

۱.۱.۱ با استقرا روی k ثابت می کنیم $V_k^* \leq R_{\max} (1 + \gamma + \dots + \gamma^{k-1})$ برای هر $k \geq 1$.

باید استقرا: می دانیم $V_1^* = \max_{\alpha} \sum_{s'} T(s, \alpha, s') [R(s, \alpha, s') + \gamma V_0^*(s')] = 0$ و $R \leq R_{\max}$ بنابراین:

$$\forall s: V_1^*(s) \leq \max_{\alpha} \sum_{s'} T(s, \alpha, s') \times R_{\max} = \max_{\alpha} R_{\max} = R_{\max}$$

پس برای $k=1$ حکم ثابت شد.

گام استقرا: فرض کنید برای $k-1$ برقرار باشد:

$$V_{k-1}^*(s') \leq R_{\max} (1 + \gamma + \dots + \gamma^{k-2})$$

$$V_k^*(s) = \max_{\alpha} \sum_{s'} T(s, \alpha, s') [R(s, \alpha, s') + \gamma V_{k-1}^*(s')] \leq \max_{\alpha} \sum_{s'} T(s, \alpha, s') [R_{\max} + \gamma V_{k-1}^*(s')]$$

$$\leq \max_{\alpha} \sum_{s'} T(s, \alpha, s') [R_{\max} + \gamma R_{\max} (1 + \gamma + \dots + \gamma^{k-2})] = R_{\max} (1 + \gamma + \dots + \gamma^{k-1})$$

پس ادعا همان ثابت شد و یک گران بالا برای V_k^* پیدا کردیم.

۱.۱.۲ اثبات راحت تر: طبق تعریف، V_k^* مقدار Optimal Value است در حالتی که $H=k$ باشد.

یعنی کلاً k گام حرکت کنیم. منطقه در هر گام Reward حداکثر R_{\max} است پس هر سارگی می تواند نتیجه گرفت:

$$V_k^*(s) = \max_{\alpha_1, \dots, \alpha_{k-1}} E \left[\sum_{t=0}^{k-1} \gamma^t R(s_t, \alpha_t, s_{t+1}) \mid s_0 = s \right] \leq \max_{\alpha} E \left[\sum_{t=0}^{k-1} \gamma^t R_{\max} \right]$$

$$\Rightarrow V_k^*(s) \leq R_{\max} (1 + \gamma + \dots + \gamma^{k-1})$$

۱.۱.۳ فرض کنید π انتخاب کند که ابتدا k گام action بهینه انجام دهد و سپس در گام $k+1$ ریزم حرکت کند (مثلاً آکشن α_k) در این صورت:

$$V_{k+1}^{\pi}(s) = V_k^*(s) + \gamma^k E[R(s_k, \alpha_k)] \geq V_k^*(s)$$

نامساوی ۱ از اینجا نتیجه گرفتیم که $V_{k+1}^{\pi}(s)$ ها نامنمی هستند. پس:

$$V_{k+1}^* \geq V_{k+1}^{\pi} \geq V_k^* \Rightarrow V_k^* \text{ is non-decreasing}$$

ما در Value-iteration در k امین مرحله قرار می دهیم $V^k(s) := \max_{\alpha} \sum_{s'} T(s, \alpha, s') [R(s, \alpha, s') + \gamma V^{k-1}(s')]$ که در واقع V^k همان V_k^* می شود. ثابت کردیم V^k ها صعودی هستند و همچنین طبق بخش قبل، یک گران بالا دارند. بنابراین دنباله متناهی گام می تواند اکیدا صعودی باشد و از جایی به بعد ثابت می ماند و در واقع Converge می کند و در آنجا داریم $V^{n+1}(s) = V^n(s)$ پس $V^{n+1}(s) = \max_{\alpha} \sum_{s'} T(s, \alpha, s') [R(s, \alpha, s') + \gamma V^n(s')]$ یعنی V ها در معادله بلمن (Optimal) صدق می کند.

