ب) درست. سه‌تایی P_3, P_4, P_2 غیرفعال است در نتیجه مسیر بین P_1 و P_2 غیرفعال بوده و این دو به شرط P_5 از هم مستقلند.
- (ج) نادرست. همه‌ی سه‌تایی‌ها در مسیر P_1, P_3, P_4, P_2 فعال هستند در نتیجه P_1 و P_2 به شرط P_8 مستقل نیستند.
- (د) درست. سه‌تایی P_4, P_7, P_8 غیرفعال است در نتیجه مسیر بین P_8 و P_5 غیرفعال بوده و این دو به شرط P_7 از هم مستقلند.

۴. (۱۴ نمره) یک مدل خطی به فرم زیر را در نظر بگیرید:

$$y(x_n, \mathbf{w}) = w_0 + \sum_{i=1}^D w_i x_{ni}$$

خطای آن را نیز به صورت زیر در نظر می‌گیریم:

$$E_D(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$

حال فرض کنید که یک نویز گوسی $\epsilon_i \sim \mathcal{N}(0, \sigma^2 I)$ به هر ورودی x_i اضافه شده است. ϵ_i ها به صورت i.i.d تولید شده‌اند. اگر $\tilde{E}_D(\mathbf{w})$ خطای مدل وقتی از $x_i + \epsilon$ استفاده می‌کنیم باشد، رابطه‌ی زیر را اثبات کنید.

$$\mathbb{E}[\tilde{E}_D(\mathbf{w})] = E_D(\mathbf{w}) + \frac{N}{N} \sum_{i=1}^D w_i^2 \sigma^2$$

حل. داریم:

$$\begin{aligned} \tilde{y}_n &= w_0 + \sum_{i=1}^D w_i (x_{ni} + \epsilon_{ni}) \\ &= y_n + \sum_{i=1}^D w_i \epsilon_{ni} \end{aligned}$$

برای $\tilde{E}_D(w)$ می‌توانیم بنویسیم:

$$\begin{aligned} \tilde{E} &= \frac{1}{N} \sum_{n=1}^N \{\tilde{y}_n - t_n\}^2 \\ &= \frac{1}{N} \sum_{n=1}^N \{\tilde{y}_n^2 - 2\tilde{y}_n t_n + t_n^2\} \\ &= \frac{1}{N} \sum_{n=1}^N \left\{ y_n^2 + 2y_n \sum_{i=1}^D w_i \epsilon_{ni} + \left(\sum_{i=1}^D w_i \epsilon_{ni} \right)^2 \right. \\ &\quad \left. - 2t_n y_n - 2t_n \sum_{i=1}^D w_i \epsilon_{ni} + t_n^2 \right\} \end{aligned}$$

اگر از رابطه‌ی بالا امید ریاضی بگیریم، جمله‌ی دوم و پنجم داخل سیگما با توجه به صفر بودن میانگین ϵ_i صفر می‌شوند.

اگر جمله‌ی سوم را باز کنیم خواهیم داشت:

$$\begin{aligned}\mathbb{E} \left[\left(\sum_{i=1}^D w_i \epsilon_{ni} \right)^2 \right] &= \mathbb{E} \left[\sum_{i=1}^D \sum_{j=1, j \neq i}^D w_i \epsilon_{ni} w_j \epsilon_{nj} + \sum_{i=1}^D w_i^2 \epsilon_{ni}^2 \right] \\ &= \sum_{i=1}^D \sum_{j=1, j \neq i}^D w_i w_j \mathbb{E}[\epsilon_{ni} \epsilon_{nj}] + \sum_{i=1}^D w_i^2 \mathbb{E}[\epsilon_{ni}^2] \\ &= 0 + \sum_{i=1}^D w_i^2 \sigma^2\end{aligned}$$

