

به نام خدا



یادگیری تقویتی

تمرین تئوری اول

یلدا شعبان زاده

۹۸۱۰۱۸۲۲

نیمسال دوم ۰۲-۰۱

فهرست

سوال ۱ ۳
سوال ۲ ۸
سوال ۳ ۱۱
سوال ۴ ۱۵

سوال ۱

(آ) می‌دانیم value iteration به فرم زیر است:

$$V_k^*(s) = \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_{k-1}^*(s'))$$

همچنین فرض کرده‌ایم که:

$$0 \leq R(s, a, s') \leq R_{max} \quad \forall s, a$$

و می‌دانیم:

$$0 \leq P(s'|s, a) \leq 1 \quad \forall s, a$$

حال می‌توان گفت:

$$\begin{aligned} V_k^*(s) &= \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_{k-1}^*(s')) \leq \sum_{s', a} P(s'|s, a) (R(s, a, s') + \gamma V_{k-1}^*(s')) \\ &\leq \sum_{s', a} R(s, a, s') + \gamma V_{k-1}^*(s') \leq R_{max} + \gamma \sum_{s'', a} R(s', a, s'') + \gamma V_{k-2}^*(s'') \\ &\leq R_{max} + \gamma R_{max} + \gamma^2 R_{max} + \dots + \gamma^{k-1} R_{max} \leq \frac{R_{max} \times (1 - \gamma^k)}{1 - \gamma} \leq \frac{R_{max}}{1 - \gamma} \end{aligned}$$

بنابراین به یک upper bound رسیدیم.

(ب) برای اینکه ثابت کنیم V_k^* نسبت به k صعودی است؛ ابتدا V_{k+1}^π را پیدا می‌کنیم به طوری که $V_{k+1}^\pi \geq V_k^*$ باشد.

برای این کار ابتدا از bellman optimality equation استفاده می‌کنیم.

$$V_{k+1}^*(s) = \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_k^*(s'))$$

از آنجا که policy بهینه همان π^* است، می‌توان گفت که:

$$V_{k+1}^* \geq V_{k+1}^\pi \quad (1)$$

حال، فرض کنید π یک policy دلخواه باشد که k تا step با π^* policy برویم و در گام $k+1$ π policy را انتخاب کرده و با آن گام برداریم. از آنجا که rewardها نامنفی هستند داریم:

$$V_{k+1}^\pi \geq V_k^* \quad (2)$$

پس می‌توان گفت که:

$$(1,2) \Rightarrow V_{k+1}^* \geq V_k^*$$

حال، برای همگرایی می‌توان گفت از آنجایی که V_k^* صعودی است و ما یک upper bound برایش بدست آوردیم، داریم:

$$V_0^* \leq \dots \leq V_k^* \leq V_{k+1}^* \leq \dots \leq \frac{R_{max}}{1-\gamma}$$

$$\Rightarrow \lim_{k \rightarrow \infty} V_k^* = \frac{R_{max}}{1-\gamma}$$

از آنجایی که این حد محدود است؛ الگوریتم ما همگرا خواهد بود.

(ج) طبق معادله بلمن داریم:

$$V_k^*(s) = \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_{k-1}^*(s'))$$

$$\Rightarrow \lim_{k \rightarrow \infty} V_k^*(s) = \lim_{k \rightarrow \infty} \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_{k-1}^*(s'))$$

$$= \max_a \sum_{s'} P(s'|s, a) \left(R(s, a, s') + \gamma \lim_{k \rightarrow \infty} V_{k-1}^*(s') \right) \quad \text{پیوستگی}$$

$$\lim_{k \rightarrow \infty} V_k^*(s) = V^*(s) \xrightarrow{k=k-1=\infty} V^*(s) = \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V^*(s))$$

پس در هر مرحله V افزایش پیدا می‌کند تا هنگامی که در بی‌نهایت convergence رخ می‌دهد و مقدارش برای آن state عوض نمی‌شود. همچنین در حالت convergence بهینه است زیرا بنا بر تعریف در حالت convergence داریم $\pi_{k+1}(s) = \pi_k(s)$ یعنی:

$$V^{\pi_k}(s) = \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V^{\pi_k}(s'))$$

از آنجایی که V^{π_k} در معادله بلمن صدق می‌کند، می‌توان گفت V^{π_k} برابر مقدار value function آپتیمال یا V^* است.