۱.۱.۳ فرضی کنید $\lim_{k \rightarrow \infty} V^k = V^*$ داریم :

$$\lim_{k \rightarrow \infty} V^k(s) = \lim_{k \rightarrow \infty} \max_a \sum_{s'} P(s, a, s') [R(s, a, s') + \gamma V^k(s')] = V^*(s)$$

$$\Rightarrow V^*(s) = \max_a \sum_{s'} P(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

پس V^* در معادله بهیگی بلین صدق می کند و بنابراین optimal value است.

۲.۱.۵ در هر حرکت یک r_0 به reward دریافت می شود بنابراین داریم :

$$\hat{V}_k^*(s) = V_k^*(s) + r_0 (1 + \gamma + \dots + \gamma^{k-1}) = V_k^*(s) + r_0 \frac{1 - \gamma^k}{1 - \gamma} = V_k^*(s) + C$$

می توان دید به تمام value ها مقدار ثابت C مستقل از s اضافه شده. بنابراین نباید تفاوتی ایجاد شود و بالایی بهینه برای \hat{V}_k^* همان بالایی بهینه برای V_k^* است.

آنگاه اگر r_0 را طوری قرار دهیم که تمام \hat{R} ها نامنتی باشند (مانند قرار دهیم $r_0 = -\min_{s,a} R(s,a)$) طبق آسانی که در بخش ۱.۱ انجام دادیم، نتیجه می شود \hat{V}_k^* ها به مقدار بهینه همگرا می شوند و همانطور که گفتیم برای سیاست بهینه، سیاست برای مسئله اصلی هم بهینه خواهد بود. برای یافتن \hat{V}^* توجه کنید :

$$\hat{V}^*(s) = \max_a E_{s'} [\hat{R}(s, a, s') + \gamma \hat{V}^*(s')]$$

آنگاه مقدار دهد $\hat{V}^* = V^* + C$ پس داریم :

$$\hat{V}^*(s) + C = \max_a E_{s'} [R(s, a, s') + r_0 + \gamma (V^*(s') + C)] \quad \text{I}$$

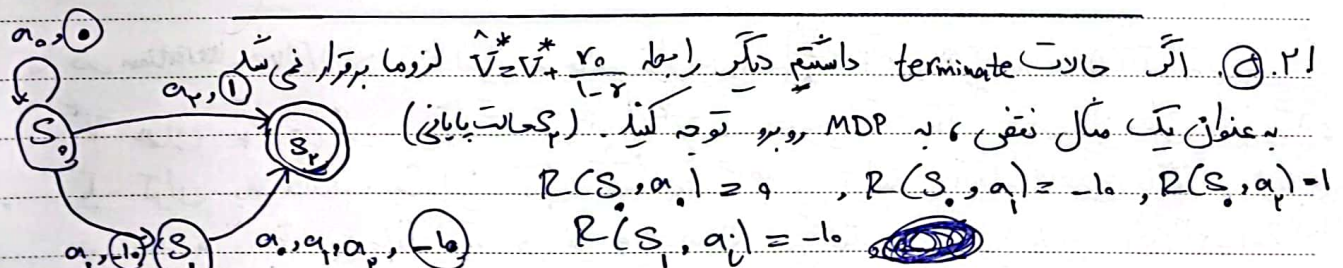
از طرفی برای V^* داریم :

$$V^*(s) = \max_a E_{s'} [R(s, a, s') + \gamma V^*(s')] \quad \text{II}$$

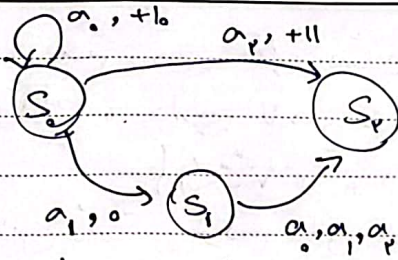
باکم کردن II، I از هم نتیجه می شود :

$$C = r_0 + \gamma C \Rightarrow C = \frac{r_0}{1 - \gamma} \Rightarrow \hat{V}^* = V^* + \frac{r_0}{1 - \gamma}$$

همچنین می شد با استفاده از $\hat{V}_{k0}^* = V_k^* + \frac{r_0}{1 - \gamma}$ و قرار دادن $k \rightarrow \infty$ همین نتیجه را گرفت.



به وضوح در اینجا $V(s_1) = 1$ و بالایی بهینه در شروع انجام می است.



