



هوش مصنوعی

بهار ۱۴۰۲

مدرس: محمد مهدی سمیعی

دانشگاه صنعتی شریف

دانشکده مهندسی کامپیوتر

گردآورندگان: پردیس زهرایی، بنیامین ملکی، امیرحسین رازلیقی

مهلت ارسال: ۱۲ فروردین

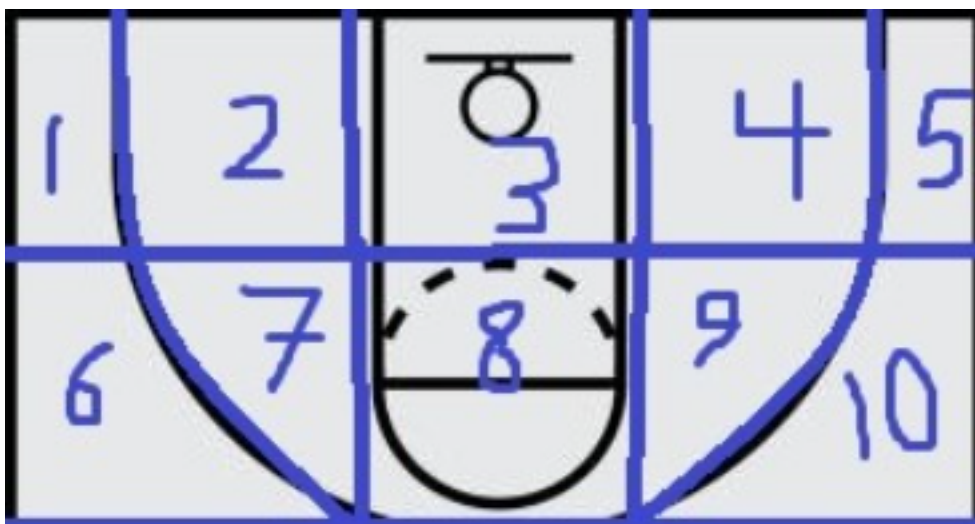
فرایند تصمیم مارکوف و یادگیری تقویتی

تمرین سوم

- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روز مشخص شده است.
- در طول ترم امکان ارسال با تاخیر پاسخ وجود ندارد و پاسخ هایی که بعد از زمان تعیین شده ارسال شوند، پذیرفته نخواهند شد.
- هم کاری و هم فکری شما در انجام تمرین مانعی ندارد اما پاسخ ارسالی هر کس حتما باید توسط خود او نوشته شده باشد.
- در صورت هم فکری و یا استفاده از هر منابع خارج درسی، نام هم فکران و آدرس منابع مورد استفاده برای حل سوال مورد نظر را ذکر کنید.
- لطفا تصویری واضح از پاسخ سوالات نظری بارگذاری کنید. در غیر این صورت پاسخ شما تصحیح نخواهد شد.
- توجه داشته باشید که تمرین نمره امتیازی ندارد.

سوالات نظری (۱۱۵ نمره)

۱. (۵۰ نمره) فرض کنید یک ربات بسکتبالیست داریم که در زمین بازی زیر قرار دارد.



شکل ۱: زمین بازی

این ربات، در هر لحظه، می تواند در یکی از ۱۰ خانه ی مشخص شده در زمین بازی حضور داشته باشد. همچنین، در هر مرحله، ربات می تواند یا به یکی از خانه های مجاور برود و یا اینکه در همان خانه ای که قرار دارد بماند و اقدام به پرتاب توپ به سمت حلقه کند. منظور از خانه ی مجاور نیز، خانه ای است که یک ضلع مشترک با خانه ی فعلی داشته باشد. فرضاً، از خانه ی ۹ می تواند به خانه های ۴، ۸ یا ۱۰ برود. در صورتی که تصمیم بگیرد به یکی از خانه های چپ، راست، بالا و یا پایینش برود، به احتمال ۱۰۰ درصد این کار با موفقیت انجام می شود. از طرفی اگر action ربات ما، حرکت کردن باشد، امتیاز (reward) 1- را

