

برنامه ریزی برای بستن آن

موضوعات درس 1 - برنامه ریزی برای 2 - ارزیابی سیاست 3 - بهبود سیاست 4 - تکرار سیاست 5 - تکرار ارزیابی 6 - تکرار سیاست

برنامه ریزی برای: مجموعه‌ای از درخت برای تعیین سیاست بهینه با داشتن مدل $EMDP$ است.

درای من زفنی کنیم $P(s', r | s, a)$ داریم.

توجه: فضای تشریف سله باید بسته و محدود باشد.

در عمل مدل اغلب دانش یادگیری تقوین **تخت** است.

کامپیوت برنامه ریزی برای فضا زیاد است.

اما برنامه ریزی برای یک تک تابع و درخت **Model Free** را ساده‌تری کند.

معادله بهینه‌ترین برای تابع ارزش حالت:

$$V_*(s) = \max_a E[R_{t+1} + \gamma V_*(s_{t+1}) | S_t = s, A_t = a]$$

$$= \max_a \sum_{s', r} P(s', r | s, a) [r + \gamma V_*(s')]$$

حل تحلیلی غیر ممکن است

در **تکرار مرتبه**

$$q_*(s, a) = E[R_{t+1} + \gamma \max_{a'} q_*(s_{t+1}, a') | S_t = s, A_t = a]$$

$$= \sum_{s', r} P(s', r | s, a) [r + \gamma \max_{a'} q_*(s', a')]$$

ارزیابی سیاست:

در برنامه ریزی برای **دسترسی عینی** برای تابع ارزش و سیاست تعیین زودهی شوند.

$$V_\pi(s) = \sum_a \pi(a|s) \sum_{s', r} P(s', r | s, a) [r + \gamma V_\pi(s')]$$

شرط وجود جواب یکتا: نرخ تخفیف (discount factor) که کمتر از یک در ارسال به حالت‌های تحت سیاست \leftarrow اگر γ به یک برسد و سیاست ثابت این سیاست π باید بتوانیم از هر حالت اولیه به هر حالت نهایی برسیم. چون می‌توانیم ارزیابی تمامی s در سیاست π داریم.

جدول من کنیم: اولین راه این می‌باشد که برای **حل عددی** استفاده از ماتریس خطی است.

$$\begin{cases} E(1) = a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1 \\ E(2) = a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2 \\ \vdots \\ E(n) = a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n = b_n \end{cases}$$

و $x^{(0)}$ را می‌توانیم

$$\begin{bmatrix} x_1^{(0)} \\ x_2^{(0)} \\ \vdots \\ x_n^{(0)} \end{bmatrix}$$

s.a.m

$$A, \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} \quad B, \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} \quad x, A^{-1}B, A^{-1}B$$

نم اول روش ژالوسکی: در مرحله نام تغییر نام بر حسب سطر جدید است.

$$x_1, \frac{b_1 - [a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n]}{a_{11}}$$

$$x_2, \frac{b_2 - [a_{21}x_1 + a_{23}x_3 + \dots + a_{2n}x_n]}{a_{22}}$$

$$\vdots$$

$$x_n, \frac{b_n - [a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n]}{a_{nn}}$$

نکته: $b_i - \sum_{j \neq i} a_{ij}x_j$

نم دوم روش ژالوسکی: کاملاً مرتب مرحله اول

$$x_1^{(1)}, \frac{b_1 - [a_{12}x_2^{(0)} + a_{13}x_3^{(0)} + \dots + a_{1n}x_n^{(0)}]}{a_{11}}$$

به همین ترتیب برای سطر دوم و سطرهای دیگر.

while $x_i^{(k)}, \frac{b_i - \sum_{j \neq i} a_{ij}x_j^{(k-1)}}{a_{ii}} \quad k=1, 2, \dots$ نکته

$$\|x^k - x^{k-1}\| < \epsilon$$

مرحله پنجم برای تابع اندکی دارای صحت نرم ژالوسکی است.