حالا باید مقدار دهیم $\gamma = 1$ و جوایز را زیاد کنیم :

اکنون اینجا اگر a_2 را انجام دهیم جایزه 11 را

دریافت می کنیم و تمام می شود. اما اگر a_1

یا a_3 را انجام دهیم جایزه $(1+\gamma x)$ را دریافت می کنیم یا $x = \frac{1}{1-\gamma}$.

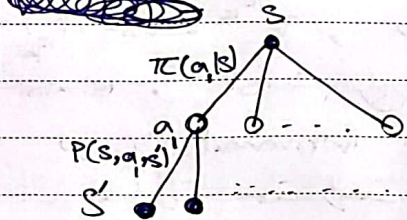
کافیست به جای $-x$ و عدد $x+1$ را قرار می دهیم بطوریکه $x \times \frac{1}{1-\gamma} > x+1$ یا معادلا

$$\frac{1-\gamma}{\gamma} > x \quad \text{یا در حالت جدید آکشن بهینه a باشد و} \quad V^*(S_i) = x \times \frac{1}{1-\gamma}$$

۳.۱ (۶) توجه کنید در هر گام ما π_{k+1} را طوری قرار می دهیم که :

$$\pi_{k+1}(S) = \operatorname{argmax}_a \sum_{S'} P(S, a, S') [R(S, a, S') + \gamma V^{\pi_k}(S')]$$

اکنون به کمک Backup diagram اثبات می کنیم که $V^{\pi_{k+1}}(S) \geq V^{\pi_k}(S)$



حالا فرض می کنیم در لایه اول و به جای $V(S)$ و صرفا

آکشن حریصانه را انتخاب کنیم. با اینکار $V(S)$ قطعا کم نمی شود (بیشتر مساوی می شود).

اکنون در لایه بعدی پالیسی را با حریصانه جایگزین کنیم. $V(S)$ ها با اینکار بهتر می شود. بنابراین مقدار لایه بالاتر (که یک Convex Combination از $V(S)$ ها است) هم کمتر نمی شود.

به همین شکل ادامه می دهیم و تا لایه n جایگزین می کنیم. بنابراین پالیسی حریصانه مقدار Value

$$V^{\pi_{k+1}}(S) \geq V^{\pi_k}(S) \quad \text{را بهتر می کند و ثابت کردیم}$$

در صورتی که برای تمام S ها حالت تساوی رخ دهد و Converge کرده ایم و

~~مقادیر $V^{\pi_k}(S)$ و $V^{\pi_{k+1}}(S)$ برابر می شوند و به همین دلیل $V^{\pi_k}(S) = V^{\pi_{k+1}}(S)$ می شود.~~

در غیر این صورت یعنی یک Value پیدا بهتر شده

(اگر $V^{\pi_{k+1}}(S) = V^{\pi_k}(S)$ یعنی $V^{\pi_{k+1}}(S) > V^{\pi_k}(S)$ ها ثابت خواهند ماند و می توانیم چرخ را پلاس بمانی

۳.۱ (۷) با توجه به اینکه تعداد متاهی سیاست مختلف وجود دارد (حد اکثر $|A|^{|S|}$ بدنام)

و در هر گام سیاست بهتر می شود و قطعا همگرا می شود. اکنون توجه کنید در آن موقع داریم :

$$\pi_{k+1}(S) = \pi_k(S) \rightarrow V^{\pi_{k+1}}(S) = \max_{a \in A} \sum_{S'} P(S, a, S') [R(S, a, S') + \gamma V^{\pi_k}(S')]$$

و یعنی V ها در معادلات بهیگی باهمین صدق می کنند و طبق اثبات کلاس، بهینه هستند.

بنابراین π_k هم سیاست بهینه است.

۱۳.۱ در بخشی های قبلی نشان دادیم که هر دو Value iteration, Policy iteration در هنگام همگرایی، دارای تلاهای هستند که در معادله بهنگی بلمن صدق می کنند.

$$V^*(s) = \max_a \sum_{s'} P(s, a, s') [R(s, a) + \gamma V^*(s')]$$

در کلاس ثابت کردیم این معادلات یک جواب یکتا دارند (با کمک Banach Fixed Point theorem) و بنابراین این دو Value function هر دو یک چیز هستند. سیاست های بهینه ممکن است متفاوت باشند، در صورتی که چند حالت با V یکسان داشته باشیم و همگی \max_a باشند. اما در هر صورت مقدار Value فنی نمی کند.

