



یادگیری تقویتی
تفاضل زمانی

محسن هوشمند
دانشکده تکنولوژی اطلاعات و علم رایانه
دانشگاه تحصیلات تکمیلی علوم پایه زنجان

مقدمه

تفاضلات زمانی یا موقتی

- استفاده از خطای بین دو سیگنال جهت یادگیری تقویتی
- اهم الگوریتم‌های ت
- ترکیبی از برنامه‌ریزی پویا و مونت کارلو
- نحوه یادگیری محدود به تعامل با محیط و بدون نیاز به مدلی از دینامیک محیط
- همانند مونت کارلو
- یادگیری از روی تخمین‌های یادگرفته شده سایر حالت‌ها و عدم انتظار رسیدن به نتیجه نهایی
- مانند برنامه‌ریزی پویا
- وصال خویشتن

تفصیل موضوع در دو بخش

- سیاست‌سنجی
- تخمین V_π برای سیاست حاضر π
- اصلاح سیاست
- اصلاح تا یافتن سیاست بهینه
- بپ، م‌ک، و ت‌ز از نوع تکرار سیاست عمومی
- \Leftarrow پس تفاوت در نحوه سیاست‌سنجی

سیاست سنجی

استفاده م ک و ت ز از تجربه جهت سنجش

کسب تخمین V

▪ با پیگیری سیاست π و کسب تجربه از طریق اجرای آن

روش مونت کارلو

▪ انتظار تا مشخص شدن بازده ملاقات حالت

▪ استفاده از بازده به عنوان هدف $V(S_t)$

▪ روش م ک هم ملاقات در محیط‌های نامانا

▪ با هدف متفاوت عدم لزوم انتظار تا رسیدن به پایان اپیزود

$$V(S_t) = V(S_t) + \alpha[G_t - V(S_t)]$$

G_t بازده واقعی پس از زمان t

روش م ک ضریب ثابت

انتظار م ک تا پایان اپیزود

سیاست سنجی

انتظار م ک تا پایان اپیزود

انتظار ت ز تا گام زمانی بعدی

ایجاد هدف و بروز کردن $V(S_{t+1})$ بر اساس پاداش دریافتی R_{t+1} در زمان $t + 1$

ساده‌ترین ت ز

$$V(S_t) = V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$
$$V(S_t) = \underbrace{V(S_t)}_{\text{ارزش فعلی}} + \alpha \underbrace{[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]}_{\text{خطا}}$$

مقایسه

- هدف م ک G_t
- هدف ت ز $R_{t+1} + \gamma V(S_{t+1})$
- مشهور به ت ز (0) یا ت ز تک-گام
- نوعی خاص از ت ز (λ) و ت ز چند-گام

سیاست سنجی - تفاضلات زمانی (صفر)

▪ الگوریتم تفاضل زمانی صفر سیاست سنج

Tabular TD(0) for estimating v_π

Input: the policy π to be evaluated

Algorithm parameter: step size $\alpha \in (0, 1]$

Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop for each episode:

 Initialize S

 Loop for each step of episode:

$A \leftarrow$ action given by π for S

 Take action A , observe R, S'

$V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$

$S \leftarrow S'$

 until S is terminal

تفاضلات زمانی

بروزرسانی ارزش‌ها پس از اعمال حالتی به محیط و دریافت پاداش موردانتظار

- بروزرسانی بر اساس تخمین موجود
- \Leftarrow روش وصال خویشتن

$$\begin{aligned}v_{\pi}(s) &= E_{\pi}[G_t | S_t = s] \\ &= E_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= E_{\pi}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s]\end{aligned}$$

معادلات بلمن تابع ارزش

تفاضلات زمانی

بروزرسانی ارزش‌ها پس از اعمال حالتی به محیط و دریافت پاداش موردانتظار
▪ بروزرسانی بر اساس تخمین موجود
▪ روش وصال خویشتن \Leftarrow

\Leftarrow روش م ک

$$\begin{aligned} v_{\pi}(s) &= E_{\pi}[G_t | S_t = s] \\ &= E_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= E_{\pi}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s] \end{aligned}$$

\Leftarrow روش ب پ

معادلات بلمن تابع ارزش

م ک

▪ از نوع تخمین

▪ دلیل؟ ارزش موردانتظار معلوم نیست، بازده نمونه به جایش در نظر گرفته می‌شود

ب پ

▪ از نوع تخمین

▪ دلیل؟ $v_{\pi}(S_{t+1})$ معلوم نیست، استفاده از $V(S_{t+1})$ به جای آن

ت ز

▪ از نوع تخمین

▪ دلیل؟ هر دو دلیل م ک و ب پ

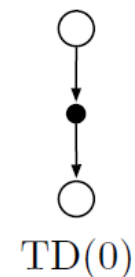
▪ $v_{\pi}(S_{t+1})$ معلوم نیست، استفاده از $V(S_{t+1})$ به جای آن

تفاضلات زمانی - قیاسی با والدان؟!!

