

تمرین سری اول

امیرحسین انتظاری

۲۲ آبان ۱۴۰۲

۱.۰ سوالات و پاسخ ها

۱. الف)

پاسخ:

هنگامی که یک تنظیم (adjustment) پاداش ثابت c به همه transition ها در MDP معرفی می شود، تابع مقدار بهینه جدید V_2 می تواند به صورت زیر بیان شود:

$$V_2(s) = V_1(s) + \frac{c}{1-\gamma}$$

که در آن $V_1(s)$ تابع مقدار بهینه اصلی و γ ضریب تخفیف (discount factor) است. استدلال پشت این عبارت این است که معرفی یک پاداش ثابت c برای همه انتقال ها به دلیل فرمول سری هندسی ارزش هر حالت را به میزان $\frac{c}{1-\gamma}$ افزایش می دهد. ضریب تخفیف γ تأثیر این تنظیم را محدود می کند.

اکنون، با توجه به استراتژی بهینه (π_2) در این سناریوی تغییر یافته، نکته مهم این است که توجه داشته باشیم که استراتژی (policy) بهینه از تابع مقدار بهینه مشتق شده است. استراتژی π_2 همچنان با انتخاب اعمالی که تابع مقدار بهینه جدید V_2 را به حداکثر می رساند تعیین می شود.

به طور کلی، اگر استراتژی بهینه π_1 برای MDP اصلی شامل اقداماتی (actions) نباشد که در مرز تصمیم گیری هستند (یعنی اقداماتی که به دلیل ساختار پاداش بی تفاوت بودند)، معرفی پاداش ثابت تنظیم c ممکن است استراتژی بهینه را تغییر ندهد. عامل همچنان اعمال یکسان را در هر حالت ترجیح می دهد زیرا تفاوت نسبی ارزش بین اقدامات در هر حالت یکسان باقی می ماند.

با این حال، اگر استراتژی اصلی شامل اعمالی در مرز تصمیم گیری باشد (جایی که عامل بین چندین اعمال بی تفاوت بود)، تعدیل پاداش ثابت c به طور بالقوه می تواند ترجیحات را تغییر دهد و منجر به تغییر در استراتژی بهینه شود. عامل کلیدی در اینجا این است که آیا تعدیل پاداش برای تغییر مقادیر نسبی اعمال در حالت های خاص کافی است و بر تصمیم گیری عامل تأثیر می گذارد. به عبارت دیگر استراتژی بهینه ممکن است بسته به ویژگی های MDP اصلی، تنظیم پاداش c و ویژگی های استراتژی بهینه اصلی π_1 تغییر کند یا تغییر نکند.

۲. ب)

پاسخ:

اینکه آیا استراتژی بهینه پس از مقیاس بندی پاداش تغییر می کند یا ثابت می ماند، به ویژگی های خاص MDP اصلی و فاکتور مقیاس c بستگی دارد. انتخاب c که برای آن همه استراتژی ها بهینه هستند، مستلزم حفظ نظم نسبی مقادیر عمل (action) در حالت ها است. وقتی همه پاداش ها در MDP با یک ثابت $c \in R$ مقیاس بندی می شوند، تابع مقدار بهینه جدید V_2 می تواند به صورت زیر بیان شود:

$$V_2(s) = c \cdot V_1(s).$$

استدلال پشت این عبارت این است که مقیاس دادن پاداش ها، تابع ارزش را به طور متناسب مقیاس می دهد.

حال در مورد استراتژی بهینه (π_2) در این سناریو:

۱. استراتژی بهینه باقی بماند π_1 - اگر استراتژی بهینه اولیه π_1 به گونه ای باشد که تصمیم گیری عامل بر اساس مقادیر نسبی اقدامات باشد نه بزرگی مطلق آن ها، ممکن است استراتژی بهینه بدون تغییر باقی بماند. این به این دلیل است که مقیاس کردن همه پاداش ها توسط یک ثابت c ترتیب نسبی مقادیر عمل را در هر حالت تغییر نمی دهد.

- اگر π_1 در MDP اصلی بهینه بود، اگر عامل به انتخاب اعمال در هر حالت بر اساس ترتیب نسبی مقادیر مقیاس شده آنها ادامه دهد، در MDP مقیاس شده بهینه می ماند.

- تابع مقدار حاصل در این مورد $V_2(s) = c \cdot V_1(s)$ است.

