

سوال ۱.

الف

روش Value Iteration یکی از الگوریتم‌های حل مسئله در یادگیری تقویتی است که برای حل مسئله‌های مارکوف تصمیم‌گیری (MDP) استفاده می‌شود. این الگوریتم به صورت تکراری و تاملاتی عمل می‌کند و به طور تدریجی مقادیر ارزش برای هر حالت را به‌روزرسانی می‌کند تا به یک تابع ارزش بهینه برسد.

برای استفاده از الگوریتم Value Iteration، نیازمند تعریف MDP هستیم. MDP توسط چهار مؤلفه اصلی تعریف می‌شود:

مجموعه حالات (States): مجموعه‌ای از حالات قابل مشاهده در مسئله را نشان می‌دهد.

مجموعه عمل‌ها (Actions): مجموعه‌ای از عمل‌های قابل انجام در هر حالت را نشان می‌دهد.

تابع انتقال (Transition Function): تابعی است که نشان می‌دهد با انجام یک عمل در یک حالت، ما به چه حالتی خواهیم رفت و با چه احتمالی.

تابع پاداش (Reward Function): تابعی است که برای هر حالت و عمل، پاداش مشخص می‌کند.

هدف Value Iteration یافتن تابع ارزش بهینه است که برای هر حالت، ارزش بهینه انتظاری را نشان می‌دهد. الگوریتم Value Iteration به صورت تکراری عمل می‌کند و به تدریج مقادیر ارزش را به‌روزرسانی می‌کند تا به تابع ارزش بهینه برسد.

حالا که تمام حالات ۰ هستند، تا عمق ۲ را محاسبه می‌کنیم. (همچنین از نوشتن یک سری محاسبات عددی چشم‌پوشی می‌کنیم چون محاسبات بسیار زیادی دارد این سوال.)

$$V(s) = \max \left[\sum p(s'|s, a) \cdot (R(s, a, s') + \gamma \cdot V(s')) \right]$$

پس داریم:

$$V_1(S_1) = \max[1 \times [-1 + V(S_2)], 1 \times [-1 + V(S_6)], 0.6 \times (3 + V_{goal}) + 0.4 \times (-10 + V_{fail})]$$

$$\rightarrow V_1(S_1) = \max[-1, -1, -2/2] = -1$$

به همین شکل برای باقی نیز خواهیم داشت:

$$\rightarrow V_1(S_2) = \max[-1, -1, -1, -1] = -1$$

$$\rightarrow V_1(S_3) = \max[-1, -1, -1, 0.8] = 0.8$$

$$\rightarrow V_1(S_{\P}) = \max[-1, -1, -1, -1] = -1$$

$$\rightarrow V_1(S_{\Delta}) = \max[-1, -1, -2/2] = -1$$

$$\rightarrow V_1(S_{\S}) = \max[-1, -1, -2/2] = -1$$

$$\rightarrow V_1(S_V) = \max[-1, -1, -1, -1] = -1$$

$$\rightarrow V_1(S_{\Lambda}) = \max[-1, -1, -1, -1] = -1$$

$$\rightarrow V_1(S_{\P}) = \max[-1, -1, -1, -1] = -1$$

$$\rightarrow V_1(S_{1.}) = \max[-1, -1, -2/2] = -1$$

حالا يك عمق ديگر پيش مي‌رويم:

$$\rightarrow V_2(S_1) = \max[-2, -2, -2/2] = -2$$

$$\rightarrow V_2(S_2) = \max[-2, -2, -1/2, -1] = -1/2$$

$$\rightarrow V_2(S_3) = \max[-2, -2, -1/8, -2] = -1/8$$

$$\rightarrow V_2(S_{\P}) = \max[-2, -2, -1, -1/2] = -1/2$$

$$\rightarrow V_2(S_{\Delta}) = \max[-2, -2, -2/2] = -2$$

$$\rightarrow V_2(S_{\S}) = \max[-2, -2, -2/2] = -2$$

$$\rightarrow V_2(S_V) = \max[-2, -2, -1/2, -1] = -1/2$$

$$\rightarrow V_2(S_{\Lambda}) = \max[-2, -2, -1/2, -1] = -1/2$$

$$\rightarrow V_2(S_{\P}) = \max[-2, -2, -2/2] = -2$$

$$\rightarrow V_2(S_{1.}) = \max[-2, -2, -1/2, -1] = -1/2$$

ب

الگوریتم Q-Learning یکی از الگوریتم‌های یادگیری تقویتی است که برای حل مسئله یادگیری تقویتی بدون مدل (Model-Free) استفاده می‌شود. این الگوریتم به صورت تکراری و تاملیتی عمل می‌کند و به تدریج تابع Q-Value را به‌روزرسانی می‌کند تا به تابع Q-Value بهینه برسد.

