

روش‌های مونت کارلو (۲ قسمت)

یادگیری تقویتی در قالب ارزش‌های

به دست می‌دهد. تخمین ارزش حالت با $Monte Carlo$ ، تخمین ارزش حالت عمل با $Monte Carlo$ ، و $Confounding Monte Carlo Control$ و $Off-policy Monte Carlo Control$

این $Model free$ است یعنی در دنیا یک برنامه دیتا داریم و نتایج $episode$ و تجربه عملی داریم و تخمین ارزش را می‌زنیم و سیاست را انتخاب می‌کنیم.

* فرض اساسی در برنامه‌ریزی پویا داشتن مدل MDP است. $p(s', r | s, a)$

در عمل در مسائل واقعی با این تقویم، دنیا یک مشخص نیست.

* در روش $Monte Carlo$ بدون نیاز به دنیا یک ارزش a و سیاست π که تقریب زده خواهند شد.

* ارزش‌های MC تنها از طریق تجربه و تقویم با مجموعه یادگیری را می‌توانیم بگیریم.

* تجربه می‌تواند نامرتب از محیط نیز باشد یا محیط به صورتی

* در برخی موارد می‌توان دنیا یک را تخمین کرد ولی به دلیل کامپیات ریاضی در برنامه‌ریزی پویا و دیگری ندارند.

* در برخی موارد به هم توصیف ریاضی دنیا یک مسئله مشخص نیست. که تقریب 50% حالت

* ارزش‌های MC از نمونه برداری تقویم جهت کامپیات ریاضی استفاده می‌کنند.

* شبیه‌سازی سیستم و کامپیات عددی بهترین کاربرد برای MC است.

$$\int_a^b h(x) dx \quad \text{که شامل انتقال} \\ \text{وگر که} \quad h(x) \sim e^{\lambda x}$$

$$w(x) \sim h(x) \cdot (b-a)$$

$$h(x) \sim \frac{1}{b-a} \rightarrow U[a, b]$$

$$\int_a^b h(x) dx = \int_a^b w(x) h(x) dx \rightarrow \text{ایستادگی برای تابع قبلی} \\ \text{انتقال} \quad h(x)$$

$$E(x) = \int x h(x) dx \quad \text{یادگیری}$$

$$x_0, x_1 \in U[a, b]$$

$$E(g(x)) = \int g(x) h(x) dx$$

برای $h(x)$ انتگرال از توزیع $h(x)$ است. یعنی $E(g(x))$ را می‌توانیم بنویسیم. $h(x)$ مجموعه $h(x)$ است که مجموعه $h(x)$ است. $h(x)$ مجموعه $h(x)$ است.

تخمین ارزش حالت با Monte Carlo:

* درخت های کامل اغلب برای تخمین ارزش استفاده می شود.

* بعد از پایان هر دوره، ارزش بازی با درخت ارزش حالت را تخمین زده خواهد شد.

$$S_0, A_0, R_0, S_1, A_1, R_1, S_2, A_2, R_2, \dots, S_T$$

$$G_t = R_{t+1} + \gamma V_{t+1} = \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1}$$

$$V_{\pi}(s) = E_{\pi}[G_t | S_t = s] \quad q_{\pi}(s, a) = E_{\pi}[G_t | S_t = s, A_t = a]$$

در صورت این state درخت ترس می شود.

* با به یادداشت کردن شماره بازی که در آن state به اولین state می رسد.

Return

تقریباً به اندازه ارزش بازی که در آن state به اولین state می رسد.

$$V_{R_t}(s) \rightarrow V_{R_{t+1}}(s)$$

$$V_{R_{t+1}}(s) = V_{R_t}(s) + \frac{1}{n} (G_t(s) - V_{R_t}(s))$$

$$Q_{R_{t+1}}(s, a) = Q_{R_t}(s, a) + \frac{1}{n} (R_t - Q_{R_t}(s, a))$$

$$Q_{R_{t+1}}(s, a) = Q_{R_t}(s, a) + \frac{1}{n} (G_t(s, a) - Q_{R_t}(s, a))$$

Input: a policy π to be evaluated

الگوریتم First-Visit

Initialize:

$V(s) \in \mathbb{R}$, arbitrarily, for all $s \in S$

$Reurns(s) \leftarrow$ an empty list, for all $s \in S$ [* * *]

Loop forever (for each episode):

Generate an episode showing $\pi, S_0, A_0, R_0, S_1, A_1, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$$G \leftarrow 0 \quad \left(G_t = R_{t+1} + \gamma G_t \right)$$

Loop for each state s in episode:

Step of episode, $t = T-1, T-2, \dots, 0$

$$G \leftarrow \gamma G + R_{t+1}$$

Unless S_t appears in S_0, S_1, \dots, S_{t-1}

Append G to Returns(S_t)

$$V(S_t) \leftarrow \text{average}(\text{Returns}(S_t))$$

روش MC برای یادگیری غیرایستایی α از آن استفاده کنیم

Constant α Monte Carlo

$$V(S_{t+1}) \leftarrow V(S_t) + \alpha [G_t - V(S_t)]$$

یادگیری تدریجی در صورت تغییر
تواند برای یادگیری در محیط‌های متغیر استفاده شود.
New Estimate
Old Estimate + step size [Target - Old Estimate]

در یادگیری در محیط‌های متغیر، این روش حالت را به روز رسانی می‌کند.

تجربین این روش را با Monte Carlo و با داشتن این روش با مدل MDP و عدم دانستن این روش با مدل بدون مدل.
در یادگیری این روش با مدل و بدون مدل به تدریج می‌توان یادگیری را بهبود داد.

$$Q_\pi(s, a) \text{ is argument } Q_\pi(s, a)$$

$$= \text{argument } E [R_{t+1} + \gamma V_\pi(s_{t+1}) | s_t, s, A_t = a]$$

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

$$Q_\pi(s, a) \text{ is } E [G_t | s_t, s, A_t = a]$$

این روش در حالت s انجام می‌دهد
 $V_\pi = E_\pi [G_t | s_t, s]$
در این روش G_t تدریجی یاد می‌گیرد.

به یادگیری **Bootstrapping**!

ممکن است برخی زوج حالت-عمل در $trajectory$ ظاهر نشوند.
اگر به یادگیری نیاز باشد، باید در هر حالت عمل را تکرار کرد.
در state تنها $action$ می‌تواند $action$ باشد؟
Exploring starts (راه):
در تمام زوج حالت-عمل با یکسان است و دانسته می‌شود.

On-policy Monte Carlo
off-policy Monte Carlo

یادگیری بهتر:

s.a.m

Initialize

$\pi(s) \in A(s)$

الزيارة الأولى

s.a.m