۲. استراتژی بهینه تغییر یابد: - اگر استراتژی بهینه اولیه π_1 شامل اعمال در مرز تصمیم‌گیری می‌شد که در آن عامل بین چندین اعمال به دلیل مقادیر پاداش خاص بی‌تفاوت بود، مقیاس‌بندی پاداش‌ها ممکن است استراتژی بهینه را تغییر دهد.

- استراتژی بهینه ممکن است تغییر کند اگر مقیاس بندی ترتیب نسبی مقادیر عمل را در ایالت‌ها تغییر دهد و بر تصمیم‌گیری عامل تأثیر بگذارد.

- در این حالت، تابع مقدار حاصل ممکن است یک رابطه خطی ساده با تابع مقدار اصلی V_1 نداشته باشد.

۳. انتخاب c برای بهینه بودن همه استراتژی‌ها: - انتخابی از c وجود دارد به طوری که همه استراتژی‌ها بهینه هستند اگر و فقط اگر MDP اصلی یک استراتژی بهینه منحصر به فرد π_1 داشته باشد و فاکتور مقیاس c به گونه ای انتخاب شود که چنین باشد. ترتیب نسبی مقادیر عمل را در هیچ حالتی تغییر ندهید.

- اگر c به گونه ای انتخاب شود که $c > 0$ و ترتیب نسبی مقادیر عملکرد در حالت‌ها پس از مقیاس بندی یکسان باقی بماند، همه سیاست‌ها، از جمله استراتژی بهینه اصلی π_1 در MDP مقیاس شده بهینه باقی بماند.

۳.ج)
پاسخ:

وجود حالت‌های پایانی در MDP در واقع می‌تواند بر پاسخ به بخش (a) تأثیر بگذارد. در بخش (الف)، ما در مورد معرفی یک تنظیم پاداش ثابت c برای همه انتقال‌ها در یک MDP بدون حالت پایانی بحث کردیم. هنگامی که حالت‌های پایانی معرفی می‌شوند، پویایی MDP تغییر می‌کند و تأثیر تعدیل پاداش ممکن است متفاوت باشد.

بیایید مثالی از MDP را در نظر بگیریم که در آن پاسخ‌ها به بخش (الف) و قسمت (ج) متفاوت است: ****MDP مثال**** - حالات: $S = \{s_1 s_2 s_3 T\}$ ، که در آن T یک حالت پایانی است. - عمل‌ها (actions): $A = \{a_1, a_2\}$. - انتقال و پاداش: $s_1 \xrightarrow{a_1} s_2$ - پاداش ۱ با $R(s_1 a_1) = 1$. - $s_2 \xrightarrow{a_1} s_3$ - پاداش ۲ با $R(s_2 a_1) = 2$. - $s_3 \xrightarrow{a_1} T$ - پاداش ۰ با $R(s_3 a_1) = 0$. - $s_1 \xrightarrow{a_2} T$ - پاداش ۱ با $R(s_1 a_2) = 1$. - $s_2 \xrightarrow{a_2} T$ - پاداش ۵ با $R(s_2 a_2) = 5$. - $s_3 \xrightarrow{a_2} T$ - پاداش ۱ با $R(s_3 a_2) = 1$. در این MDP، سیاست بهینه π_1 بدون حالت‌های پایانه ممکن است شامل انجام عمل a_1 در s_1 و s_2 و انجام عمل a_2 در s_3 . تابع مقدار بهینه اصلی V_1 بر اساس این سیاست تعیین می‌شود.

حال، اگر یک تنظیم پاداش ثابت c را برای همه انتقال‌ها در MDP، از جمله مواردی که به حالت‌های پایانی منتهی می‌شوند، معرفی کنیم، تابع مقدار بهینه جدید V_2 تحت تأثیر قرار می‌گیرد. تأثیر روی استراتژی بهینه به ویژگی‌های تعدیل پاداش و عامل تخفیف بستگی دارد.

در مقابل، اگر حضور حالت‌های پایانی را در قسمت (c) در نظر بگیریم، مقیاس کردن همه پاداش‌ها با یک ثابت c همچنان بر تابع مقدار تأثیر می‌گذارد، اما استراتژی بهینه ممکن است تا زمانی که ترتیب نسبی عمل بدون تغییر باقی بماند. ارزش‌های درون ایالت‌ها حفظ می‌شود. با این حال، معرفی حالت‌های پایانه به MDP پیچیدگی می‌افزاید و تأثیر آن بر استراتژی بهینه باید در چارچوب مقیاس پاداش و *discount factor* خاص تحلیل شود.