تذکره ترکیب نمونه برداری م ک با وصال خویشتن ب پ
 امکان رسیدن به مزایای هر دو

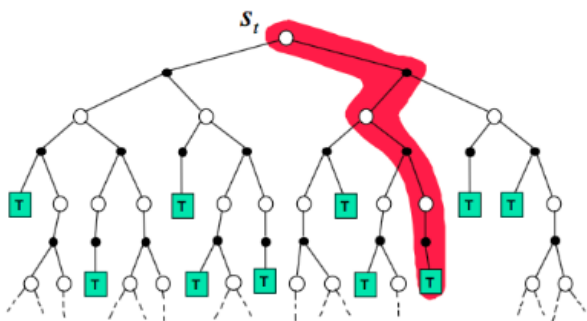
بروزرسانی

- بروزرسانی نمونه (م ک و تذ) صرفا تک نمونه
- بروزرسانی امید مورد انتظار (ب پ) توزیع کامل تمامی جانشین ها



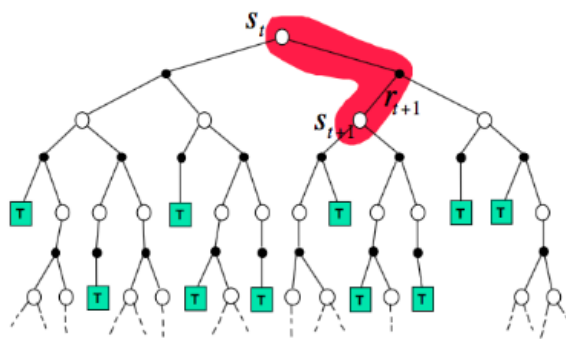
مونت کارلو

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$$



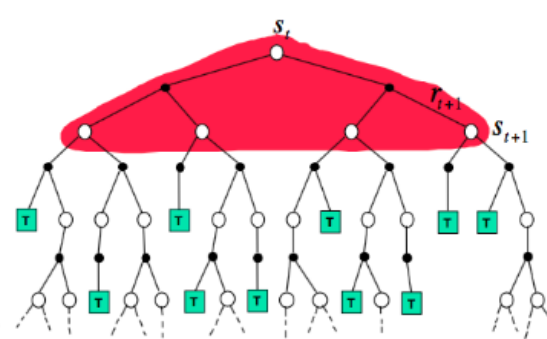
تفاضل زمانی

$$V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$



برنامه ریزی پویا

$$V(S_t) \leftarrow \mathbb{E}_\pi [R_{t+1} + \gamma V(S_{t+1})]$$



تفاضلات زمانی

$$V(S_t) = \underbrace{V(S_t)}_{\text{ارزش فعلی}} + \alpha \underbrace{[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]}_{\text{خطا}}$$

خطای مورد استفاده برای بروزرسانی ارزش‌ها

▪ $s \rightarrow s'$

▪ $r + \gamma v_\pi(s') - v_\pi(s)$ معروف به خطای تز

▪ $v_\pi(s)$ ارزش فعلی

▪ $v_\pi(s')$ هدف

▪ بروزرسانی با استفاده از میزان خطای بالا

تفاضلات زمانی

خطای تز

$$\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$$

امکان نوشتن «تقریبی» خطای م ک بر اساس جمع خطاهای تز

$$G_t - V(S_t)$$

$$\begin{aligned} G_t - V(S_t) &= R_{t+1} + \gamma G_{t+1} - V(S_t) + \gamma V(S_{t+1}) - \gamma V(S_{t+1}) \\ &= \delta_t + \gamma(G_{t+1} - V(S_{t+1})) \\ &= \delta_t + \gamma\delta_{t+1} + \gamma^2(G_{t+2} - V(S_{t+2})) \\ &= \delta_t + \gamma\delta_{t+1} + \gamma^2\delta_{t+2} + \cdots + \gamma^{T-t-1}\delta_{T-1} + \gamma^{T-t}(G_T - V(S_T)) \\ &= \delta_t + \gamma\delta_{t+1} + \gamma^2\delta_{t+2} + \cdots + \gamma^{T-t-1}\delta_{T-1} + \gamma^{T-t}(0 - 0) \\ &= \sum_{k=t}^{T-1} \gamma^{k-t} \delta_k. \end{aligned}$$

تفاضل زمانی - مزایا

بهبود تخمین با استفاده از بخشی از دیگر تخمین‌ها

فراگرفتن حدسی از حدس دیگر
▪ وصال خویشتن

بی‌نیاز از مدل محیط
▪ برعکس بپ

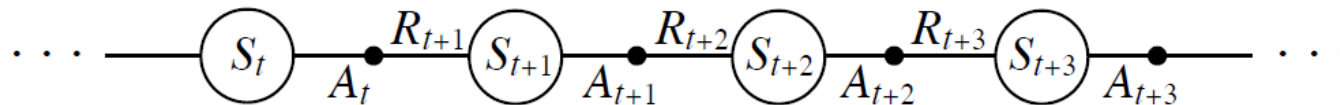
بی‌نیاز از رسیدن به پایان اپیزود و امکان اجرای برخط
▪ برعکس مونت کارلو