پس:

$$\begin{aligned}\mathbb{E}[\tilde{E}_D(w)] &= \frac{1}{N} \sum_{n=1}^N \left\{ y_n^2 + 0 + \sum_{i=1}^D w_i^2 \sigma^2 - 2 t_n y_n - 0 + t_n^2 \right\} \\ &= \frac{1}{N} \sum_{n=1}^N [y_n^2 - 2 t_n y_n + t_n^2] + \frac{N}{N} \sum_{i=1}^D w_i^2 \sigma^2 \\ &= E_D(w) + \sum_{i=1}^D w_i^2 \sigma^2\end{aligned}$$

۵. (۱۰ نمره) با توجه به جدول زیر قصد داریم یک درخت تصمیم بسازیم که مقدار Y را براساس X_1, X_2 و X_3 تعیین کند.

X_1	X_2	X_3	Y
No	No	No	No
Yes	No	Yes	Yes
No	Yes	Yes	Yes
Yes	Yes	Yes	No

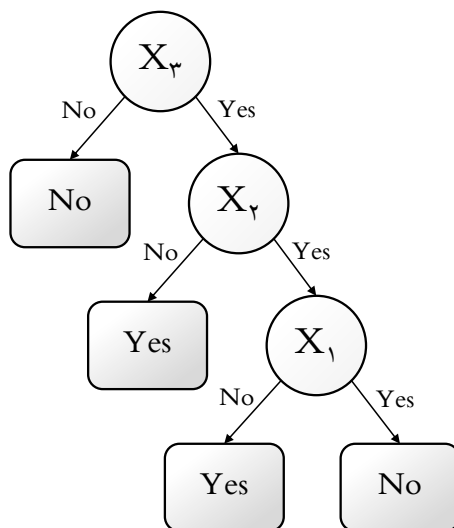
(آ) درخت تصمیم را رسم کنید. در هر مرحله ریشه را براساس Information Gain انتخاب کنید. (لزومی به محاسبه‌ی عددی IG نیست اما در صورت نیاز از آن استفاده کنید.)

(ب) آیا درخت به‌دست آمده (با انتخاب ریشه‌ها از طریق IG) بهینه است؟ اگر بهینه است علت آن را بیان کنید؛ در غیر این صورت درخت بهینه را رسم کنید. (منظور از بهینه، درختی با کوتاه‌ترین ارتفاع ممکن است که نمونه‌های سازگار را جداسازی کند.)

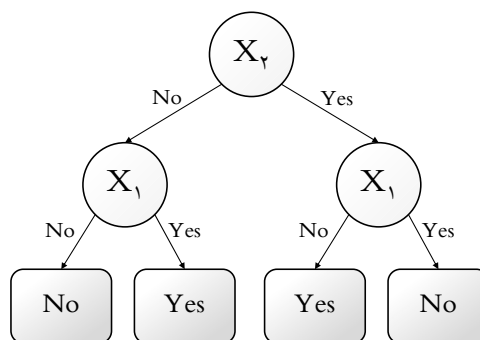
حل.

(آ) درخت در شکل ۳ رسم شده است.

(ب) با توجه به اینکه $Y = X_1 \oplus X_2$ درخت شکل ۴ ارتفاع کم‌تری دارد. به‌طور کلی لزوماً انتخاب ریشه با IG بهینه نیست؛ زیرا در این حالت ویژگی‌ها جداگانه بررسی می‌شوند و ارتباط توأم آن‌ها بررسی نمی‌شود.



شکل ۳: درخت با کمک IG



شکل ۴: درخت بهینه

۶. (۱۵ نمره) می‌دانیم در شبکه‌های عصبی از توابع مختلفی برای فعال‌سازی استفاده می‌شود. در این سوال می‌خواهیم به کمک تابع پله، شبکه‌های عصبی ساده‌ای برای پیش‌بینی مقادیر توابع منطقی بدست آوریم. توجه داشته باشید در عبارات زیر، \wedge نماد عملگر منطقی AND، \vee نماد عملگر منطقی OR و \neg نماد عملگر منطقی NOT است.