(د) در MDP جدید فرض نامنفی بودن پاداش‌ها را نداریم. اما می‌توان با اضافه کردن اندازه کمترین پاداش به همه پاداش‌ها باند reward ها را نامنفی کرد. یعنی داریم:

$$0 \leq R'(s, a, s') \leq R_{max} + |R_{min}| \quad \forall s, a$$

حال طبق معادله بلمن داریم:

$$\begin{aligned}
V_k'^*(s) &= \max_a \sum_{s'} P(s'|s, a) (R'(s, a, s') + \gamma V_{k-1}'^*(s')) \\
&= \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + |R_{min}| + \gamma V_{k-1}'^*(s')) \\
&= \max_a \left(\sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_{k-1}'^*(s')) + \sum_{s'} P(s'|s, a) |R_{min}| \right) \\
&= \max_a \left(\left(\sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_{k-1}'^*(s')) \right) + |R_{min}| \times 1 \right) \\
&= \max_a \left(\sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_{k-1}'^*(s')) \right) + |R_{min}|
\end{aligned}$$

از آنجایی که داریم $V_0'^*(s) = V_0^*(s)$ و $V_1^*(s) = \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_0^*(s'))$ می توان گفت:

$$\begin{aligned}
V_1'^*(s) &= \max_a \left(\sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_0^*(s')) \right) + |R_{min}| = V_1^*(s) + |R_{min}| \\
V_2'^*(s) &= \max_a \left(\sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_1'^*(s')) \right) + |R_{min}| \\
&= \max_a \left(\sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma (V_1^*(s) + |R_{min}|)) \right) + |R_{min}| \\
&= \max_a \left(\sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_1^*(s)) + \gamma |R_{min}| \sum_{s'} P(s'|s, a) \right) + |R_{min}| \\
&= \max_a \left(\sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_1^*(s)) \right) + \gamma |R_{min}| + |R_{min}|
\end{aligned}$$

به همین صورت می توان به صورت افزایشی نشان داد که داریم:

$$\begin{aligned}
V_k'^*(s) &= \max_a \sum_{s'} P(s'|s, a) (R'(s, a, s') + \gamma V_{k-1}'^*(s')) \\
&= \max_a \left(\sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_{k-1}^*(s)) \right) + \sum_{i=0}^{k-1} \gamma^i |R_{min}| \\
&= V_{k-1}^*(s) + \sum_{i=0}^{k-1} \gamma^i |R_{min}| = V_{k-1}^*(s) + \sum_{i=0}^{k-1} \gamma^i r_0
\end{aligned}$$

در حالت convergence نیز با شرط نداشتن terminating state داریم:

$$V'^*(s) = \lim_{k \rightarrow \infty} V'_k(s) = V^*(s) + \frac{r_0}{1 - \gamma}$$

Policy بهینه قبلی در convergence برابر بود با: $\pi^*(s) = \arg \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V^*(s'))$

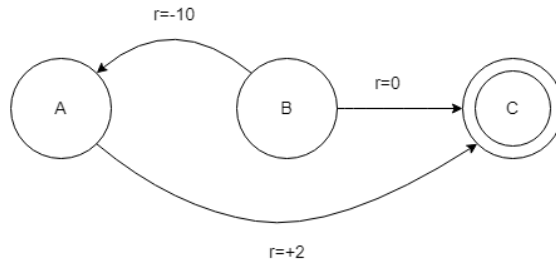
در این حالت نیز داریم:

$$\begin{aligned} \pi'^*(s) &= \arg \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V'^*(s')) \\ &= \arg \max_a \sum_{s'} P(s'|s, a) \left(R(s, a, s') + \gamma \left(V^*(s) + \frac{r_0}{1 - \gamma} \right) \right) \\ &= \arg \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V^*(s)) + \frac{\gamma r_0}{1 - \gamma} \\ &= \arg \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V^*(s)) = \pi^*(s) \end{aligned}$$

(ه) شرط نداشتن terminating state به این دلیل لازم است که در عبارت بالا بتوانیم k را به سمت بی‌نهایت میل دهیم و به عبارتی که به k وابسته نیست برسیم تا بتوانیم policy بهینه را بیابیم.

برای مثال نقض می‌توان گفت:

یک MDP را با سه حالت در نظر بگیرید: A و B و C. حالت C حالت پایانی است. حالت B نیز حالت شروع است.



در اینجا با شروع از B، دو اکشن چپ و راست دارد که احتمال هر کدام 0.5 است. در صورت انتخاب راست reward هایش 0 است و در صورت حرکت به سمت چپ به A می‌رود که reward=-10 و سپس تنها اکشنش که حرکت به سمت C باشد را انجام می‌دهد و reward=+2 را می‌گیرد. پس reward هایش برابر 8- است. پس policy بهینه با شروع از B حرکت به راست و دریافت reward 0 است.

$$V^*(A) = \max_a (1 \times (2 + 1 \times V^*(C))) = 2 \Rightarrow \pi^*(A) = right$$

$$V^*(B) = \max_a \left(\frac{1}{2} \times (-10 + 1 \times V^*(A)), \frac{1}{2} \times (0 + 1 \times V^*(C)) \right) = \max_a (-4, 0) = 0 \Rightarrow \pi^*(B) = right$$

$$V^*(C) = 0$$

اما اگر $r_0 = |R_{min}| = 10$ را به همه پاداش‌ها اضافه کنیم، سیاست بهینه تغییر می‌کند. اکنون، اگر از B به سمت چپ حرکت و سپس به C برویم، پاداش 12 میگیریم که با policy بهینه در مرحله قبل متفاوت است.

$$V^*(A) = \max_a \left(1 \times (12 + 1 \times V^*(C)) \right) = 12 \Rightarrow \pi^*(A) = right$$

$$V^*(B) = \max_a \left(\frac{1}{2} \times (0 + 1 \times V^*(A)), \frac{1}{2} \times (10 + 1 \times V^*(C)) \right) = \max_a(6, 5) = 6 \Rightarrow \pi^*(B) = left$$

$$V^*(C) = 0$$

بنابراین policy بهینه در B با policy بهینه در مرحله قبل تفاوت دارد.

