

روش تفاضل‌های زمانی

سید سینا ایروانیان

دانشگاه صنعتی شریف

۵ خرداد ۱۳۸۸

روش تفاضل‌های زمانی

روش تفاضل‌های زمانی یک روش یادگیری است

- یادگیری برای پیش‌بینی وضعیت یک سیستم پویا در دراز مدّت
- یادگیری برای کنترل \Leftarrow کاربرد در یادگیری تقویتی

انواع روش‌های یادگیری:

- **یادگیری بانظارت:** عامل با یک سری ورودی و پاسخ صحیح از قبل تعیین شده آموزش می‌بیند
- **یادگیری بی‌نظارت:** عامل هیچ بازخوردی دریافت نمی‌کند؛ ورودی‌ها را بصورت خوشه‌ها، یا طبقه‌بندی‌ها بازسازی می‌کند
- **یادگیری تقویتی:** یادگیری از طریق تعامل با محیط، و سعی و خطا توسط خود عامل

یادگیری به روش تفاضل‌های زمانی از نوع یادگیری تقویتی به حساب می‌آید

روش تفاضل‌های زمانی

روش تفاضل‌های زمانی یک روش یادگیری است

- یادگیری برای پیش‌بینی وضعیت یک سیستم پویا در دراز مدّت
- یادگیری برای کنترل \Leftarrow کاربرد در یادگیری تقویتی

انواع روش‌های یادگیری:

- **یادگیری بانظارت:** عامل با یک سری ورودی و پاسخ صحیح از قبل تعیین شده آموزش می‌بیند
- **یادگیری بی‌نظارت:** عامل هیچ بازخوردی دریافت نمی‌کند؛ ورودی‌ها را بصورت خوشه‌ها، یا طبقه‌بندی‌ها بازسازی می‌کند
- **یادگیری تقویتی:** یادگیری از طریق تعامل با محیط، و سعی و خطا توسط خود عامل

یادگیری به روش تفاضل‌های زمانی از نوع یادگیری تقویتی به حساب می‌آید

روش تفاضل‌های زمانی

روش تفاضل‌های زمانی یک روش یادگیری است

- یادگیری برای پیش‌بینی وضعیت یک سیستم پویا در دراز مدت
- یادگیری برای کنترل \Leftarrow کاربرد در یادگیری تقویتی

انواع روش‌های یادگیری:

- **یادگیری بانظارت:** عامل با یک سری ورودی و پاسخ صحیح از قبل تعیین شده آموزش می‌بیند
- **یادگیری بی‌نظارت:** عامل هیچ بازخوردی دریافت نمی‌کند؛ ورودی‌ها را بصورت خوشه‌ها، یا طبقه‌بندی‌ها بازسازی می‌کند
- **یادگیری تقویتی:** یادگیری از طریق تعامل با محیط، و سعی و خطا توسط خود عامل

یادگیری به روش تفاضل‌های زمانی از نوع یادگیری تقویتی به حساب می‌آید

پیش‌بینی چند مرحله‌ای

مسئله‌ی پیش‌بینی چند مرحله‌ای

- **ورودی:** دنباله‌ی مشاهدات-نتیجه: $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m, z$
- **خروجی:** P_1, P_2, \dots, P_m
- \mathbf{x}_t بردار مشاهده در مرحله‌ی t با درایه‌های حقیقی (برداری از اندازه‌گیری‌ها یا خصوصیات محیط) است
- z یک اسکالر حقیقی و نتیجه‌ی دنباله‌ی مشاهدات است
- P_t پیش‌بینی مقدار z که در مرحله‌ی t انجام شده، P تابعی از \mathbf{x}_t و بردار وزن‌ها \mathbf{w} ، (برداری از پارامترهای تغییر پذیر) است و آن را با $P(\mathbf{x}_t, \mathbf{w})$ نشان می‌دهند، به P تابع پیش‌بینی می‌گویند.