در این سوال ورودی‌ها (x_i) صفر یا یک هستند و تابع فعال‌ساز هم به صورت زیر است:

$$f(x) = \begin{cases} 0 & - (w_0 + \sum_i w_i x_i) < 0 \\ 1 & - (w_0 + \sum_i w_i x_i) \geq 0 \end{cases}$$

در این فرمول w_0 نمایانگر بایاس بوده و بقیه w_i ها وزن‌های نسبت داده شده به یال‌های درونی شبکه عصبی هستند.

(آ) یک شبکه عصبی با یک پرسپترون تشکیل دهید که بتواند تابع منطقی $y = x_1 \vee x_2$ را به درستی محاسبه کند.

(ب) یک شبکه عصبی با یک پرسپترون تشکیل دهید که بتواند تابع منطقی $y = x_1 \wedge x_2$ را به درستی محاسبه کند.

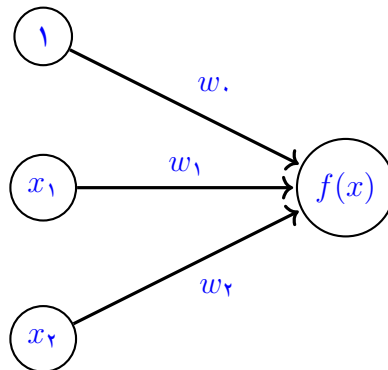
(ج) یک شبکه عصبی متشکل از یک لایه‌ی نهان تشکیل بدهید که تابع منطقی زیر را به درستی محاسبه کند.

$$y = (x_1 \wedge x_2 \wedge x_3) \vee (\neg x_1 \wedge x_3)$$

حل.

(آ) تابع OR در زمان‌هایی که حداقل یکی از ورودی‌ها ۱ است، برابر ۱ می‌شود. باید مقدار: $w_0 + w_1x_1 + w_2x_2$ را طوری تعیین کنیم که به ازای تمامی سه حالت $x_1 = 0, x_2 = 1 | x_1 = 1, x_2 = 0 | x_1 = 1, x_2 = 1$ مقدار آن بزرگتر مساوی ۰ باشد و در حالت $x_1 = 0, x_2 = 0$ کمتر ۰ باشد. جواب‌های متفاوتی می‌توان بر وزن‌ها یافت که یکی از آن‌ها اعداد زیر هستند:

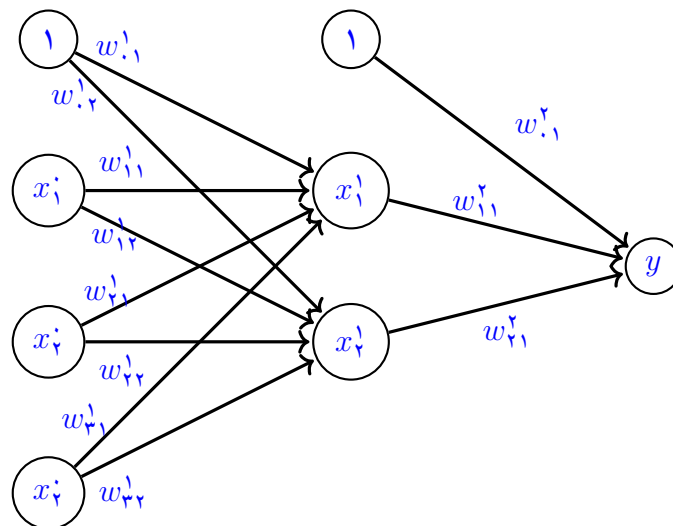
$$w_0 = 0.5, w_1 = -1, w_2 = -1$$



(ب) این قسمت مشابه قسمت قبلی حل می‌شود. یکی از جواب‌های قابل قبول این اعداد هستند:

$$w_0 = 1/5, w_1 = -1, w_2 = -1$$

(ج) در زیر ساختار شبکه و وزن‌های مربوط به آن برای ساخت عبارت خواسته شده قابل مشاهده است:



$$w^1_0 = 2/1, w^1_1 = -1, w^1_2 = -1, w^1_3 = -1$$

$$w^1_0 = 0.5, w^1_1 = 1, w^1_2 = 0, w^1_3 = -1$$

$$w^2_0 = 0.5, w^2_1 = -1, w^2_2 = -1$$

توجه کنید که به نوعی در عبارت های متشکل از AND، هر جا خود متغیر ظاهر شده در وزن نظیر آن ۱ قرار داده ایم. هر جا نقیض عبارت ظاهر شده در وزن معادل آن ۱- و هر جا هم اصلا ظاهر نشده ۰ قرار داده ایم.