سوال ۲

(آ) فرض کنید $V^{\pi_t}(s)$ و $V^{\pi_{t+1}}(s)$ به ترتیب توابع مقدار حالت برای سیاست های π_t و π_{t+1} باشند و فرض کنید π^* پالیسی بهینه باشد. فرض شده که $V_{\infty}^{\pi_{t+1}}(s) = V^{\pi_{t+1}}(s) = V_{\infty}^{\pi_t}(s) = V^{\pi_t}(s)$ حال اثبات می کنیم $\pi_{t+1} = \pi^*$

برای تناقض فرض کنید π_{t+1} برابر π^* نیست. سپس یک حالت s وجود دارد که $\pi_{t+1}(s) \neq \pi^*(s)$

همچنین داریم:

$$Q^{\pi_t}(s, \pi^*(s)) = R(s, \pi^*(s)) + \gamma \sum_{s'} P(s'|s, \pi^t(s)) V^{\pi_t}(s')$$

$$Q^{\pi_{t+1}}(s, \pi^{(t+1)}(s)) = R(s, \pi^{(t+1)}(s)) + \gamma \sum_{s'} P(s'|s, \pi^{(t+1)}(s)) V^{\pi_{t+1}}(s')$$

با فرض، $V_{\infty}^{\pi_{t+1}}(s) = V_{\infty}^{\pi_t}(s)$ بنابراین دو مقدار Q برابر هستند. از آنجایی که π_{t+1} و π_t برابر نیستند، داریم:

$$R(s, \pi^*(s)) < R(s, \pi^{t+1}(s))$$

بنابراین، می توانیم معادله فوق را به صورت زیر بازنویسی کنیم:

$$V^{\pi_t}(s) + \varepsilon < V^{\pi_{t+1}}(s)$$

بنابراین به تناقض رسیدیم و نتیجه می گیریم اگر تابع ارزش-حالت برای یک policy بین دو تکرار تغییر نکند، به پالیسی بهینه رسیده ایم.

(ب) می توان تعداد مراحل policy evaluation را حداکثر برابر با $|A|^{|S|}$ در نظر گرفت که در آن $|S|$ تعداد حالت ها و $|A|$ تعداد actionها است. زیرا در نظر بگیرید که هر مرحله ارزیابی تابع مقدار را برای هر حالت با استفاده از معادله بلمن محاسبه می کند:

$$V_{k+1}^{\pi}(s) = \sum_{s'} P(s'|s, \pi(s)) (R(s, \pi(s), s') + \gamma V_k^{\pi}(s'))$$

از آنجا که π از فضای stateها به actionها است و صعودی است می توان گفت که $|A|^{|S|}$ policy داریم پس این الگوریتم تمام می شود. زیرا با توجه به قسمت بعد نیز در صورت برابر شدن دو V به پالیسی بهینه رسیده ایم و در لوپ نمی افتیم.

برای اثبات اینکه policy iteration در نهایت خاتمه می یابد و به مقدار بهینه همگرا می شود، می توانیم از استدلال زیر استفاده کنیم:

policy iteration شامل دو مرحله است: policy evaluation و policy improvement. در مرحله policy evaluation. value function را تا زمانی که به تابع ارزش واقعی policy فعلی همگرا شود، به روز می کنیم. در مرحله policy improvement. policy جدیدی را انتخاب می کنیم که با توجه به value function فعلی انتخاب شده است.

از آنجایی که value function تضمین شده است که در تعداد محدودی از تکرارها به تابع ارزش واقعی برای policy فعلی همگرا شود، مرحله ارزیابی policy در نهایت خاتمه می یابد. علاوه بر این، سیاست جدید انتخاب شده در مرحله بهبود policy تضمین شده است که بهتر یا برابر با policy فعلی باشد. اگر policy جدید کاملاً بهتر از policy فعلی باشد، مرحله ارزیابی policy را با policy جدید تکرار می کنیم تا زمانی که تابع ارزش به تابع ارزش واقعی policy جدید همگرا شود.

از آنجایی که فقط تعداد محدودی از policy ها وجود دارد، و هر policy را می توان در تعداد محدودی از مراحل بهبود بخشید، الگوریتم policy iteration تضمین شده است که خاتمه یابد و به policy بهینه همگرا شود.