برای استفاده از الگوریتم Q-Learning، نیازمند تعریف یک MDP هستیم. MDP توسط چهار مؤلفه اصلی تعریف می‌شود:

مجموعه حالات (States): مجموعه‌ای از حالات قابل مشاهده در مسئله را نشان می‌دهد.

مجموعه عمل‌ها (Actions): مجموعه‌ای از عمل‌های قابل انجام در هر حالت را نشان می‌دهد.

تابع انتقال (Transition Function): تابعی است که نشان می‌دهد با انجام یک عمل در یک حالت، ما به چه حالتی خواهیم رفت و با چه احتمالی.

تابع پاداش (Reward Function): تابعی است که برای هر حالت و عمل، پاداش مشخص می‌کند.

هدف Q-Learning یافتن تابع Q-Value بهینه است که برای هر حالت و عمل، ارزش بهینه انتظاری را نشان می‌دهد. الگوریتم Q-Learning به صورت تکراری و مبتنی بر تجربه (Experience-Based) عمل می‌کند و به تدریج مقادیر Q-Value را به‌روزرسانی می‌کند تا به تابع Q-Value بهینه برسد.

حالا الگوریتم را ران می‌کنیم:

Episode ۱ :

$$Q(s_8, R) = 0.5 \times -1 = -0.5$$

$$Q(s_9, U) = 0.5 \times -1 = -0.5$$

$$Q(s_4, S) = 0.5 \times 1 = 0.5$$

Episode ۲ :

$$Q(s_8, S) = 0.5 \times 1 = 0.5$$

$$Q(s_v, R) = 0.5 \times 1 = 0.5$$

$$Q(s_8, R) = 0.5 \times 0.5 + 0.5 = 0.75$$

Episode ۳ :

$$Q(s_1, S) = 0.5 \times 2 = 1$$

$$Q(s_4, R) = 0.5 \times 1 = 0.5$$

$$Q(s_8, R) = 0.75 + 0.5 \times 0.25 = 0.875$$

Episode 4 :

$$Q(s_8, S) = 0.5 \times 4/5 = 2/5$$

$$Q(s_7, L) = 0.5 \times 0.5 = 0.25$$

$$Q(s_8, R) = 0.875 + 0.5 \times 0.125 = 0.9375$$

Episode 5 :

$$Q(s_1, S) = 1 + 0.5 \times 2 = 2$$

$$Q(s_4, U) = 0.5 + 0.5 \times -0.5 = 0.25$$

$$Q(s_8, R) = 0.9375 + 0.5 \times 0.625 = 0.96875$$

سوال ۲.

الف

صحیح

ب

نه، این جمله درست نیست. در عمل، سیاست‌های پیداشده توسط الگوریتم‌های Value و Policy Iteration ممکن است با هم تفاوت داشته باشند و بسته به مسئله و شرایط مختلف، نتایج ممکن است متفاوت باشند. اما در حالت کلی، نمی‌توان گفت که سیاست‌های پیداشده توسط Policy Iteration ضعیف‌تر از سیاست‌های پیداشده توسط Value Iteration هستند یا برعکس.

در Policy Iteration، به‌روزرسانی‌ها در دو مرحله انجام می‌شود: در مرحله اول، تابع ارزش بهینه برای سیاست فعلی به‌روزرسانی می‌شود و در مرحله دوم، سیاست جدید بر اساس تابع ارزش بهینه جدید تعیین می‌شود. به این ترتیب، الگوریتم Policy Iteration بهبودی پیوسته در سیاست و تابع ارزش بهینه ایجاد می‌کند.

