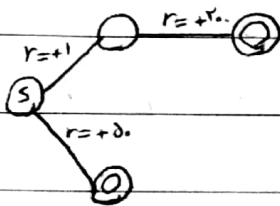


① (آ) عملیات. به میزان مثال نقض می‌توان MDP را در نظر گرفت. اگر $\gamma = 0.9$ باشد آنگاه $Q^*(s, down) = 5$ و $Q^*(s, up) = 1 + 0.9 \times 20 = 21$ می‌باشد. بنابراین اگر $\gamma = 0.9$ در حالی که اگر $\gamma = 0.9$ باشد آنگاه $Q^*(s, down) = 5$ و $Q^*(s, up) = 1 + 0.9 \times 20 = 21$ می‌باشد. بنابراین اگر $\gamma = 0.9$ در حالی که اگر $\gamma = 0.9$ باشد آنگاه $Q^*(s, down) = 5$ و $Q^*(s, up) = 1 + 0.9 \times 20 = 21$ می‌باشد. مشاهده می‌شود که سیاست بهینه در حالت s با تغییر ضریب تخفیف، تغییر می‌کند.



ب) مسئله: طراحی یک سیستم پیشنهاد دهنده برای یک ربات زودتی کتاب. سیستم باید بر اساس تاریخچه خرید و استیازهای کاربر آن به آنها کتاب پیشنهاد کند.

مدل ۱: عامل را سیستم پیشنهاد دهنده در نظر می‌گیریم و محیط را مشتریان. در این مدل حالات محیط برای تاریخچه خرید و استیازهای مشتریان است. اگر کسی می‌خواهد که عامل می‌تواند انجام دهد، کتاب‌هایی است که می‌تواند پیشنهاد دهد. پاداش نیز فیدبکی است که مشتری (محیط) به سیستم (عامل) می‌دهد.

مدل ۲: عامل را مشتری در نظر می‌گیریم و محیط را سیستم پیشنهاد دهنده. در این مدل حالات محیط برای لیست کتاب‌های موجود است. اگر کسی می‌خواهد که عامل می‌تواند انجام دهد کتاب‌هایی است که حاضر در پاداش نیز میزان رضایت مشتری از خرید است.

اگر بخواهیم سیستم پیشنهاد دهنده را بر اساس بازخورد مشتری بهینه کنیم، مدل ۱ بهتر است. اما اگر بخواهیم رفتار مشتری و علایق وی را در یاد بگیریم، مدل ۲ بهتر است.

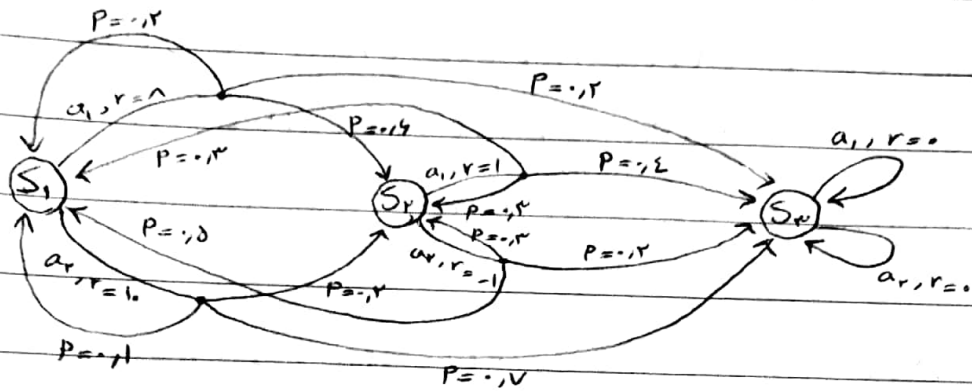
ج) الگوریتم Value Iteration با استفاده از معادلات بلمن، مقادیر سودمندی بهینه را به دست می‌آورد. پس با استفاده از این مقادیر سودمندی بهینه، سیاست بهینه را استخراج می‌کند. این روش اگر تعداد حالات زیاد باشد یا تابع پاداش به خوبی تعریف نشده باشد، به کندی همگرا می‌شود. الگوریتم Policy Iteration ابتدا با استفاده از یک سیاست ثابت π مقادیر سودمندی بهینه را محاسبه می‌کند و سپس با استفاده از این مقادیر، سیاست خود را به روزرسانی می‌کند. در این روش اگر آپدیت سیاست منجر به افزایش سودمندی نشود ممکن است به کندی همگرا شود. به علاوه یکی از این دو الگوریتم به دلیل

