



تمرین هفتم - بخش دوم

شماره دانشجویی: ۹۸۱۰۱۰۷۴

محمدجواد هزاره

۱ سوال ۱

آ) استفاده از Q-Value ها مناسب تر خواهد بود. با دانستن Q-Value ها بهبود دادن سیاست داده شده بدون نیاز به دانستن T و R شدنی خواهد بود. چرا که اگر π سیاست داده شده باشد، داریم:

$$\hat{\pi}(s) = \arg \max_a Q^\pi(s, a)$$

بنابراین برای سنجش یک سیاست بهتر است از Q-Value استفاده کنیم.

ب) استفاده از policy iteration باعث می شود نتوانیم کل فضا را explore کنیم. در این روش با استفاده از یک policy اولیه شروع به سنجش آن و آپدیت آن می کنیم و این باعث می شود اگر policy اولیه به یک سو گرایش داشته باشد، الگوریتم بهترین سیاستی را پیدا کند که به بهترین ارزش های آن سو بگراید. در حالی که ممکن است در گوشه ی دیگری از فضای حالت، حالت های با پاداش بهتر وجود داشته باشد که با استفاده از این روش عامل هرگز به آن حالت ها نخواهد رفت. به عبارتی دیگر عامل گمان می کند که به بهترین سیاست با دریافت بهترین پاداش رسیده در حالی که اصلاً تمام حالات مسئله را بررسی نکرده است. برای رفع این مشکل، از روش ϵ -greedy استفاده می شود که با انجام حرکت های تصادفی با احتمال ϵ ، به عامل این امکان را می دهد که در بعضی از مواقع به حالت های جدید رفته و بتواند تمام فضای حالت مسئله را کشف کند و به بهترین پاداش ممکن دست یابد.

ج) می خواهیم اثبات کنیم که اگر از s شروع کرده و سیاست π' را دنبال کنیم، امید ریاضی مطلوبیت یا همان $V^{\pi'}(s)$ ، بیش تر از حالتی خواهد شد که سیاست π را دنبال کنیم. برای اثبات، از استقرا استفاده می کنیم. تعریف می کنیم:

$$V_n(s) := \begin{pmatrix} \text{امید ریاضی مطلوبیت، اگر از } s \text{ شروع کرده و } n \text{ قدم اول را از سیاست } \pi' \\ \text{و در ادامه از سیاست } \pi \text{ پیروی کنیم.} \end{pmatrix}$$

پایه: پایه ی استقرا برای V_1 همان فرضی است که در سوال داده شده. می دانیم:

$$V_1 = \mathbb{E}_{a \sim \pi'}[Q^\pi(s, a)] \geq \mathbb{E}_{a \sim \pi}[Q^\pi(s, a)] = V^\pi(s)$$

گام استقرا: فرض کنید بدانیم V_n خواسته‌ی مسئله را برآورده می‌کند یا به عبارتی $V_n(s) \geq V^\pi(s)$. اثبات می‌کنیم $V_{n+1}(s)$ نیز حداقل به اندازه‌ی $V^\pi(s)$ است. برای این منظور، اگر از s شروع کرده و با دنبال کردن π' برای n قدم، به ترتیب پاداش‌های R_1 تا R_n را دریافت کرده و حالت‌های s_1 تا s_n را دیده باشیم، (مشخصاً R_i ها و s_i ها متغیرهایی تصادفی هستند که از توزیع T و R و فرض این‌که از سیاست π' استفاده کرده‌ایم پیروی می‌کنند) آن‌گاه برای V_n داریم:

$$V_n(s) = \mathbb{E}_{\pi'} [R_1 + \gamma R_2 + \dots + \gamma^{n-1} R_n + \gamma^n Q^\pi(s_n, \pi(s_n))] \geq V^\pi(s) \quad (1)$$

هم‌چنین برای V_{n+1} نیز داریم:

$$V_{n+1}(s) = \mathbb{E}_{\pi'} [R_1 + \gamma R_2 + \dots + \gamma^{n-1} R_n + \gamma^n \mathbb{E}_{a \sim \pi'} [Q^\pi(s_n, a)]] \quad (2)$$

از طرفی چون سیاست‌های π در روش ϵ -greedy خود سیاست‌هایی تصادفی هستند، داریم:

$$Q^\pi(s, \pi(s)) = \mathbb{E}_{a \sim \pi} [Q^\pi(s, a)] \quad (*)$$

هم‌چنین با توجه به فرض مسئله، داریم:

$$\begin{aligned} \forall s : E_{a \sim \pi'} [Q^\pi(s, a)] &\geq E_{a \sim \pi} [Q^\pi(s, a)] \\ \implies E_{a \sim \pi'} [Q^\pi(s_n, a)] &\geq E_{a \sim \pi} [Q^\pi(s_n, a)] \quad (**) \end{aligned}$$

در نتیجه با جایگذاری رابطه‌ی (*) در (1) و استفاده از (**) در مقایسه‌ی روابط (1) و (2) به راحتی می‌توان نتیجه گرفت که: $V_{n+1}(s) \geq V_n(s) \geq V^\pi(s)$.

بنابراین مراحل استقرا تکمیل شده و حکم اثبات می‌شود. از آنجایی که به ازای هر n داریم $V_n(s) \geq V^\pi(s)$ ، اگر تمام قدم‌ها را نیز از π' پیروی کنیم باز هم نابرابری برقرار و $V^{\pi'}(s) \geq V^\pi(s)$ خواهد بود.

The END!