(ج) دو الگوریتم پالیسی ایتريشن (Policy Iteration) و ولیو ایتريشن (Value Iteration) جزو الگوریتم های متداول برای حل مسئله فرایند تصمیم گیری مارکوف (MDP) به منظور یافتن سیاست بهینه هستند. این الگوریتم ها دارای مزایای مختلفی هستند. در ادامه به بررسی مزایای الگوریتم پالیسی ایتريشن نسبت به الگوریتم ولیو ایتريشن پرداخته می شود:

۱- همگرایی سریعتر: الگوریتم پالیسی ایتريشن معمولاً سریعتر همگرا می شود. زیرا در هر مرحله بین مراحل ارزیابی سیاست و بهبود سیاست جابجا شده می شود. در مقابل، الگوریتم ولیو ایتريشن در هر مرحله ارزیابی سیاست را انجام می دهد که ممکن است زمان بر باشد.

۲- همگرایی تضمین شده: الگوریتم پالیسی ایتريشن تضمین می کند که در تعداد محدودی از مراحل به سیاست بهینه دست پیدا می کند، در حالی که الگوریتم ولیو ایتريشن تنها همگرایی بهینه را در حالت نامحدود تضمین می کند. این به این معناست که الگوریتم پالیسی ایتريشن ممکن است به صورت دقیق تر و سریع تر به راه حل MDP برسد.

۳- عملکرد بهتر در MDP های بزرگ: الگوریتم پالیسی ایتريشن برای MDP های بزرگ کارایی محاسباتی بیشتری دارد. چرا که به تعداد کمتری از مراحل نسبت به الگوریتم ولیو ایتريشن نیاز دارد. این به این معناست که الگوریتم پالیسی ایتريشن به صورت مستقیم سیاست را به به جای به روزرسانی مقادیر همه حالت ها و اقدامات در هر تکرار، پالیسی را مستقیماً به روزرسانی می کند.

۴- اجرای ساده تر: اجرای پالیسی ایتريشن از نظر مفهومی ساده تر از value iteration است، زیرا شامل دو مرحله اصلی است: policy evaluation و بهبود policy. در مقابل، ولیو ایتريشن فقط یک مرحله را شامل می شود، اما می تواند پیچیده تر باشد زیرا به محاسبه تابع مقدار بهینه نیاز دارد.

به طور کلی، policy iteration می تواند انتخاب بهتری نسبت به value iteration از نظر سرعت همگرایی، کارایی محاسباتی و سهولت اجرا، به ویژه برای MDP های بزرگتر باشد. با این حال، انتخاب بین دو الگوریتم در نهایت به ویژگی های خاص مسئله MDP در حال حل بستگی دارد.

(د) اگر داشته باشیم:

$$V_0^{\pi_{t+1}}(s) = \sum_{s'} P(s' | \pi_{t+1}(s), s) [R(s', \pi_{t+1}(s), s) + \gamma V_{\infty}^{\pi_t}(s')]$$

ابتدا ثابت می کنیم که:

$$\forall s \quad V_0^{\pi_{t+1}}(s) \geq V_{\infty}^{\pi_t}(s)$$

برای اینکار می توان گفت که با توجه به policy evaluation داریم:

$$V_{\infty}^{\pi_t}(s) = \sum_{s'} P(s' | \pi_t(s), s) [R(s', \pi_t(s), s) + \gamma V_{\infty}^{\pi_t}(s')]$$

همچنین می‌توان گفت که بهترین policy ما در این لحظه π_{t+1} است. پس می‌توان گفت که:

$$\sum_{s'} P(s' | \pi_{t+1}(s), s) [R(s', \pi_{t+1}(s), s) + \gamma V_{\infty}^{\pi_t}(s')] \geq \sum_{s'} P(s' | \pi_t(s), s) [R(s', \pi_t(s), s) + \gamma V_{\infty}^{\pi_t}(s')] \\ \Rightarrow V_0^{\pi_{t+1}}(s) \geq V_{\infty}^{\pi_t}(s) \quad \forall s \quad (1)$$

حال فرض می‌کنیم مقادیر در policy evaluation صعودی است. یعنی: $\forall s \quad V_{k+1}^{\pi_{t+1}}(s) \geq V_k^{\pi_{t+1}}(s)$

برای اثبات $V_{\infty}^{\pi_{t+1}}(s) \geq V_{\infty}^{\pi_t}(s)$ از آنجا که policy evaluation صعودی است می‌توان گفت:

$$V_0^{\pi_{t+1}}(s) \leq \dots \leq V_k^{\pi_{t+1}}(s) \leq V_{k+1}^{\pi_{t+1}}(s) \leq \dots \leq V_{\infty}^{\pi_{t+1}}(s) \quad (2)$$

پس داریم:

$$(1,2) \Rightarrow V_{\infty}^{\pi_{t+1}}(s) \geq V_0^{\pi_{t+1}}(s) \geq V_{\infty}^{\pi_t}(s) \quad \forall s$$

پس در حالت تغییر یافته این عبارت برقرار است.