نیاز به ذخیره مقادیر سودمندی سیاست ما نیاز به حافظه زیادی هستند. همچنین به علاوه یکی از الگوریتم Policy Iteration نسبت به Value Iteration از لحاظ محاسباتی بهینه‌تر است و معمولاً زودتر همگرا می‌شود.

$$Q(1) = \begin{matrix} / \\ -2 \end{matrix} \quad Q(2) = \begin{matrix} / \\ 3 \\ 2,5 \end{matrix} \quad Q(3) = \begin{matrix} / \\ -1 \\ -5,5 \end{matrix} \quad Q(4) = \begin{matrix} / \\ 5 \end{matrix} \quad Q(5) = 0 \quad (2)$$

	ϵ random	$1-\epsilon$ greedy	
$(A_1 = 1, R_1 = -2)$	✓	✓	
$(A_2 = 2, R_2 = 3)$	✓	✓	
$(A_3 = 3, R_3 = -1)$	✓	X	در اینجا بهترین حالت اگر ۲ بود
$(A_4 = 2, R_4 = 2)$	✓	✓	
$(A_5 = 3, R_5 = 0)$	✓	X	در اینجا بهترین حالت اگر ۲ بود
$(A_6 = 4, R_6 = 5)$	✓	X	در اینجا بهترین حالت اگر ۲ بود

در مراحل ۳، ۵ و ۶ حتماً ϵ رخ داده و (گشای تعدادی انتخاب شده در بار اول در مرحله ۳ و ۵ و ۶ یعنی انتخاب اگر تعدادی در میانه ممکن بوده انتخاب شده باشند



$\gamma = 1$ (2)

$$V_0(S_1) = 1.0$$

$$V_0(S_2) = 1$$

$$V_0(S_3) = 0$$

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$

$$V_1(S_1) = \max \left\{ \begin{array}{l} 0.2(1 + V_0(S_1)) + 0.4(1 + V_0(S_2)) + 0.2(1 + V_0(S_3)) \\ 0.1(1 + V_0(S_1)) + 0.5(1 + V_0(S_2)) + 0.1(1 + V_0(S_3)) \end{array} \right\}$$

$$\Rightarrow V_1(S_1) = 1.1, \gamma \quad Q_1(S_1, a_1) = 1.0, \gamma \quad Q_1(S_1, a_2) = 1.1, \gamma \quad \pi_1(S_1) = a_1$$

$$V_1(S_2) = \max \left\{ \begin{array}{l} 0.4(1 + V_0(S_2)) + 0.3(1 + V_0(S_2)) + 0.1(1 + V_0(S_3)) \\ 0.3(-1 + V_0(S_1)) + 0.4(-1 + V_0(S_2)) + 0.2(-1 + V_0(S_3)) \end{array} \right\}$$

$$\Rightarrow V_1(S_2) = 0.5, \gamma \quad Q_1(S_2, a_1) = 0.5, \gamma \quad Q_1(S_2, a_2) = 0.5, \gamma \quad \pi_1(S_2) = a_1 \text{ or } a_2$$

$$V_1(S_3) = \max \left\{ 1(0 + V_0(S_3)), 1(0 + V_0(S_3)) \right\}$$

$$\Rightarrow V_1(S_3) = 0 \quad Q_1(S_3, a_1) = 0 \quad Q_1(S_3, a_2) = 0 \quad \pi_1(S_3) = a_1 \text{ or } a_2$$

$$V_r(S_1) = \max \left\{ \begin{array}{l} 12,82 \\ 0,2(1 + V_r(S_1)) + 0,4(1 + V_r(S_2)) + 0,2(1 + V_r(S_3)) \\ 0,1(1 + V_r(S_1)) + 0,2(1 + V_r(S_2)) + 0,7(1 + V_r(S_3)) \end{array} \right\}$$

$$\hookrightarrow V_r(S_1) = 12,82 \quad Q_r(S_1, a_1) = 12,82 \quad Q_r(S_1, a_2) = 11,98 \quad \pi_r(S_1) = a_1$$

$$V_r(S_2) = \max \left\{ \begin{array}{l} 8,48 \\ 0,3(1 + V_r(S_1)) + 0,3(1 + V_r(S_2)) + 0,4(1 + V_r(S_3)) \\ 0,5(-1 + V_r(S_1)) + 0,3(-1 + V_r(S_2)) + 0,2(-1 + V_r(S_3)) \end{array} \right\}$$

$$\hookrightarrow V_r(S_2) = 8,89 \quad Q_r(S_2, a_1) = 8,48 \quad Q_r(S_2, a_2) = 8,89 \quad \pi_r(S_2) = a_2$$

$$V_r(S_3) = \max \left\{ 1(1 + V_r(S_1)) , 1(1 + V_r(S_2)) \right\}$$

$$\hookrightarrow V_r(S_3) = 0 \quad Q_r(S_3, a_1) = 0 \quad Q_r(S_3, a_2) = 0 \quad \pi_r(S_3) = a_1 \text{ or } a_2$$