$$\rightarrow U_k(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma U_k(s')]$$

$$U_{k+1}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma U_k(s')]$$

تکرار می شود $k=1, 2, \dots$

s.a.m

الگوریتم ارزیابی سیاست به مقدار تعیین تابع ارزش و سیاست
 تکراری انجام می‌دهیم \rightarrow **iterative policy evaluation**

Input π , the policy to be evaluated

Algorithm, parameter a small threshold $\theta > 0$ determining accuracy of estimation
 دقت تعیین ارزش \leftarrow برای توقف الگوریتم استفاده می‌شود ϵ در زکات

$$|V_{\pi}(s) - V_{\pi'}(s)| < \theta$$

Initialize $V(s)$, for all $s \in S$, arbitrarily except that $V(\text{terminal}) = 0$
 مقدار اولیه باید بدهیم

Loop:

$$\Delta \leftarrow 0$$

Loop for each $s \in S$:

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \sum_a \pi(a|s) \sum_{s'} P(s'|s,a) [r + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

$$\text{until } \Delta < \theta$$

$$\pi_1^{n+1} \leftarrow \pi_1^n, \pi_2^{n+1} \leftarrow \pi_2^n$$

$$\pi_1^{n+1} \leftarrow g(\pi_1^n, \pi_2^n)$$

$$\pi_1^{n+1} \leftarrow \pi_1^n, \pi_2^{n+1} \leftarrow \pi_2^n$$

$$\pi_2^{n+1} \leftarrow g(\pi_1^n, \pi_2^n)$$

برای state برابر با s و a این به ما π می‌دهد و وجود دارد.

به دست آمدن: در هر حالت احتمال انجام عمل که وجود دارد.

$$V_{k+1}(s) = \sum_a \pi(a|s) \sum_{s'} P(s'|s,a) [r + \gamma V_k(s')]$$

$$s_i \rightarrow a_i$$

$$s_i \rightarrow a_i$$

به دست آمدن: در هر حالت π به عمل مشخص می‌شود انجام است.

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

اینکه به دست می‌آید $\pi(s)$ است $\pi(s)$ \rightarrow $\pi(s)$ \rightarrow $\pi(s)$

همه‌ی توان از π به دست می‌آید به دست می‌آید به دست می‌آید

به دست می‌آید $\pi(s)$ به دست می‌آید $\pi(s)$ به دست می‌آید $\pi(s)$

$$V_{\pi'}(s) \geq V_{\pi}(s) \rightarrow \text{policy improvement theorem}$$

$$V_{\pi'}(s) \geq V_{\pi}(s)$$

* با تغییر ارزشی حالت نسبت به سیاست اولیه، می توان سیاست بهبود یافته را تعیین کرد.
اگر q_k یا $q_{k'}$ را داشته باشیم می تونه به k برسیه.

$$q_k(s, a) \doteq E[R_{t+1} + \gamma V_k(S_{t+1}) | S_t = s, A_t = a]$$

$$= \sum_{s', r} p(s', r | s, a) [r + \gamma V_k(s')]$$

حداکثر باید انتخاب شود تا این عبارت \max شود.

← سیاست جدید نسبت به تابع ارزشی حالت عمل greedy می شود.

$$k'(s) = \arg \max_a q_k(s, a)$$

$$= \arg \max_a E[R_{t+1} + \gamma V_k(S_{t+1}) | S_t = s, A_t = a]$$

$$= \arg \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_k(s')]$$

آنها q_k را داشته باشیم تبدیل MDP می تونه است.

به همین دلیل سیاست جدید بهتر از سیاست قبلی است.

برای بهبود نهایی است ارزشی حالت **سیاست** بهتر شود.

چون داریم بهبودی گزینی ارزشی را تعیین می کنیم، $k'(s)$ با بهره مندی از روشی که به دست می آید.
← جدید greedy از طریق انتخاب عمل برای **بزرگترین بازدهی کوتاه مدت** است.