دریافت می‌کند (مستقل از وضعیتی که در آن قرار دارد). از طرفی اگر انتخاب ربات، اقدام به پرتاب (shoot) باشد، بازی با تمام می‌رسد. همچنین، با توجه به خانه‌ای که در آن قرار دارد (فرضا s)، به احتمال $P_{goal}(s)$ پرتابش گل می‌شود. متناظراً، به احتمال $1 - P_{goal}(s)$ اقدامش با شکست مواجه می‌شود. برای وضعیت‌های مختلف، احتمالات زیر را داریم:

$$P_{goal}(s_1) = P_{goal}(s_5) = P_{goal}(s_6) = P_{goal}(s_{10}) = 0.6$$

$$P_{goal}(s_2) = P_{goal}(s_7) = P_{goal}(s_8) = P_{goal}(s_9) = P_{goal}(s_4) = 0.75$$

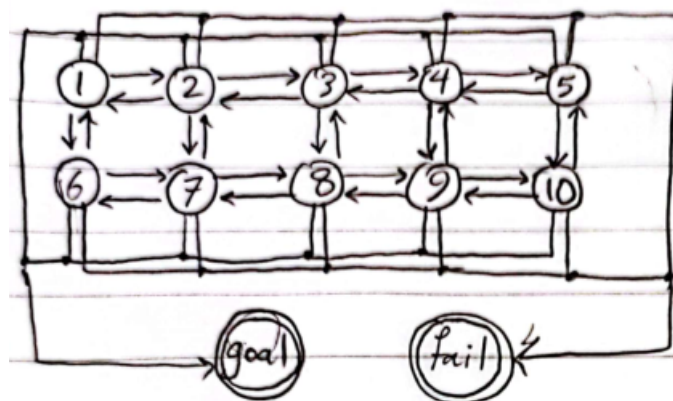
$$P_{goal}(s_3) = 0.9$$

در صورتی که پرتابش گل نشود، امتیاز 10- را دریافت می‌کند. اما اگر گل شود، بسته به اینکه از کجا اقدام به پرتاب کرده باشد، امتیاز متفاوتی دریافت می‌کند:

$$\forall i \in \{1, 5, 6, 10\} : R(s_i, shoot, goal) = 3$$

$$\forall i \in \{2, 3, 4, 7, 8, 9\} : R(s_i, shoot, goal) = 2$$

برای واضح‌تر شدن صورت سوال، شماتیک MDP شرح داده شده، در زیر آمده است:



شکل ۲: MDP

الف) به کمک روش Value Iteration، تابع ارزش وضعیت‌ها ($V(s_i)$) را برای تمامی استیت‌های MDP محاسبه کنید. توجه کنید که برای محاسبات خود، صرفاً ۲ بار ایتريشن انجام دهید، به این معنا که Value‌هایی که محاسبه می‌کنید، مقادیر درخت Expectimax تا عمق ۲ باشند. ($\gamma = 1$) (۳۰ نمره)

ب) حال فرض کنید که از توابع R و T اطلاعی نداریم و به جای آن، تجربه‌هایی از محیط به دست آورده ایم که در جدول ۱ (در صفحه‌ی بعد) آمده‌اند.

حال الگوریتم Q-learning را اجرا نموده و $Q(s,a)$ را به ازای state‌ها و action‌هایی که در تجربه‌های ما دیده شده‌اند و تغییر می‌کنند، محاسبه کنید. در حقیقت، در محاسبات مربوط به هر Episode، صرفاً $Q(s,a)$ ‌های حاضر در آن Episode را محاسبه کنید و اگر $Q(s,a)$ یک خاص دچار تغییر نمی‌شود، نیازی به ذکر کردن مجدد مقدار آن نیست. با این فرض که همه‌ی $Q(s,a)$ ‌ها در ابتدا صفر باشند و با نرخ یادگیری $\alpha = 0.5$ و ضریب تخفیف $\gamma = 1$ ، محاسبات خود را انجام دهید. (۲۰ نمره)

۲. (۱۵ نمره) عبارات صحیح و غلط را مشخص و به صورت مختصر برای عبارات غلط دلیل بیاورید.

الف) در فرایند value iteration اگر مقدار تخفیف بین ۰ و ۱ باشد value‌ها حتماً همگرا می‌شوند. (۳ نمره)
 ب) سیاست‌های پیدا شده توسط policy iteration از سیاست‌های پیدا شده توسط value iteration ضعیف‌تر هستند. (۳ نمره)

پ) در Q-learning می‌توان بدون رسیدن به سیاست بهینه، به Q^* بهینه رسید. (۳ نمره)

ت) یک تخفیف کوچک‌تر از ۱ همواره می‌تواند به عنوان یک پاداش منفی در نظر گرفته شود. (۳ نمره)