مک‌نیاز به انتظار تا پایان اپیزود - تز نیاز به انتظار در صرفاً یک گام
▪ اساسی برای کاربردهای با اپیزود طولانی

همگرایی به ارزش واقعی با آهنگ یادگیری کوچک

الگوریتم سارسا

تخمین ارزش حالت-کنش بهینه بر اساس تفاضل زمانی سیاست مدار
سنجش و اصلاح اندر کنار هم
بررسی انتقال از حالت-کنش اولیه به حالت-کنش بعدی



$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

فرزند وصال خویشتن
▪ استفاده از خطای موقت به جای امید ریاضی هر حالت-کنش

نیاز به نرمش اپسیلون
▪ چرا؟

سارسا- اصلاح تز سیاستمدار

یافتن بهینگی

استفاده از الگوی تکرار سیاست عمومی

رقابت کاوش و بهره‌برداری

▪ دو رده ← سیاستمدار و سیاست‌ن‌دار

نیاز به یادگیری تابع ارزش-کنش به جای ارزش-حالت

▪ روش سیاستمدار

▪ تخمین $q_{\pi}(s,a)$ رفتار فعلی سیاست π و تمامی حالت‌ها و کنش‌ها

سارسا- اصلاح تز سیاست مدار

سیاست سنجی (مرحله قبل)

- گذر از حالت به حالت
- یادگیری از ارزش حالتها

اصلاح (مرحله کنونی مدنظر)

- انتقال از زوج حالت-کنش به زوج دیگری از حالت-کنش
- یادگیری از ارزش زوجهای حالت-کنش

هر دو مورد بالا زنجیرههای مارکوفی با فرایند پاداش دهی

$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

اجرای تغییر بالا به ازای هر انتقال در حالت‌های غیر پایانی

اگر S_{t+1} حالت پایانی، آنگاه $Q(S_{t+1}, A_{t+1}) = 0$

اعمال قانون بالا بر هر پنج تایی $(S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1})$
▪ سارسا!

سارسا- اصلاح تز سیاست مدار

امکان و سادگی ایجاد الگوریتم اصلاح سیاست مدار سارسا

تخمین پیوسته q_π از سیاست و سپس تغییر π حریمانه با توجه به q_π

امکان استفاده از حریمانه-اپسیلون یا نرمش-اپسیلون جهت حصول به بهینگی

الگوریتم سارسا

Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

 Initialize S

 Choose A from S using policy derived from Q (e.g., ε -greedy)

 Loop for each step of episode:

 Take action A , observe R, S'

 Choose A' from S' using policy derived from Q (e.g., ε -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A';$

 until S is terminal

یادگیری ک

مبنی بر روش تفاضل زمانی سیاست‌نادر

واتکین ۱۹۸۹

اهم روش‌های بدون مدل

$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha \left[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \right]$$

تخمین مستقیم q_* مستقل از سیاست فعلی

▪ ساده‌سازی تحلیل الگوریتم و

▪ امکان‌دهی به اثبات همگرایی زودهنگام

بررسی انتقال از حالت-کنش اولیه به حالت-کنش بعدی

یادگیری ک

الگوریتم

Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

 Initialize S

 Loop for each step of episode:

 Choose A from S using policy derived from Q (e.g., ε -greedy)

 Take action A , observe R, S'

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$

 until S is terminal

یادگیری ک

تخمین ارزش بهینه مستقل از سیاست رفتاری

$$\begin{aligned} & q_*(s, a) \\ = & E \left[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a \right] \\ = & \sum_{s'} \sum_r p(s', r \mid s, a) \left[r + \max_{a'} q_*(s', a') \right] \end{aligned}$$

حالت انتقال قطعی

$$q_*(s, a) = r + \max_{a'} q_*(s', a')$$

نداشتن مقدار

▪ پس تبدیل به تابع خطای موقت (البته بدون ستاره)

$$r + \max_{a'} q_*(s', a') - q_*(s, a)$$

یادگیری ک

سارسا

$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

یادگیری ک

$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha \left[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \right]$$

تفاوت

- در تابع محاسبه خطا
- یادگیری ک سیاست‌نادر
- سارسا سیاست‌مدار

سیاست رفتاری

- نرمش اپسیلون

سخن کوتاه

یادگیری تفاضلی زمانی و امکان استفاده از آن در مسائل ی ت

سارسا

سیاستمدار

یادگیری ک

سیاست‌ن‌دار

روش‌های دیگر

یادگیری ک-دوبل

سارسای مورد انتظار ← سیاست‌ن‌دار

روش‌های منتقد-بازیگر

پیشبرد دانش و فن

ریشه‌گیری تز

از روانشناسی یادگیری حیوانات و هوش مصنوعی

خاصه کارهای کلوفف و ساموئل

تحقیقات هالند درباره سازگاری بین پیش‌بینی‌های ارزش

مدرس در دان‌میش

موثر بر بارتو؟ بارتو کیست؟

دانشجوی تحصیلات تکمیلی از ۱۹۷۰ تا ۱۹۷۵ در دان‌میش

یافته‌های هالند موثر در تعداد از نظام‌های تز-محور

منابع

ساتن

زندى