برای حالت تغییر نیافته نیز داریم:

$$V_0^{\pi_{t+1}}(s) = V_0^{\pi_t}(s) = 0$$

$$V_{\infty}^{\pi_t}(s) = \sum_{s'} P(s' | \pi_t(s), s) [R(s', \pi_t(s), s) + \gamma V_{\infty}^{\pi_t}(s')] \\ V_{\infty}^{\pi_{t+1}}(s) = \sum_{s'} P(s' | \pi_{t+1}(s), s) [R(s', \pi_{t+1}(s), s) + \gamma V_{\infty}^{\pi_{t+1}}(s')]$$

اکنون، برای اثبات اینکه در هر تکرار از policy iteration، مقدار همگرایی با مقدار اولیه صفر صعودی است، باید نشان دهیم که تابع مقدار برای هر حالت با هر تکرار ارزیابی پالیسی در حال افزایش است.

با فرض اینکه تابع مقدار برای هر حالت به صفر مقداردهی شود، می‌توانیم نشان دهیم که تابع مقدار برای هر حالت با هر تکرار ارزیابی سیاست افزایش می‌یابد. زیرا معادله بلمن یک معادله بازگشتی است که به تابع مقدار حالت بعدی بستگی دارد. از آنجایی که تابع مقدار برای هر حالت به صفر مقداردهی می‌شود، تابع مقدار برای هر حالت تنها با هر تکرار ارزیابی سیاست افزایش می‌یابد.

هنگامی که تابع مقدار همگرا شد، مرحله بهبود پالیسی، پالیسی را به روز می‌کند تا عملی را انتخاب کند که تابع مقدار همگرا را برای هر حالت به حداکثر می‌رساند. این به این معنی است که سیاست جدید بهتر از سیاست قبلی خواهد بود، که تضمین می‌کند که ارزش همگرایی صعودی است.

بنابراین، می‌توان نتیجه گرفت که در هر تکرار از تکرار سیاست، مقدار همگرایی با مقدار اولیه صفر صعودی است.

سوال ۳

(آ) بر اساس سوال داریم:

$$R = E \left[\sum_{t=0}^H r_t \right]$$

$$Gini(P) = \sum_k p_k(1 - p_k)$$

حال اگر در نظر بگیریم action ۲ داریم خواهیم داشت:

$$Gini(P) = p_1(1 - p_1) + p_2(1 - p_2)$$

$$s.t. \quad p_1 + p_2 = 1, \quad p_1, p_2 \geq 0$$

$$\Rightarrow Gini(P) = p_1 p_2 + p_2 p_1 = 2p_1 p_2 = 2p_1(1 - p_1)$$

حال داریم:

$$\arg \max Gini(P) = \arg \max 2p_1(1 - p_1) = \arg \max p_1 - p_1^2$$

$$\Rightarrow 1 - 2p_1^* = 0 \Rightarrow p_1^* = \frac{1}{2} \Rightarrow p_2^* = \frac{1}{2}$$

برای کمینه کردن اختلاف می توان گفت:

$$\arg \min |p_1 - p_2| = \arg \min |p_1 - (1 - p_1)| = \arg \min |1 - 2p_1|$$

$$\Rightarrow 1 - 2p_1^* = 0 \Rightarrow p_1^* = \frac{1}{2} \Rightarrow p_2^* = \frac{1}{2}$$

از دو عبارت بالا می توان نتیجه گرفت که:

$$\arg \max Gini(P) = \arg \min |p_1 - p_2|$$

بنابراین این توزیع به کم شدن تفاوت احتمال انتخاب شدن action ها کمک می کند.

(ب) می خواهیم با استفاده از KKT تابع لاگرانژ مربوط به بهینه سازی عبارت زیر را بنویسیم: (با توجه به فرم کلی قیدها را به این صورت تغییر می دهیم: منبع) همچنین از آنجایی که Gini به فرم concave است و اکسپکتد یک عدد است؛ می توان در نظر گرفت تابع f به صورت convex است. (در این سوال a و a نشان دهنده ی ایندکس آام و action آام هستند و A همان مجموعه G است.)

$$\min_{\pi_A} f(\pi_A) = -E_{\pi_A}[r(a)] - \beta Gini(\pi_A)$$

$$subject \ to \quad l(\pi_A) = \sum_a \pi_A(a) - 1 = 0$$

$$h_i(\pi_A) = -\pi_A(a) \leq 0 \quad \text{for all } i = 0, 1, \dots, |A| - 1$$

حال مسئله زیر را به عنوان فرم dual مسئله بالا در نظر می گیریم:

$$\begin{aligned} \max_{u,v} \min_{\pi_A} & f(\pi_A) + \sum_i u_i h_i(\pi_A) + v l(\pi_A) \\ \text{subject to} & u \geq 0 \end{aligned}$$