۱۳.۱ Policy iteration: فاز Evaluation: $V_{\pi_k}(s) := \sum_{s'} P(s, \pi_k(s), s') [R(s, \pi_k(s)) + \gamma V_{\pi_k}(s')]$ را تکراری کنیم تا

ما نتایج کند. هر گام $O(|S|)$ طول می کشد و بطور متوسط به $O\left(\frac{\log(\frac{1}{\epsilon})}{1-\gamma}\right)$ برای دقت ϵ نیاز داریم. $O(|S| \frac{\log(\frac{1}{\epsilon})}{1-\gamma})$ می تواند بسیار زیاد باشد. راه دیگر حل با کمک Matrix inversion است. قرار است V_{π_k} را پیدا کنیم که در $V_{\pi_k} = T V_{\pi_k}$ صدق می کند (نقطه ثابت)، $T V_{\pi_k} = r_{\pi_k} + \gamma P_{\pi_k} V_{\pi_k}$ بطور سر راست داریم، $V_{\pi_k} = (I - \gamma P_{\pi_k})^{-1} r_{\pi_k}$ ، در $O(|S|^3)$ محاسبه می شود (برای inverse گرفتن ماتریس $O(|S|^3)$ ، $|S| \times |S|$ ، $|S|$ هم برای سایر عملیات).

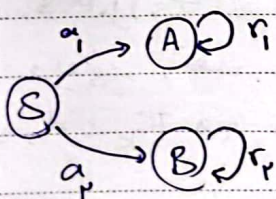
در فاز improvement نیز قرار می دهیم $Q(s, a) := \sum_{s'} P(s, a, s') [R(s, a) + \gamma V_{\pi_k}(s')]$ که محاسبه عبارت سمت راست $O(|S||A|)$ طول می کشد و $|S|$ بار باید اینکار را انجام دهیم (تمام مقادیر s). پس $O(|S|^2|A|)$ طول می کشد. پس هر iteration در کل $O(|S|^2|A| + O(|S|))$ طول می کشد.

Value iteration: در هر گام قرار می دهیم $V(s) := \max_a \sum_{s'} P(s, a, s') [R(s, a) + \gamma V(s')]$ که یعنی برای تمام مقادیر s ، $O(|S||A|)$ طول می کشد.

پس هر iteration از روش Value iteration پیچیدگی محاسباتی کمتری دارد، اما در عوض نیاز به iteration های بیشتری برای همگرایی دارد و در عمل روش Policy iteration سریعتر به پایایی بهینه می رسد.

۱۳.۱ در این حالت، دیگر Bellman Operator یک contraction mapping نیست و اثبات

همگرایی کمتر کلاسی داشتهیم. صحیح نیست. در حقیقت اگر جمع ریواردها کران دار نباشد (یعنی R جمع) هیچ کدام از این دو الگوریتم تضمینی برای همگرایی ندارند.



حتی برای خود ما هم دشوار است در این حالت اکتش بهینه را پیدا کنیم. مثلاً در MDP روبرو، در هر دو حالت Return برابر ∞ خواهد بود و طبیعتاً اگر تمایلی به ریوارد کوتاه مدت نداشته باشیم (یعنی $\gamma = 1$) هیچ کدام را به دیگری

نباید ترجیح دهیم.

$$\|B^\pi V - B^\pi V'\| = \max_s |B^\pi V(s) - B^\pi V'(s)| = \gamma \max_s \left| \sum_{s'} P(s'|s, \pi(s)) [V(s') - V'(s')] \right| \quad (1).2$$

نامساوی مثلث

$$\leq \gamma \max_s \sum_{s'} |P(s'|s, \pi(s)) [V(s') - V'(s')]| = \gamma \max_s \sum_{s'} P(s'|s, \pi(s)) |V(s') - V'(s')|$$