مثال

پیش‌بینی این که آیا روز جمعه باران می‌بارد یا خیر، از طریق مشاهده‌ی وضعیت هوای دوشنبه، و انجام پیش‌بینی برای روزهای متوال تا جمعه

پیش‌بینی چند مرحله‌ای

مسئله‌ی پیش‌بینی چند مرحله‌ای

- ورودی: دنباله‌ی مشاهدات-نتیجه: $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m, z$
- خروجی: P_1, P_2, \dots, P_m
- \mathbf{x}_t بردار مشاهده در مرحله‌ی t با درایه‌های حقیقی (برداری از اندازه‌گیری‌ها یا خصوصیات محیط) است
- z یک اسکالر حقیقی و نتیجه‌ی دنباله‌ی مشاهدات است
- P_t پیش‌بینی مقدار z که در مرحله‌ی t انجام شده، P تابعی از \mathbf{x}_t و بردار وزن‌ها \mathbf{w} ، (برداری از پارامترهای تغییر پذیر) است و آن را با $P(\mathbf{x}_t, \mathbf{w})$ نشان می‌دهند، به P تابع پیش‌بینی می‌گویند.

مثال

پیش‌بینی این که آیا روز جمعه باران می‌بارد یا خیر، از طریق مشاهده‌ی وضعیت هوای دوشنبه، و انجام پیش‌بینی برای روزهای متوالی تا جمعه

پیش‌بینی تک‌مرحله‌ای

مسئله‌ی پیش‌بینی تک‌مرحله‌ای

- پیش‌بینی چندمرحله‌ای به ازای $m = 1$
- دیگر در زمره‌ی مسائل تفاضل‌های زمانی قرار نمی‌گیرد
- در دسته‌ی یادگیری با نظارت قرار می‌گیرد (چون بلافاصله پس از پیش‌بینی، نتیجه نیز ارائه می‌شود)

مثال

پیش‌بینی این که آیا روز جمعه باران می‌بارد یا خیر، از طریق مشاهده‌ی وضعیت هوای پنج‌شنبه

پیش‌بینی تک‌مرحله‌ای

مسئله‌ی پیش‌بینی تک‌مرحله‌ای

- پیش‌بینی چندمرحله‌ای به ازای $m = 1$
- دیگر در زمره‌ی مسائل تفاضل‌های زمانی قرار نمی‌گیرد
- در دسته‌ی یادگیری با نظارت قرار می‌گیرد (چون بلافاصله پس از پیش‌بینی، نتیجه نیز ارائه می‌شود)

مثال

پیش‌بینی این که آیا روز جمعه باران می‌بارد یا خیر، از طریق مشاهده‌ی وضعیت هوای پنج‌شنبه

یادگیری با نظارت

- منظور از یادگیری، تنظیم مقادیر بردار وزن‌ها \mathbf{w} است، به‌طوری‌که خطای تابع پیش‌بینی در هر مرحله کمینه شود.
- به‌ازای هر مرحله می‌توان تغییر در بردار وزن‌ها $\Delta \mathbf{w}_t$ را در نظر گرفت

$$\mathbf{w} \leftarrow \mathbf{w} + \sum_{t=1}^m \Delta \mathbf{w}_t \quad (1)$$

- در یادگیری با نظارت تمام تغییرات در انتهای دنباله‌ی مشاهدات اعمال می‌شود؛ نه در هر مرحله

یادگیری با نظارت (ادامه)

- یادگیری با نظارت به هر دنباله‌ی مشاهدات-نتیجه، به صورت دنباله‌ای از زوج‌های مشاهده-نتیجه برخورد می‌کند

$$(\mathbf{x}_1, z), (\mathbf{x}_2, z), \dots, (\mathbf{x}_m, z)$$

- تغییرات در بردار وزن‌ها به‌ازای مشاهده در لحظه‌ی t :

$$\Delta \mathbf{w}_t = \alpha(z - P_t) \nabla_{\mathbf{w}} P_t \quad (2)$$

- α : نرخ یادگیری
- $\nabla_{\mathbf{w}} P_t$: بردار مشتقات جزئی P ، نسبت به \mathbf{w} ؛ مشخص‌کننده‌ی جهتی که در آن تغییرات در \mathbf{w} بیشترین تأثیر را در P دارد.