تابع لاگرانژ به صورت زیر است که در اینجا u, v ضرایب لاگرانژ هستند:

$$\begin{aligned} \mathcal{L}(\pi_A, u, v) &= f(\pi_A) + \sum_i u_i h_i(\pi_A) + v l(\pi_A) \\ &= -E_{\pi_A}[r(a)] - \beta Gini(\pi_A) + \sum_{i,a} -u_i \pi_A(a) + v \left(\sum_a \pi_A(a) - 1 \right) \end{aligned}$$

(ج) سپس می توان شرایط KKT را برای این مسئله بیان کنیم:

Stationarity Condition:

$$0 \in \partial \left(f(\pi_A) + \sum_i u_i h_i(\pi_A) + v l(\pi_A) \right)$$

Primal feasibility:

$$\begin{aligned} h_i(\pi_A) &\leq 0 \quad \text{for all } i = 0, 1, \dots, |A| - 1 \\ l(\pi_A) &= \sum_a \pi_A(a) - 1 = 0 \end{aligned}$$

Dual feasibility:

$$u_i \geq 0 \quad \text{for all } i = 0, 1, \dots, |A| - 1$$

Complementary slackness:

$$u_i h_i(\pi_A) = 0, \text{ for all } i = 0, 1, \dots, |A| - 1$$

حال می توان گفت:

$$\begin{aligned} \frac{\partial \mathcal{L}(\pi_A, u, v)}{\partial \pi_A} &= -r(a) - \beta(1 - 2\pi_A^*(a)) - u_i^* + v^* = 0 \\ \Rightarrow \pi_A^*(a) &= \frac{r(a) + \beta + u_i^* - v^*}{2\beta} \end{aligned}$$

همچنین داریم:

$$\begin{aligned}
l(\pi_A) &= \sum_a \pi_A(a) - 1 = 0 \Rightarrow \sum_a \pi_A^*(a) = 1 \Rightarrow \sum_a \frac{r(a) + \beta + u_i^* - v^*}{2\beta} = 1 \Rightarrow \sum_{a,i} r(a) + \beta + u_i^* - v^* \\
&= 2\beta \Rightarrow |A|\beta - |A|v^* + \sum_{a,i} r(a) + u_i^* = 2\beta \Rightarrow |A|v^* - \sum_i u_i^* \\
&= \sum_a r(a) + \beta(|A| - 2)
\end{aligned}$$

از آنجایی که احتمالات غیر صفر هستند، می توان گفت که $h_i(\pi_A) < 0$. همچنین چون در شروط KKT داریم که $u_i h_i(\pi_A) = 0$ می توان گفت که:

$$\begin{aligned}
u_i &= 0 \quad \text{for all } i = 0, 1, \dots, |A| - 1 \\
\Rightarrow v^* &= \frac{(\sum_a r(a)) + \beta(|A| - 2)}{|A|}
\end{aligned}$$

پس داریم:

$$\begin{aligned}
\pi_A^*(a) &= \frac{r(a) + \beta + u_i^* - v^*}{2\beta} = \frac{r(a)|A| + \beta|A| - (\sum_a r(a)) - \beta(|A| - 2)}{2\beta|A|} \\
&= \frac{r(a)|A| - (\sum_a r(a))}{2\beta|A|} + \frac{1}{|A|}
\end{aligned}$$

(د) ب رای تبدیل مسئله داده شده به یک مسئله فرم درجه دوم بدون داشتن A ، می توانیم از روش زیر استفاده کنیم:

ابتدا Gini را بر حسب مقدار مورد انتظار توزیع احتمال مجذور، به صورت زیر بیان می کنیم:

$$Gini(\pi) = \left(\pi(a) - \frac{1}{|A|} \right)^2$$

این یک نتیجه استاندارد در تئوری احتمال است که نشان می دهد معیار ناخالصی جینی معادل واریانس توزیع احتمال است، زمانی که میانگین ثابت شود. حال مجموعه جدیدی از متغیرها را مشخص می کنیم که با y_a نشان داده می شوند، به طوری که $\pi(a) = \frac{\exp(y_a)}{Z}$ که در آن Z یک normalization constant است به طوری که $\sum_a \exp(y_a) = 1$. این تبدیل تضمین می کند که جمع احتمالات ۱ شوند، یعنی $\sum_a \pi(a) = 1$.