* با توجه به ذکر تابع پله در صورت سوال، در صورتی که به جای تابع داده شده f که قرینه پله است از خود تابع پله استفاده شده باشد نیز مشکلی وجود ندارد. در این حالت تمامی وزن های گفته شده در این راه حل باید در منفی ضرب شود.

۷. (۱۴ نمره) می خواهیم مسئله ی پک من را با استفاده از مدل MDP حل کنیم. در مسئله ی پک من، یک موجود به نام پک من داریم که در یک جدول قرار دارد، می تواند به چهار جهت حرکت کند و هدف آن خوردن غذاهای موجود در صفحه ی بازی است. در مسئله ی مدل شده، حالت ها همان خانه های جدول و کنش ها حرکت به چهار جهت مختلف هستند. هر حرکتی که به سمت دیوار انجام شود، چه دیوارهای داخل محیط بازی و چه دیوارهای مرزی محیط بازی، بی تاثیر بوده و پک من در همان خانه ی قبلی خود باقی می ماند. هم چنین فرض کنید در کل محیط بازی فقط یک غذا وجود دارد (که با \circ در نقشه ی بازی نشان داده شده است) و پک من به محض رسیدن به آن یک امتیاز دریافت کرده و بازی تمام می شود.

• در این قسمت محیط بازی را مطابق شکل زیر در نظر بگیرید. (حروف موجود در خانه ها نام هر خانه هستند.)

A	B	C
D	E	F, \circ

ضریب تخفیف برابر 0.5 بوده و پاداش زنده بودن صفر است. با در نظر گرفتن این موارد به سوال های زیر پاسخ دهید:

(آ) سیاست بهینه را برای هر حالت مشخص کنید.

(ب) ارزش بهینه ی خانه ی A یا همان $V^*(A)$ چه قدر است؟

توجه کنید امتیاز موجود در یک خانه هنگام ورود به خانه ی دارای امتیاز محاسبه می شود. یعنی به طور مثال در شکل بالا، $R(E, \text{راست}, F) = 1$ می باشد.

• حال محیط دیگری مانند شکل زیر را برای بازی در نظر بگیرید:

A	B, \circ		
C	D, *	E	F, *

در این محیط علاوه بر غذای اصلی (\circ) دو ماده ی غذایی دیگر (*) نیز در نقشه وجود دارند که خوردن آنها ۵ امتیاز مثبت خواهد داشت. این امتیازها نیز هنگام ورود به خانه ی دارای غذا کسب می شوند. در این حالت هم چنان غذای اصلی همان \circ است که خوردن آن باعث تمام شدن بازی می شود و پاداش ۱ دارد. در این حالت به سوال های زیر پاسخ دهید:

(ج) اگر ضریب تخفیف برابر ۱ باشد و هزینه ی زنده ماندن برابر ۱-، سیاست بهینه را برای حالت های مختلف پیدا کنید.

(د) هم چنان فرض کنید ضریب تخفیف برابر ۱ است، هزینه ی زنده ماندن در چه بازه ای می تواند باشد تا پک من با شروع از خانه ی A دقیقا یک * را بخورد؟

حل.