یادگیری با نظارت (ادامه) حالت خاص خطی

- حالت خاص: P_t تابع خطی از \mathbf{x}_t و \mathbf{w} است:

$$P_t = \mathbf{w}^T \mathbf{x}_t = \sum_i \mathbf{w}(i) \mathbf{x}(i)$$

- خواهیم داشت: $\nabla_{\mathbf{w}} P_t = \mathbf{x}_t$

قانون Widrow-Hoff برای بروزرسانی وزن‌ها - قانون دلتا

$$\Delta \mathbf{w}_t = \alpha (z - \mathbf{w}^T \mathbf{x}_t) \mathbf{x}_t \quad (۳)$$

- مقدار z در انتهای دنباله‌ی مشاهدات مشخص می‌شود، بنابراین تمامی \mathbf{x}_t ها باید ذخیره شوند، و تمامی بروزرسانی‌ها در انتهای دنباله محاسبه شوند.

یادگیری با نظارت (ادامه) حالت خاص خطی

- حالت خاص: P_t تابع خطی از \mathbf{x}_t و \mathbf{w} است:

$$P_t = \mathbf{w}^T \mathbf{x}_t = \sum_i \mathbf{w}(i) \mathbf{x}(i)$$

- خواهیم داشت: $\nabla_{\mathbf{w}} P_t = \mathbf{x}_t$

قانون Widrow-Hoff برای بروزرسانی وزن‌ها - قانون دلتا

$$\Delta \mathbf{w}_t = \alpha (z - \mathbf{w}^T \mathbf{x}_t) \mathbf{x}_t \quad (۳)$$

- مقدار z در انتهای دنباله‌ی مشاهدات مشخص می‌شود، بنابراین تمامی \mathbf{x}_t ها باید ذخیره شوند، و تمامی بروزرسانی‌ها در انتهای دنباله محاسبه شوند.

یادگیری با نظارت (ادامه) حالت خاص خطی

- حالت خاص: P_t تابع خطی از \mathbf{x}_t و \mathbf{w} است:

$$P_t = \mathbf{w}^T \mathbf{x}_t = \sum_i \mathbf{w}(i) \mathbf{x}(i)$$

- خواهیم داشت: $\nabla_{\mathbf{w}} P_t = \mathbf{x}_t$

قانون Widrow-Hoff برای بروزرسانی وزن‌ها - قانون دلتا

$$\Delta \mathbf{w}_t = \alpha (z - \mathbf{w}^T \mathbf{x}_t) \mathbf{x}_t \quad (۳)$$

- مقدار z در انتهای دنباله‌ی مشاهدات مشخص می‌شود، بنابراین تمامی \mathbf{x}_t ها باید ذخیره شوند، و تمامی بروزرسانی‌ها در انتهای دنباله محاسبه شوند.

محاسبه‌ی افزایشی

- نمایش خطای $z - P_t$ به صورت مجموع تغییرات در پیش‌بینی‌های متوالی:

$$z - P_t = \sum_{k=t}^m (P_{k+1} - P_k) \quad \text{و} \quad P_{m+1} \stackrel{\text{تعریف}}{=} z$$

- با ترکیب با (۱) و (۲)، روابط زیر بدست می‌آیند:

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} + \sum_{t=1}^m \alpha (z - P_t) \nabla_{\mathbf{w}} P_t \\ &= \mathbf{w} + \sum_{t=1}^m \alpha \sum_{k=t}^m (P_{k+1} - P_k) \nabla_{\mathbf{w}} P_t \end{aligned}$$

محاسبه‌ی افزایشی

- نمایش خطای $z - P_t$ به صورت مجموع تغییرات در پیش‌بینی‌های متوالی:

$$z - P_t = \sum_{k=t}^m (P_{k+1} - P_k) \quad \text{و} \quad P_{m+1} \stackrel{\text{تعریف}}{=} z$$