این عبارت یک Convex Combination از $|V(s') - V'(s')|$ هاست و بنابراین می توان نوشت:

$$\gamma \max_s \sum_{s'} P(s'|s, \pi(s)) |V(s') - V'(s')| \leq \gamma \max_s \max_{s'} |V(s') - V'(s')|$$

$$= \gamma \max_s \|V - V'\| = \gamma \|V - V'\|$$

$$\Rightarrow \|B^\pi V - B^\pi V'\| \leq \gamma \|V - V'\|$$

۱.۲ (۲) با توجه به موال قبل ثابت کردیم B^π یک Contraction mapping است و می توانی از قضیه نقطه ثابت باناخ استفاده کرد و گفت B^π دقیقاً یک نقطه ثابت دارد همچنین می توان به سادگی با فرض خلاف ثابت کرد.

فرض کن V_1 و V_2 دو نقطه ثابت متفاو باشند یعنی $B^\pi V_1 = V_1$ و $B^\pi V_2 = V_2$.

طبق نامساوی بخش قبل:

$$\|V_1 - V_2\| = \|B^\pi V_1 - B^\pi V_2\| \leq \gamma \|V_1 - V_2\| \Rightarrow 0 \leq \gamma < 1$$

پس نقطه ثابت B^π یکتا است.

$$B^\pi V'(s) - B^\pi V(s) = \left[r(s, \pi(s)) + \gamma \sum_{s'} P(s'|s, \pi(s)) V'(s') \right] - \left[r(s, \pi(s)) + \gamma \sum_{s'} P(s'|s, \pi(s)) V(s') \right] \quad (3).1.2$$

$$= \gamma \sum_{s'} P(s'|s, \pi(s)) \times [V'(s') - V(s')] \geq 0 \Rightarrow B^\pi V'(s) \geq B^\pi V(s) \quad \forall s \in S$$

* توجه کن این که $V'(s') - V(s') \geq 0$ است مستقیماً از فرض سوال نتیجه می شود.

۲.۲ (۴) فقط به ازای V^* (Optimal Value Function) می تواند ۰ شود:

$$\|BV - V\| = 0 \Leftrightarrow BV - V = 0 \Leftrightarrow BV = V \Leftrightarrow V \text{ is fixed point of } B$$

$$\Leftrightarrow V = V^* \quad (\text{Since the fixed point is unique})$$

۲.۲ (۵) ابتدا طبق نامساوی مثلث داریم: (I)

$$\|B^\pi V - B^\pi V^\pi\| \leq \gamma \|V - V^\pi\| \quad (II)$$

با توجه به آنکه V^π نقطه ثابت B^π است داریم: (III)

$$\textcircled{II}, \textcircled{III} \Rightarrow \|B^\pi V - V^\pi\| \leq \gamma \|V - V^\pi\|$$

حالا اگر کران بست آمده را در نامساوی \textcircled{I} استفاده کنیم، نتیجه می شود:

$$\|V - V^\pi\| \leq \|V - B^\pi V\| + \gamma \|V - V^\pi\| \Rightarrow \|V - V^\pi\| \leq \frac{\|V - B^\pi V\|}{1 - \gamma}$$

این نامساوی درم نیز کاملاً مشابه همین است! مراحل را بدون توضیحات می نویسم:

$$\|BV - BV^*\| \leq \gamma \|V - V^*\| \xrightarrow{BV^* = V^*} \|BV - V^*\| \leq \gamma \|V - V^*\| \quad \textcircled{IV}$$

$$\|V - V^*\| \leq \|V - BV\| + \|BV - V^*\| \stackrel{\textcircled{IV}}{\leq} \|V - BV\| + \gamma \|V - V^*\|$$

$$\Rightarrow \|V - V^*\| \leq \frac{\|V - BV\|}{1 - \gamma}$$

۲.۲ $\textcircled{5}$ توجه کنید اینکه سیاست حریصانه استخراج شده از V است، ما این نتیجه را می دهد که $BV = B^\pi V$. با نگرش کردن به روابط $BV(s)$ ، $B^\pi V(s) = \arg\max_{\pi} [V(s) + \gamma \sum_{s'} P(s'|s, \pi) V(s')]$ این نتیجه به سادگی قابل دریافت است. حالا به کمک این نامساوی معادله سوال را باز نویسی می کنیم:

$$\|V - V^\pi\| \leq \frac{\|V - B^\pi V\|}{1 - \gamma} = \frac{\|V - BV\|}{1 - \gamma} = \frac{\epsilon}{1 - \gamma}$$

$$\|V - V^*\| \leq \frac{\|V - BV\|}{1 - \gamma} = \frac{\epsilon}{1 - \gamma}$$

$$\frac{\epsilon}{1 - \gamma} \geq \|V - V^\pi\| + \|V^* - V\| \geq \|V^* - V^\pi\|$$