- (آ) سیاست بهینه به صورت زیر خواهد بود:

$$\begin{cases} \pi(A) = \text{راست یا پایین} \\ \pi(B) = \text{راست یا پایین} \\ \pi(C) = \text{پایین} \\ \pi(D) = \text{راست} \\ \pi(E) = \text{راست} \end{cases}$$

(ب) دو مسیر بهینه برای رسیدن به غذا از خانه A وجود دارد. یکی از آن‌ها به صورت زیر است:

$$A \rightarrow B \rightarrow C \rightarrow F$$

با توجه به این که ضریب تخفیف برابر ۰/۵ است یعنی اگر در یک گام به غذا (F) برسیم، یک امتیاز، در دو گام ۰/۵ امتیاز، و در سه گام ۰/۲۵ امتیاز دریافت می‌کنیم. پس مقدار بهینه خانه‌ی A برابر ۰/۲۵ خواهد بود.

- (ج) سیاست بهینه به شکل زیر خواهد بود:

$\pi(s)$	s
پایین	A
راست	C
راست	D و در F غذا است
بالا	D و در F غذا نیست
راست	E و در F غذا است
چپ	E و در F غذا نیست
چپ	F

مابقی حالت‌ها یا تکراری هستند و یا به آن‌ها نمی‌رسیم. در کل سه‌تایی (مکان، بودن غذا در F، بودن غذا در D) حالت سیستم را مشخص می‌کند.

(د) اگر x هزینه زنده ماندن باشد، پاداش نخوردن هیچ کدام از *ها برابر $x + 1$ خواهد بود. پاداش خوردن فقط یکی از *ها برابر $3x + 6$ خواهد بود. و پاداش خوردن دوتا * برابر $7x + 11$ است. برای اینکه فقط یک * خورده شود باید داشته باشیم:

$$\left. \begin{array}{l} 3x + 6 > x + 1 \quad (\text{خوردن یکی صرفه داشته باشد}) \\ 3x + 6 > 7x + 11 \quad (\text{خوردن دو تا صرفه نداشته باشد}) \end{array} \right\} \Rightarrow -2/5 \leq x \leq -1/25$$

۸. (۱۲ نمره) فرض کنید بازی‌ای فقط شامل دو حالت A و B بوده و کنش‌های قابل انجام از هر کدام از این حالت‌ها نیز حرکت به سمت بالا یا پایین است. عاملی در این محیط با استفاده از سیاست π به انجام بازی پرداخته و دنباله‌ی حالت‌ها و پاداش‌های زیر را مشاهده کرده است:

t	s_t	a_t	s_{t+1}	r_t
۰	A	پایین	B	۲
۱	B	پایین	B	-۴
۲	B	بالا	B	۰
۳	B	بالا	A	۳
۴	A	بالا	A	-۱

با در نظر گرفتن ضریب تخفیف برابر ۰/۵ و $\alpha = 0/5$ به سوال‌های زیر پاسخ دهید:

(آ) می‌خواهیم از روش Q-learning برای بروزرسانی Q-value ها استفاده کنیم. ابتدا رابطه‌ی بروزرسانی را برای Q-value ها بنویسید. سپس با در نظر گرفتن مقدار صفر برای مقدار اولیه‌ی Q-value ها مقدار $Q(A)$ و $Q(B)$ را بدست آورید.

(ب) در روش model-based ابتدا تابع گذار $T(s, a, s')$ و تابع پاداش $R(s, a, s')$ را تخمین می‌زنیم. برای این منظور با استفاده از دنباله‌ی حالت‌های مشاهده شده این توابع را به‌ازای مقادیر مختلف s ، a و s' محاسبه کنید (در مجموع مقادیر هشت حالت را باید تخمین بزنید). اگر داده‌های لازم برای تخمین یکی از ورودی‌های هر کدام از توابع وجود ندارد به این موضوع اشاره کنید.