در Value Iteration، تابع ارزش بهینه به صورت تکراری به‌روزرسانی می‌شود و در هر مرحله، تابع ارزش بهینه به تدریج به مقدار نهایی همگرا می‌شود. در نهایت، می‌توان از تابع ارزش بهینه برای تولید سیاست بهینه استفاده کرد. به طور کلی، می‌توان گفت که Value Iteration و Policy Iteration به دست آوردن سیاست بهینه را هدف دارند و در نهایت به نتایج مشابهی می‌رسند. اما در مسائلی که سیاست‌ها بیش از یک سیاست بهینه دارند، ممکن است سیاست‌های پیداشده توسط Value Iteration و Policy Iteration با هم تفاوت داشته باشند. در این حالت، می‌توان گفت که هیچ‌کدام از این الگوریتم‌ها برتری مطلق را دارا نیستند و نتیجه نهایی بسته

پ

صحیح

ت

نه، این جمله درست نیست. تخفیف در الگوریتم‌های یادگیری تقویتی برای کنترل تأثیر پاداش‌های آینده استفاده می‌شود و به عنوان یک عامل تعدیل‌کننده (discounting factor) عمل می‌کند. ارزش تخفیف باید بین ۰ و ۱ قرار بگیرد و انتخاب مقدار مناسب آن بسته به مسئله و شرایط خاص است.

تخفیف باعث می‌شود پاداش‌های آینده کمتر ارزش داشته باشند و به ارزش پاداش‌های فعلی بیشتری توجه شود. با تخفیف بالا، عامل به دنبال دریافت پاداش‌های نزدیک در زمان حال است و با تخفیف کم، عامل تمایل دارد به پاداش‌های آینده بیشتر توجه کند. اما استنتاج کلی این نیست که تخفیف کوچکتر از ۱ همواره به عنوان یک پاداش منفی در نظر گرفته می‌شود.

ث

صحیح

الف

از راست شروع می‌کنیم:

$$\begin{aligned} p(1 \cdot 0 + \gamma \times 0) + (1-p)(0 + \gamma \times 0) &= 1 \cdot p \\ p(0 + \gamma \times 1 \cdot p) + (1-p)(0 + \gamma \times 0) &= 1 \cdot p^2 \gamma \\ 1(0 + \gamma \times 1 \cdot p^2 \gamma) &= 1 \cdot p^2 \gamma^2 \\ 1(0 + \gamma \times 1 \cdot p^2 \gamma^2) &= 1 \cdot p^2 \gamma^3 \end{aligned}$$

ب

از چپ شروع می‌کنیم.

$$\begin{aligned} 1(5 + \gamma \times 0) &= 5 \\ 1(0 + \gamma \times 5) &= 5\gamma \\ p(0 + \gamma \times 5\gamma) + (1-p)(0 + \gamma \times 0) &= 5p\gamma^2 \\ p(0 + \gamma \times 5p\gamma^2) + (1-p)(0 + \gamma \times 0) &= 5p^2\gamma^3 \end{aligned}$$

ج

باید معادله‌ی زیر را حل کنیم. به صورت کلی نیز می‌دانیم که $0 < p, \gamma < 1$ است.

$$\begin{aligned} 1 \cdot p^2 \gamma^2 > 5\gamma &\implies 2p^2 \gamma > 1 \implies p^2 > \frac{1}{2\gamma} \implies \\ p &> \frac{1}{\sqrt{2\gamma}} \end{aligned}$$

د

در یادگیری تقویتی (Reinforcement Learning) ،

Policy Iteration و Value Iteration دو روش اصلی برای یافتن بهترین راه حل (optimal policy) در یک محیط تصادفی هستند.

Value Iteration: این الگوریتم از تخمین مقدار فعلی یک حالت شروع می‌کند و سپس برای هر حالت امتیاز تمام اقدامات ممکن را محاسبه می‌کند. این امتیازات به روزرسانی می‌شوند و این فرایند تا زمانی که تغییرات کمتر از یک آستانه خاص شود، ادامه می‌یابد. در نهایت، سیاست بهینه براساس تابع ارزش به روز شده تولید می‌شود.

Policy Iteration: در این روش، ما از یک سیاست تصادفی شروع می‌کنیم و سپس در یک حلقه دو قدم را تکرار می‌کنیم: (a) ثابت کردن سیاست و محاسبه تابع ارزش برای آن سیاست (Policy Evaluation) ، و

تفاوت اصلی این دو روش در نحوه‌ی به‌روزرسانی سیاست است. در Value Iteration ، ما فقط یک بار سیاست را به روز می‌کنیم (پس از اینکه تابع ارزش کافیا ثابت شده باشد)، در حالی که در Policy Iteration ، ما سیاست را در هر مرحله به روز می‌کنیم.