توجه کنید می دانیم $V^* \geq V^\pi$ و بنابراین $V^*(s) - V^\pi(s) = |V^*(s) - V^\pi(s)| \leq \|V^* - V^\pi\|$

$$\forall s : V^*(s) - V^\pi(s) \leq \frac{\epsilon}{1 - \gamma} \Rightarrow V^\pi(s) \geq V^*(s) - \frac{\epsilon}{1 - \gamma}$$

۲.۲ \textcircled{V} یک کران پایین روی V^π می تواند تضمین کند که سیاست ما از یک حری بدتر

نیست، و در مسائلی مانند Robotics (با فرض کنید هلیکوپتر پرنده!) که سیاست

خوبی بد می تواند باعث آسیب جزی به تجهیزات ما شود (شلاریات معقول کند و ضربه بخورد)

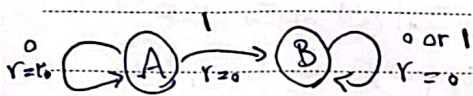
داشتن کران پائین روی $V^\pi(s)$ می تواند یک راه برای تضمین آسیب نپذیرن وسیله باشد.

مثلاً یک threshold تضمین می کنیم و اگر $V^\pi(s)$ از آن کران پائین بهتری داشت می توانیم

اطمینان حاصل کنیم با سیاست فعلی رفتن به حالت خطرناک نیست.

۲.۲ (۱) خیر، لزومی ندارد. فقط می توان نتیجه گرفت که آن یاسین های که برای $V^\pi(s)$ داریم برابر هستند، اما لزومی ندارد برابر باشند.

در ادامه یک مثال ساده می زنیم که ادعای سوال را نقض می کند:



همچنین در این مثال $\delta = 0.9$ در نظر بگیرید.

فرض کنید $V(A) = 2$ ، $V(B) = 0$. در این حالت خواهیم داشت $1.8V - V = 1$

همچنین اگر $V(A) = 0$ ، $V(B) = 1$ ، باز هم خواهیم داشت $1.8V - V = 1$

اما در حالت اول $\pi(A) = 0$ خواهد بود. در حالت دوم $\pi(A) = 1$ و بنابراین خواهیم داشت $V^\pi(A) = 1$ اما $V^\pi(A) = 0$ و برابر نیستند. این مثال هم نشان می دهد لزوماً ادعا برقرار نیست.

۲.۲ (۹) در ابتدا همانطور که در سوال گفته می شود $BV = BV^\pi$ و همچنین نامعادی اول سوال ۵ به شکل زیر در می آید:

$$\|V - V^\pi\| \leq \frac{\varepsilon}{1-\gamma} \quad (I)$$

حالا ادعای کنیم $\|V - V^\pi\| \leq 1.1V^*$

توجه کنید در صورت اثبات این ادعا، نتیجه می شود (طبق (I)) $\|V - V^\pi\| \leq \frac{\varepsilon}{1-\gamma}$

و بنابراین نامعادی خواسته شده سوال ثابت می شود: $-\frac{\varepsilon}{1-\gamma} + V^*(s) \leq V^\pi(s)$

آنگاه به اثبات ادعایمان می پردازیم. توجه کنید $V^\pi \leq V^*$ (چون $V^* = \max_{\pi} V^\pi$ است) و بنابراین آنگاه داریم $V^\pi \leq V^* \leq V$ حالا توجه کنید:

$$\|V - V^\pi\| = \max_s |V(s) - V^\pi(s)| = \max_s V(s) - V^\pi(s)$$

$$\|V - V^\pi\| = \max_s |V^*(s) - V^\pi(s)| = \max_s V^*(s) - V^\pi(s)$$

با توجه به آنکه می داریم $\forall s: V(s) \leq V^*(s)$ ، به وضوح نتیجه می شود:

$$\max_s V^*(s) - V^\pi(s) \leq \max_s V(s) - V^\pi(s) \Rightarrow \|V^* - V^\pi\| \leq \|V - V^\pi\|$$

بنابراین ادعایمان ثابت شد و کارن بهتر هم نتیجه می شود.