- با ترکیب با (۱) و (۲)، روابط زیر بدست می‌آیند:

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} + \sum_{t=1}^m \alpha (z - P_t) \nabla_{\mathbf{w}} P_t \\ &= \mathbf{w} + \sum_{t=1}^m \alpha \sum_{k=t}^m (P_{k+1} - P_k) \nabla_{\mathbf{w}} P_t \end{aligned}$$

محاسبه‌ی افزایشی (ادامه)

- با جابجایی دو جمع و تبدیل حدود

$$\begin{aligned}\mathbf{w} &\leftarrow \mathbf{w} + \sum_{k=1}^m \alpha \sum_{t=1}^k (P_{k+1} - P_k) \nabla_{\mathbf{w}} P_t \\ &= \mathbf{w} + \sum_{t=1}^m \alpha (P_{t+1} - P_t) \sum_{k=1}^t \nabla_{\mathbf{w}} P_k\end{aligned}$$

رابطه‌ی بروز رسانی افزایشی وزن‌ها - TD(1)

$$\Delta \mathbf{w}_t = \alpha (P_{t+1} - P_t) \sum_{k=1}^t \nabla_{\mathbf{w}} P_k \quad (4)$$

محاسبه‌ی افزایشی (ادامه)

- با جابجایی دو جمع و تبدیل حدود

$$\begin{aligned}\mathbf{w} &\leftarrow \mathbf{w} + \sum_{k=1}^m \alpha \sum_{t=1}^k (P_{k+1} - P_k) \nabla_{\mathbf{w}} P_t \\ &= \mathbf{w} + \sum_{t=1}^m \alpha (P_{t+1} - P_t) \sum_{k=1}^t \nabla_{\mathbf{w}} P_k\end{aligned}$$

رابطه‌ی بروز رسانی افزایشی وزن‌ها - TD(1)

$$\Delta \mathbf{w}_t = \alpha (P_{t+1} - P_t) \sum_{k=1}^t \nabla_{\mathbf{w}} P_k \quad (4)$$

TD(1)

- اگر بیشترین طول دنباله‌ی مشاهدات M باشد، آن‌گاه الگوریتم TD(1) نیاز به حافظه و سرعت پردازشگر لازم برای یادگیری با نظارت دارد.
- اگر تابع پیش‌بینی را تابع خطی در نظر بگیریم

رابطه‌ی بروز رسانی TD(1) خطی

$$\Delta \mathbf{w}_t = \alpha (P_{t+1} - P_t) \sum_{k=1}^t \mathbf{x}_k \quad P_t = \mathbf{w}^T \mathbf{x}_t$$

الگوریتم TD(1) خطی

ورودی: $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m, z$

خروجی: \mathbf{w}

- 1: $\mathbf{w} \leftarrow \mathbf{w}$ ▷ دنباله‌ی مشاهدات- نتیجه
- 2: $P_t \leftarrow \mathbf{w}^T \mathbf{x}_1$ ▷ بردار وزن‌ها
- 3: $\mathbf{S}_t \leftarrow \mathbf{x}_1$ ▷ بردار وزن‌ها را با مقادیر تصادفی مقداردهی اولیه کن
- 4: for all $\mathbf{x}_t, t = 2, \dots, m$ do ▷ \mathbf{x}_1 مشاهده شد
- 5: $P_{t-1} \leftarrow P_t$ ▷ \mathbf{S}_t همان جمع گرادیان‌هاست
- 6: $P_t \leftarrow \mathbf{w}^T \mathbf{x}_t$ ▷ برای مشاهدات \mathbf{x}_2 تا \mathbf{x}_m
- 7: $\Delta \mathbf{w}_{t-1} \leftarrow \alpha(P_t - P_{t-1})\mathbf{S}_t$ ▷ میزان افزایش بردار وزن‌ها برای مشاهده‌ی قبلی
- 8: $\mathbf{S}_t \leftarrow \mathbf{S}_t + \mathbf{x}_t$
- 9: $\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}_{t-1}$ ▷ بروزرسانی بردار وزن‌ها
- 10: end for
- 11: $\Delta \mathbf{w}_{t-1} \leftarrow \alpha(z - P_t)\mathbf{S}_t$
- 12: $\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}_{t-1}$ ▷ بروزرسانی نهایی بردار وزن‌ها