حی) برای اینکه ثابت کنیم که سیاست بهینه در S_1 برای a_1 است، در یک معنی نمی شود باید داشته باشیم:

$$0,2(1 + V_K(S_1)) + 0,4(1 + V_K(S_2)) + 0,2(1 + V_K(S_3)) \geq$$

$$0,1(1 + V_K(S_1)) + 0,2(1 + V_K(S_2)) + 0,7(1 + V_K(S_3))$$

$$\Rightarrow 0,1 V_K(S_1) + 0,2 V_K(S_2) \geq 2$$

از آنجایی که $V_K(S_1) \geq 11,2$ و $V_K(S_2) \geq 4,3$ بنابراین راه حل بالا به ازای $K \geq 2$ همواره برقرار است.

برای اینکه ثابت کنیم که سیاست بهینه در S_2 برای a_2 است، در یک معنی نمی شود باید داشته باشیم:

$$0,5(-1 + V_K(S_1)) + 0,3(-1 + V_K(S_2)) + 0,2(-1 + V_K(S_3)) \geq$$

$$0,3(1 + V_K(S_1)) + 0,3(1 + V_K(S_2)) + 0,4(1 + V_K(S_3))$$

$$\Rightarrow V_K(S_1) \geq 10$$

از آنجایی که $V_K(S_1) \geq 11,2$ بنابراین راه حل بالا به ازای $K \geq 2$ همواره برقرار است.

از آنجایی که همواره سود منفی در S_3 همواره برای هر سیاست بهینه در آن هیچگاه تغییری نمی کنند.

④ معماری I2A از ترکیب دو روش Model based و Model free بهره می برد که روش مبتنی بر مدل آن به پایبند مدل های محیطی است که مطابق Imagination و دورنگری عامل است. هدف از این جنبش این است که عامل قبل از اینکه در یک حالت خاص، یک اکشن خاصی را انجام دهد، بتواند با استفاده از این سازول Imagination یک تصویر از عواقب و حالاتی که پس از انجام یک اکشن مشخص به آن می رسد دست پیدا کند.

ساختار این سازول به این صورت است که حالت فعلی یا \hat{O}_t را به یک Policy net می دهیم و یک اکشن منتخب \hat{a}_t می گیریم. حال برای اینکه بتوانیم به یک تصویر از عواقب انجام این اکشن بگیریم، \hat{O}_t و \hat{a}_t را به یک Environment Model که اساساً یک شبکه عصبی است داد و \hat{O}_{t+1} و \hat{r}_{t+1} را به دست می آوریم. به ستمی که این بدست می آید را انجام می دهیم Imag. core می توان با کنار هم قرار دادن تعدادی از این Imag. core ما پست رسم، time step های بیشتری از آینده را تصور کنیم. در مرحله بعد فرضی حرکتی از این Imag. core ما را به یک Encoder که می تواند بلاک LSTM یا یک سیستم فرضی این encoder های پست رسم، یک Rollout ending می نامند که نشان دهنده عواقب انجام اکشن a_t در O_t است. به حرکتی از این بدست می آید Rollout encoder می گویند. با کنار هم قرار دادن تعدادی از این بدست می آید به تعدادات مختلفی رسید و عملکرد را از این روش دار. در نهایت فرضی این Rollout encoder ها تسلسل را می به نام Aggregator با هم Concat می کنند و به همراه فرضی شبکه Model free به شبکه استخراج گر سیاست داده می شوند.

* مزایا: دسترسی نازل به تخمینی از عواقب انجام اکشن های می تواند منجر به انجام exploration بهینه تر و درنتیج همگرا شدن سریع تر شود. / استفاده از این سازول می تواند به عامل در انتخاب و پنهان کردن اکشن های که منجر به یادداشت نهایی بیشتری می شوند کمک کند.

x معایب: استفاده از این سازول افشان دارای هزینه محاسباتی بالاتر و آشنایی مدل را کندتر میکند /

اگر این سازول با استفاده از یک مدل محیطی دارای پایاس آشنایی دیده باشد، ممکن است منجر به عملکرد غیر بهینه شود.