مسئله را بر حسب y_a با استفاده از عبارت بالا برای $\pi(a)$ بازنویسی می کنیم:

$$\begin{aligned}
\min_y f(y) &= -E[r(a)] - \beta \sum_a \left(\frac{\exp(y_a)}{Z} - \frac{1}{|A|} \right)^2 \\
\text{subject to } l(y) &= \sum_a \exp(y_a) - Z = 0 \\
h_{i(y)} &= -y_a \leq 0 \quad \text{for all } i = 0, 1, \dots, |A| - 1
\end{aligned}$$

مسئله به دست آمده یک مسئله فرم درجه دوم با شکل زیر است:

$$\begin{aligned} \text{minimize: } Q(y) &= -E[r(a)] - \beta \sum_a \left(\frac{\exp(y_a)}{Z} - \frac{1}{|A|} \right)^2 \\ \text{subject to: } \sum_a \exp(y_a) - Z &= 0 \end{aligned}$$

برای یافتن راه حل بهینه برای مجموعه اقدام A ، می توانیم با استفاده از تکنیک های بهینه سازی استاندارد، مانند gradient descent یا روش نیوتن، مسئله شکل درجه دوم بالا را حل کنیم. هنگامی که جواب بهینه y را به دست آوردیم، می توانیم توزیع احتمال بهینه π را با استفاده از فرمول بازیابی کنیم:

$$\pi(a) = \frac{\exp(y_a)}{Z}$$

در نهایت، ما می توانیم مجموعه اقدامات بهینه A را با انتخاب اقدامات با احتمال غیر صفر بازیابی کنیم، یعنی $A = \{a: \pi(a) > 0\}$.

سوال ۴

(آ) طبق روابط داده شده می‌دانیم:

$$E_{-1}(s) = 0$$

$$E_t(s) = \gamma\lambda E_{t-1}(s) + I_{ss_t}$$

حال می‌خواهیم ثابت کنیم:

$$E_t(s) = \sum_{k=0}^t (\gamma\lambda)^{t-k} I_{ss_k}$$

برای اینکار از استقرا استفاده می‌کنیم:

فرض استقرا:

$$E_{-1}(s) = 0$$

$$E_0(s) = \gamma\lambda E_{-1}(s) + I_{ss_0} = I_{ss_0} = \sum_{k=0}^0 (\gamma\lambda)^{0-k} I_{ss_k}$$

و همینطور برای $t=1$ داریم:

$$E_1(s) = \gamma\lambda E_0(s) + I_{ss_1} = \gamma\lambda I_{ss_0} + I_{ss_1} = \sum_{k=0}^1 (\gamma\lambda)^{1-k} I_{ss_k}$$

پس برای فرض اثبات شد.

گام استقرا: فرض می‌کنیم برای $t-1$ فرض برقرار است. حال نشان می‌دهیم برای t نیز برقرار است.

$$E_t(s) = \gamma\lambda E_{t-1}(s) + I_{ss_t} = \gamma\lambda \sum_{k=0}^{t-1} (\gamma\lambda)^{t-1-k} I_{ss_k} + I_{ss_t} = \sum_{k=0}^{t-1} (\gamma\lambda)^{(t-1-k)+1} I_{ss_k} + I_{ss_t}$$

$$= \sum_{k=0}^t (\gamma\lambda)^{t-k} I_{ss_k}$$

بنابراین فرض اثبات شد.

(ب) می‌دانیم:

$$\delta_t = R_{t+1} + \gamma V(s_{t+1}) - V(s_t)$$

$$V(s) \leftarrow V(s) + \alpha \delta_t E_t(s)$$

و همچنین:

$$\begin{aligned}
\sum_{t=0}^{T-1} \Delta V_t^{TD}(s) &= \sum_{t=0}^{T-1} \alpha \delta_t E_t(s) = \sum_{t=0}^{T-1} \alpha \delta_t \sum_{k=0}^t (\gamma \lambda)^{t-k} I_{ss_k} \\
&= \alpha \delta_0 I_{ss_0} \\
&\quad + \alpha \delta_1 (\gamma \lambda) I_{ss_0} + \alpha \delta_1 I_{ss_1} \\
&\quad + \alpha \delta_2 (\gamma \lambda)^2 I_{ss_0} + \alpha \delta_2 (\gamma \lambda) I_{ss_1} + \alpha \delta_2 I_{ss_2} \\
&\quad + \dots \\
&\quad + \alpha \delta_{T-1} (\gamma \lambda)^{T-1} I_{ss_0} + \alpha \delta_{T-1} (\gamma \lambda)^{T-2} I_{ss_1} + \dots + \alpha \delta_{T-1} I_{ss_{T-1}} \\
&= \sum_{t=0}^{T-1} \alpha I_{ss_t} \sum_{i=0}^{T-1-t} (\gamma \lambda)^i \delta_{i+t} = \sum_{i=0}^{T-1} \alpha I_{ss_i} \sum_{k=i}^{T-1} (\gamma \lambda)^{k-i} \delta_k
\end{aligned}$$

تساوی مورد نظر اثبات شد.