۲.۲ (۱۰) ابتدا به استقراری n ثابت می کنیم $B^n V \leq V$. حالت $n=1$ به عنوان فرض داده شده. $B^1 V \leq V$ استقرای

استقرای ثابت $B(B^n V) \leq BV \Rightarrow B^{n+1} V \leq BV \leq B$

سوال ۳: $V \leq V' \Rightarrow BV \leq BV'$

حالا توجه کنید $\lim_{n \rightarrow \infty} B^n V = V^*$. از ادعای که اثبات کردیم هم استفاده می کنیم:

$$V^* = \lim_{n \rightarrow \infty} B^n V \leq V \Rightarrow V^* \leq V$$

(توجه کنید در سوال ۳ داشتیم اگر $V \leq V'$ ، آنکه $B(V) \leq B(V')$ اما اگر قرار دهیم $\pi = \pi_V^*$ ،

که بالایی بهینه با مقادیر V است، نتیجه می شود $BV = B^{\pi} V \leq B^{\pi} V' \leq BV'$ که همان نتیجه ای است که ما استفاده کردیم).

با براین ثابت کردیم $BV \leq V$ یک شرط کافی برای $V^* \leq V$ است و به وضوح چون دیگر نیاز به دانستن V^* نداریم، BV به راحتی قابل محاسبه است، چک کردن این شرط برایمان بسیار راحت تر است.

۲.۲ (۱۱) در سوال ۵ نتیجه زیر را گرفتیم:

$$\|B^{\pi} V - V^{\pi}\| \leq \gamma \|V - V^{\pi}\| \quad (I)$$

همچنین با توجه به آنکه $BV = B^{\pi} V$ ، در سوال ۹ گفتیم: (نتیجه سوال ۵)

$$\|V - V^{\pi}\| \leq \frac{\varepsilon}{1-\gamma} \quad (II)$$

$$(I), (II) \Rightarrow \|B^{\pi} V - V^{\pi}\| \leq \frac{\gamma \varepsilon}{1-\gamma} \quad (*)$$

الگزن ادما می کنیم $\|V^* - B^{\pi} V\| \leq \frac{\gamma \varepsilon}{1-\gamma}$. توجه کنید در صورت اثبات این ادعا، با نامساوی مثلث نتایج دلخواهتان را می گیریم:

$$\frac{\gamma \varepsilon}{1-\gamma} \geq \|V^* - B^{\pi} V\| + \|B^{\pi} V - V^{\pi}\| \geq \|V^* - V^{\pi}\|$$

که مستقیماً نامساوی اول را نتیجه می دهد. الگزن به اثبات ادماهای پردلترم. در سوال ۵ اگر بجای V مقدار BV را قرار دهیم نتیجه می گیریم: (نامساوی دوم)

$$\|BV - V^*\| \leq \frac{\|BV - B^{\pi} V\|}{1-\gamma} \quad (III)$$

$$\text{Contraction Property} \rightsquigarrow \|BV - B^2 V\| \leq \gamma \|V - BV\| = \gamma \varepsilon \quad (IV)$$

$$BV = B^{\pi} V \Rightarrow \|V^* - B^{\pi} V\| = \|BV - V^*\| \quad (V)$$

$$(II), (IV), (V) \Rightarrow \|V^* - B^{\pi} V\| \leq \frac{\gamma \varepsilon}{1-\gamma} \rightarrow \text{ادما ثابت شد}$$

نامساوی دوم نیز راحت است. توجه کنید $V \geq V^*$ نتیجه می دهد $BV \geq BV^* = V^*$

$$V^* - V^{\pi} \leq BV - V^{\pi} = B^{\pi} V - V^{\pi} \stackrel{(*)}{\leq} \frac{\gamma \varepsilon}{1-\gamma} \quad \text{پس:}$$

و نامساوی دوم هم اثبات شد. ■