یکی از نکاتی که باید توجه داشته باشید این است که در برخی موارد، Policy Iteration ممکن است سریع‌تر از Value Iteration به سیاست بهینه برسد، زیرا هر بروزرسانی در Policy Iteration می‌تواند به سیاست بهینه

نزدیکی تر شود، در حالی که در Value Iteration، سیاست بهینه فقط در انتها تولید می شود. با این حال، Value Iteration می تواند در محیط هایی با فضای حالت یا عمل بزرگ یا در مواردی که محاسبات Policy Evaluation زمان بر هستند، کارآمدتر باشد.

سوال ۴.

الف

وقتی از عبارت «در صورتی که discount را نداشته باشیم» استفاده می کنیم، منظور از discount میزان تاثیر یا اهمیتی است که به مراحل آینده در محاسبه reward نسبت داده می شود. در الگوریتم ها و مدل های مختلف یادگیری تقویتی، از این عامل برای تعیین اهمیت مراحل آینده در محاسبه reward استفاده می شود. با داشتن discount، مقدار reward برای مراحلی که در آینده قرار خواهند گرفت، کاهش می یابد.

بنابراین، در حالتی که discount را نداشته باشیم، ممکن است که به دلیل عدم توجه به اهمیت مراحل آینده، reward کلی که در ۱۰۰۰ مرحله بدست می آید، بیشتر به نظر برسد و ما بیشتر تمایل داشته باشیم آن را انتخاب کنیم. اما وقتی discount وجود دارد، reward ی که بعد از ۱۰ مرحله بدست می آید ممکن است اهمیت بیشتری داشته باشد و ما تمایل داشته باشیم آن را انتخاب کنیم.

همچنین، ممکن است که در مواقعی که دو حالت به تعداد مراحل یکسانی کمیت مشابهی از reward را ارائه می دهند، ما متردد شویم و تصمیم گیری را به تعویق بیندازیم. این به معنای این است که می توانیم در انتخاب بین این دو حالت کمی اضافه زمان صرف کنیم تا مطمئن شویم کدام گزینه بهتر است.

به طور کلی، در صورتی که عواملی مانند discount در نظر گرفته نشوند، اهمیت مراحل آینده و تفاوت بین reward ها کاهش می یابد و این می تواند در فرآیند تصمیم گیری تاثیر داشته باشد.

ب

بله، بالاخره این مراحل باید طی شوند و در نهایت به یک جواب همگرا می شوند.

ج

در رابطه ای ارائه شده برای جفت حالت ها، $R_{max|min}$ نشان دهنده ی حداکثر یا حداقل مقدار reward است که در آن جفت حالت بدست می آید. همچنین، γ مقداری بین ۰ و ۱ است که به آن عامل تخفیف گفته می شود و نشان دهنده ی اهمیت مراحل آینده در محاسبه ی reward است.

در صورتی که $\gamma < 1$ باشد، جمع هندسی $\sum \gamma^i$ به مقدار $\frac{1}{1-\gamma}$ همگرا می شود. این به این معنی است که اهمیت مراحل آینده در محاسبه ی reward با افزایش تعداد مراحل کاهش می یابد، اما به طور متناسب با عامل تخفیف گفته می شود. به عبارت دیگر، هر چه γ کوچک تر باشد، اهمیت مراحل آینده بیشتر خواهد بود.

اگر $\gamma \geq 1$ باشد، مجموعه $\sum \gamma^i$ به نام همگرایی دارای مقدار بی نهایت می شود. این به این معنی است که در این حالت، اهمیت مراحل آینده بی نهایت است و تاثیر مراحل از لحاظ زمانی نامحدود است.

بنابراین، برای اطمینان از همگرایی وجود عامل تخفیف، مقدار γ باید بین ۰ و ۱ قرار گیرد. در نهایت، مقدار U که نشان دهنده ی utility حالت است، به $\frac{R_{max|min}}{1-\gamma}$ همگرا می شود. این به این معنی است که با در نظر گرفتن عامل

تخفیف، utility حالت برابر با تقسیم مقدار حداکثر یا حداقل reward بر تفاوت یک منهای γ خواهد بود. این مقدار نشان‌دهنده‌ی ارزش حالت در نظر گرفته شده با اهمیت مراحل آینده است.

موفق باشید.