(ج) بر اساس نگاه forward داریم:

$$\begin{aligned}
\frac{1}{\alpha} \Delta V_t^\lambda(s) &= \frac{1}{\alpha} \left[\alpha \left(G_t^\lambda - V_t(s_t) \right) \right] = G_t^\lambda - V_t(s_t) \\
G_t^\lambda &= (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)} = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \left[\left(\sum_{k=0}^{n-1} r_{t+1+k} \gamma^k \right) + \gamma^n V_t(s_{t+n}) \right] \\
&= (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \left[\left(\sum_{k=1}^n r_{t+k} \gamma^{k-1} \right) \right] + (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \gamma^n V_t(s_{t+n}) \\
&= (1 - \lambda) \sum_{k=1}^{\infty} \gamma^{k-1} r_{t+k} \left[\left(\sum_{n=k-1}^{\infty} \lambda^n \right) \right] + (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \gamma^n V_t(s_{t+n}) \\
&= (1 - \lambda) \sum_{k=1}^{\infty} \gamma^{k-1} r_{t+k} \frac{\lambda^{k-1}}{1 - \lambda} + (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \gamma^n V_t(s_{t+n}) \\
&= \sum_{k=1}^{\infty} (\gamma \lambda)^{k-1} r_{t+k} + (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \gamma^n V_t(s_{t+n}) \\
&= \sum_{n=1}^{\infty} (\gamma \lambda)^{n-1} r_{t+n} + \sum_{n=1}^{\infty} \lambda^{n-1} \gamma^n V_t(s_{t+n}) - \sum_{n=1}^{\infty} (\gamma \lambda)^n V_t(s_{t+n})
\end{aligned}$$

بنابراین می توان گفت:

$$\begin{aligned}
\frac{1}{\alpha} \Delta V_t^\lambda(s) &= \sum_{n=1}^{\infty} (\gamma\lambda)^{n-1} r_{t+n} + \sum_{n=1}^{\infty} \lambda^{n-1} \gamma^n V_t(s_{t+n}) - \sum_{n=1}^{\infty} (\gamma\lambda)^n V_t(s_{t+n}) - V_t(s_t) \\
&= \sum_{n=1}^{\infty} (\gamma\lambda)^{n-1} r_{t+n} + \sum_{n=1}^{\infty} \lambda^{n-1} \gamma^n V_t(s_{t+n}) - \sum_{n=0}^{\infty} (\gamma\lambda)^n V_t(s_{t+n}) = \\
&= \sum_{n=1}^{\infty} (\gamma\lambda)^{n-1} r_{t+n} + \sum_{n=1}^{\infty} \lambda^{n-1} \gamma^n V_t(s_{t+n}) - \sum_{n=0}^{\infty} (\gamma\lambda)^n V_t(s_{t+n}) \\
&= \sum_{n=1}^{\infty} (\gamma\lambda)^{n-1} (r_{t+n} + \gamma V_t(s_{t+n})) - \sum_{n=1}^{\infty} (\gamma\lambda)^{n-1} V_t(s_{t+n-1})
\end{aligned}$$

پس داریم:

$$\begin{aligned}
\frac{1}{\alpha} \Delta V_t^\lambda(s) &= \sum_{n=1}^{\infty} (\gamma\lambda)^{n-1} (r_{t+n} + \gamma V_t(s_{t+n}) - V_t(s_{t+n-1})) \\
&= \sum_{n=0}^{\infty} (\gamma\lambda)^n (r_{t+n+1} + \gamma V_t(s_{t+n+1}) - V_t(s_{t+n})) \\
&= \sum_{k=t}^{\infty} (\gamma\lambda)^{k-t} (r_{k+1} + \gamma V_t(s_{k+1}) - V_t(s_k))
\end{aligned}$$

تساوی مورد نظر اثبات شد.

(د)

$$\begin{aligned}
\frac{1}{\alpha} \Delta V_t^\lambda(s) &= \sum_{k=t}^{\infty} (\gamma\lambda)^{k-t} (r_{k+1} + \gamma V_t(s_{k+1}) - V_t(s_k)) \Rightarrow \Delta V_t^\lambda(s) \\
&= \alpha \sum_{k=t}^{\infty} (\gamma\lambda)^{k-t} (r_{k+1} + \gamma V_t(s_{k+1}) - V_t(s_k)) \xrightarrow{\text{offline}} \sum_{t=0}^{T-1} \Delta V_t^\lambda(s) \\
&= \sum_{t=0}^{T-1} \alpha \sum_{k=t}^{\infty} (\gamma\lambda)^{k-t} \delta_t \stackrel{2}{\Rightarrow} \sum_{t=0}^{T-1} \Delta V_t^{TD}(s) = \sum_{t=0}^{T-1} \alpha I_{ss_t} \sum_{k=t}^{T-1} (\gamma\lambda)^{k-t} \delta_t = \sum_{t=0}^{T-1} \Delta V_t^\lambda(s) I_{ss_t}
\end{aligned}$$

در حالت **offline** به روزرسانی ها روی هم انباشته و در انتهای **episode** اعمال می شوند اما در حالت آنلاین پس از هر **step** در هر **episode** اعمال می شود. حالت تساوی در حالت **offline** دقیق است، زیرا در این حالت V_t تغییر نمی کند ولی در حالت آنلاین تغییر می کند. همه مراحل پس از **terminal state** دارای پاداش صفر و ارزش صفر هستند. بنابراین همه آنها نیز صفر هستند. (۲)