(ج) فرض کنید تجربه‌ی جدیدی از بازی کسب کرده‌ایم و از روی آن توابع \hat{T} و \hat{R} را تخمین زده‌ایم. نتیجه جدول زیر شده است:

s	a	s'	$\hat{T}(s, a, s')$	$\hat{R}(s, a, s')$
A	بالا	A	۱	۱۰
A	پایین	A	۰/۵	۲
A	پایین	B	۰/۵	۲
B	بالا	A	۱	-۵
B	پایین	B	۱	۸

با استفاده از این جدول سیاست بهینه $\hat{\pi}^*(s)$ و تابع ارزش بهینه $\hat{V}^*(s)$ را بدست آورید. (راهنمایی: برای هر x حقیقی که $|x| < 1$ داریم: $\frac{1}{1-x} = 1 + x + x^2 + x^3 + \dots$)

حل.

(آ) رابطه‌ی بروزرسانی Q-value ها به صورت زیر است:

$$Q(s_t, a_t) \leftarrow (1 - \alpha)Q(s_t, a_t) + \alpha \left[r_t + \gamma \max_{a'} Q(s_{t+1}, a') \right]$$

با استفاده از این رابطه:

$$Q(A, \text{پایین}) \leftarrow (1 - 0.5) \times 0 + 0.5 [2 + 0.5 \times 0] = 1$$

$$Q(B, \text{پایین}) \leftarrow (1 - 0.5) \times 0 + 0.5 [-4 + 0.5 \times 0] = -2$$

$$Q(B, \text{بالا}) \leftarrow (1 - 0.5) \times 0 + 0.5 [0 + 0.5 \times 0] = 0$$

$$Q(B, \text{بالا}) \leftarrow (1 - 0.5) \times 0 + 0.5 [3 + 0.5 \times 1] = \frac{7}{4}$$

$$Q(A, \text{بالا}) \leftarrow (1 - 0.5) \times 0 + 0.5 [-1 + 0.5 \times 1] = \frac{-1}{4}$$

پس $Q(A, \text{پایین})$ برابر ۱ و $Q(B, \text{بالا})$ برابر $1/75$ است.

(ب) با استفاده از شمارش تعداد حالت‌ها:

$$\begin{cases} \hat{T}(A, \text{بالا}, A) = ۱ \\ \hat{T}(A, \text{بالا}, B) = ۰ \\ \hat{T}(A, \text{پایین}, A) = ۰ \\ \hat{T}(A, \text{پایین}, B) = ۱ \\ \hat{T}(B, \text{بالا}, A) = ۰/۵ \\ \hat{T}(B, \text{بالا}, B) = ۰/۵ \\ \hat{T}(B, \text{پایین}, A) = ۰ \\ \hat{T}(B, \text{پایین}, B) = ۱ \end{cases} \quad \begin{cases} \hat{R}(A, \text{بالا}, A) = -۱ \\ \hat{R}(A, \text{بالا}, B) = na \\ \hat{R}(A, \text{پایین}, A) = na \\ \hat{R}(A, \text{پایین}, B) = ۲ \\ \hat{R}(B, \text{بالا}, A) = ۳ \\ \hat{R}(B, \text{بالا}, B) = ۰ \\ \hat{R}(B, \text{پایین}, A) = na \\ \hat{R}(B, \text{پایین}, B) = -۴ \end{cases}$$

(ج) ابتدا سیاست بهینه را با توجه به مقدار پاداش‌های داده شده پیدا می‌کنیم:

$$\begin{cases} \hat{\pi}^*(A) = \text{بالا} \\ \hat{\pi}^*(B) = \text{پایین} \end{cases}$$

حال می‌توان مقادیر بهینه‌ی V را با توجه به روابط بلمن حساب کرد:

$$\begin{aligned} \hat{V}^*(A) &= ۱۰ + \lambda * \hat{V}^*(A) \\ &= ۱۰ + ۱۰\lambda + ۱۰\lambda^2 + \dots = \frac{۱۰}{۱ - \lambda} = \frac{۱۰}{۰/۵} = ۲۰ \end{aligned}$$

$$\begin{aligned} \hat{V}^*(B) &= ۸ + \lambda * \hat{V}^*(B) \\ &= ۸ + ۸\lambda + ۸\lambda^2 + \dots = \frac{۸}{۱ - \lambda} = \frac{۸}{۰/۵} = ۱۶ \end{aligned}$$