اما چون **نیاز نیست** Return در الگوریتم بهبود سیاست حضور دارند، در بلند مدت هم greedy خواهد بود.

* اگر تابع ارزشی سیاست اولیه را داشته باشیم می تونه به **سیاست بهینه** برسیه ایم.

$$V_{k'}(s) = \max_a E[R_{t+1} + \gamma V_{k'}(S_{t+1}) | S_t = s, A_t = a]$$

$$= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_{k'}(s')]$$

$$V_{k'}(s) = \sum_a \pi(a|s) q_k(s, a)$$

$$V_{k'}(s) = \max_a q^*(s, a)$$

تکرار سیاست فیدبک تعیین سیاست جدید داشته می تونیم تابع ارزشی حالت آن را هم می تونه نیمه و می تونه از آن جدا به سیاست بهتری برسیم.
میکنیم **ارزشی سیاست** و **بهره سیاست** تا رسیدن به سیاست بهینه می تونه تکرار شود.

s.a.m



هر سیاست جدید اغلب بهتر از سیاست قبلی است. بنابراین به سیاست بهینه برسیم.
چون MDP دارای حالت دو سیاست دی (عددی) است، تکرار سیاست به سیاست بهینه منجر فراموش می‌شود.

1. Initialization

$U(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation

Loop,

$\Delta \leftarrow 0$

Loop for each $s \in \mathcal{S}$:

$U \leftarrow U(s)$

$U(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) [r + \gamma U(s')]$

$\Delta \leftarrow \max(\Delta, |U - U(s)|)$

until $\Delta < \theta$ Caswell positive number determining the accuracy of evaluation.

3. Policy Improvement

Policy stable \leftarrow true

For each $s \in \mathcal{S}$:

old-action $\leftarrow \pi(s)$

$\pi(s) \leftarrow \arg \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma U(s')]$

if old-action $\neq \pi(s)$, then policy-stable \leftarrow false

if policy-stable, then stop and return $U \approx U_*$ and $\pi \approx \pi_*$; else go to 2

S^+ - state دی نه "مختار" است
 S - state دی نه "مختار" است
 S^+ - state دی نه "مختار" است
Terminal state

معمولاً به صورت زیر است:

تکرار این است:

الگوریتم بهینه‌سازی

برای هر $s \in \mathcal{S}$ و $a \in \mathcal{A}(s)$ به صورت زیر:

$\pi_{k+1}(s) = \arg \max_a Q_k(s,a)$

این به ما می‌دهد بهترین سیاست برای هر s

s.a.m

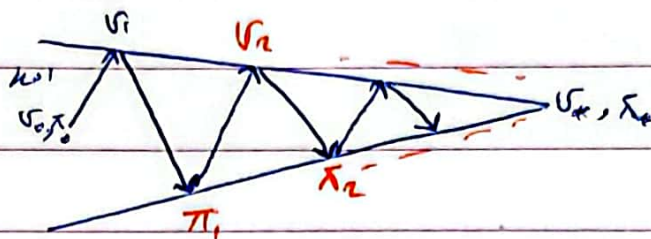
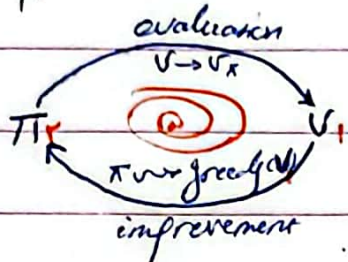
تکرار سیاست تقسیم یافته

Generalized Policy Iteration (GPI)

در تکرار سیاست باید مسئله ارزیابی و بهتری تابع ارزش حاصل تکرار شوند.

در فصل ارزیابی سیاست، خود سیاست ثابت می ماند.

در ادامه تکرار سیاست تقسیم یافته (GPI) در پایان هر حلقه ارزیابی، خود سیاست هم بهبود دارد، خواهد شد. (مانند تقسیم درج زوالی به دو سری سایل)



s.a.m