ث) یک پاداش بزرگ منفی می‌تواند باعث رفتار حریصانه شود. (۳ نمره)

	S	A	S'	Reward
Episode 1	s_8	move R	s_9	-1
	s_9	move U	s_6	-1
	s_6	shoot	goal	+1
Episode 2	s_8	move R	s_7	-1
	s_7	move R	s_8	-1
	s_8	shoot	goal	+1
Episode 3	s_8	move R	s_9	-1
	s_9	move R	s_{10}	-1
	s_{10}	shoot	goal	+2
Episode 4	s_8	move R	s_7	-1
	s_7	move L	s_8	-1
	s_8	shoot	fail	-5
Episode 5	s_8	move R	s_9	-1
	s_9	move U	s_{10}	-1
	s_{10}	shoot	fail	-5

جدول ۱: تجربه‌های بدست آمده از محیط

۳. (۲۵ نمره) در یک شب سرد زمستانی، شما توسط گروه مافیایی Tarasht، اسیر می‌شوید. رئیس گروه، که از قضا انسان فرهیخته‌ای است، پس از بررسی لپ‌تاپ شما و پاسخ‌هایتان برای تمرین هوش مصنوعی، از توانایی شما در حل مسئله‌های قبلی خوشش می‌آید، و با طراحی چالشی به شما فرصتی می‌دهد تا فرار کنید. ابتدا در زندانی مانند شکل زیر قرار دارید:

5		S			10
			0	0	

اگر بر روی یک سلول شماره‌گذاری شده (با شماره مثبت) باشید، تنها اقدام موجود، خروج از آن است و زمانی که از آن خارج می‌شوید، پاداشی برابر با عدد روی سلول دریافت می‌کنید (در صورتی که از خانه با پاداش ۱۰ خارج شوید، دری به رویتان باز می‌شود که سر از خیابان تیموری در می‌آورید). در هر سلول دیگر (بدون شماره)، اقدامات موجود، حرکت به سمت شرق یا غرب است. رئیس مافیا، برای اینکه کار را سخت‌تر کند، چندین چاله هم در مسیر قرار داده است (خانه‌های با شماره صفر). اگر در سلول‌های بالایی این چاله‌ها قرار داشته باشید، با حرکت کردن به شرق یا غرب، به احتمال $1-p$ در چاله‌ها می‌افتید و به احتمال p حرکتی که قصد داشتید انجام دهید، به درستی انجام می‌شود.

سیاست‌های مختلفی به ذهن شما می‌رسند. در هر حالت، $V_{\pi}(s)$ را برای هر سلول بدون شماره به دست آورید.

- الف) همواره به سمت شرق حرکت کنید. (۵ نمره)
 ب) همیشه به سمت غرب حرکت کنید. (۵ نمره)
 ج) برای هر کدام از سیاست های بالا دامنه p را تعیین کنید. (۱۰ نمره)
 د) در ادامه تفاوت policy iteration با value iteration را بنویسید. (۵ نمره)

۴. (۲۵ نمره) می‌دانیم که اگر بخواهیم یک تابع utility برای ارزش‌دهی به دنباله‌ای از state‌ها ارائه دهیم که به صورت stationary (مستقل از زمان رسیدن به state)، آن دنباله را ارزش‌گذاری کند، باید یکی از دو روش زیر باشد:

$$(۱) \text{ Additive} : U_1[s_0, s_1, s_2, \dots] = \sum_{i=0}^{\infty} R(s_i)$$

$$(۲) \text{ Discounted} : U_2[s_0, s_1, s_2, \dots] = \sum_{i=0}^{\infty} \gamma^i \cdot R(s_i)$$

- الف) دو مورد از مشکلات تابع U_1 را ذکر کنید. (۵ نمره)
 ب) اگر دنباله وضعیت‌های ما محدود باشد (مثلاً $[s_0, s_1, \dots, s_n]$)، آیا همچنان روش‌های با stationary preference می‌توانند خوب باشند؟ (۵ نمره)
 ج) فرض کنیم تمامی $R(s_i)$ ‌ها متعلق به بازه $[R_{min}, R_{max}]$ باشند. حال، برای تابع U_2 ، کران بالا و پایینی بیابید و نشان دهید که همواره کران‌دار است. (۱۵ نمره)

سوالات عملی (۸۵ نمره)

سوالات عملی در فایل جویتر نوت‌بوک موجود هستند.

- تمرین MDP (۴۵ نمره)
- تمرین RL (۴۰ نمره)