خانواده‌ی روش‌های یادگیری $TD(\lambda)$

- در $TD(1)$ ، به ازای هر مشاهده، بردار وزن‌ها طوری تغییر می‌کند، که تمام پیش‌بینی‌های گذشته را به یک میزان تغییر می‌دهد
- کلاس روش‌های $TD(\lambda)$ ، پیش‌بینی‌های جدیدتر را بیش‌تر از پیش‌بینی‌های گذشته تغییر می‌دهد
- با وزن‌دهی نمایی با تأخر، تغییر در پیش‌بینی مشاهداتی که در k مرحله پیش انجام شدند متناسب است با λ^k برای $0 \leq \lambda \leq 1$

رابطه‌ی بروزرسانی وزن‌ها $TD(\lambda)$

$$\Delta \mathbf{w}_t = \alpha (P_{t+1} - P_t) \sum_{k=1}^t \lambda^{t-k} \nabla_{\mathbf{w}} P_k \quad (5)$$

خانواده‌ی روش‌های یادگیری $TD(\lambda)$

- در $TD(1)$ ، به ازای هر مشاهده، بردار وزن‌ها طوری تغییر می‌کند، که تمام پیش‌بینی‌های گذشته را به یک میزان تغییر می‌دهد
- کلاس روش‌های $TD(\lambda)$ ، پیش‌بینی‌های جدیدتر را بیش‌تر از پیش‌بینی‌های گذشته تغییر می‌دهد
- با وزن‌دهی نمایی با تأخر، تغییر در پیش‌بینی مشاهداتی که در k مرحله پیش انجام شدند متناسب است با λ^k برای $0 \leq \lambda \leq 1$

رابطه‌ی بروزرسانی وزن‌ها $TD(\lambda)$

$$\Delta \mathbf{w}_t = \alpha (P_{t+1} - P_t) \sum_{k=1}^t \lambda^{t-k} \nabla_{\mathbf{w}} P_k \quad (5)$$

خانواده‌ی روش‌های یادگیری $\text{TD}(\lambda)$ (ادامه)

- مزیت وزندهی نمایی در این است که رابطه را می‌توان به صورت افزایشی محاسبه کرد
- مثلاً اگر مقدار جمع در (۵) را برای مرحله‌ی t ، با s_t نمایش دهیم، در این صورت مقدار s_{t+1} را می‌توان به صورت افزایشی محاسبه کرد

$$\begin{aligned}
 s_{t+1} &= \sum_{k=1}^{t+1} \lambda^{t+1-k} \nabla_w P_k \\
 &= \nabla_w P_{t+1} + \sum_{k=1}^t \lambda^{t+1-k} \nabla_w P_k \\
 &= \nabla_w P_{t+1} + \lambda s_t
 \end{aligned}$$

خانواده‌ی روش‌های یادگیری $TD(\lambda)$ (ادامه)

- به ازای $\lambda < 1$ ، خانواده‌ی روش‌های تفاضل‌های زمانی، بردار وزن‌ها را به شکلی کاملاً متمایز با تمامی روش‌های یادگیری با نظارت تغییر می‌دهد
- به ازای $\lambda = 0$ ، این مطلب مشهودتر است. در $TD(0)$ ، میزان افزایش در بردار وزن‌ها متناسب است با تأثیر آن در آخرین پیش‌بینی (با فرض $\gamma = 1$)

روش بروز رسانی وزن‌ها $TD(0)$

$$\Delta \mathbf{w}_t = \alpha (P_{t+1} - P_t) \nabla_w P_t \quad (۶)$$

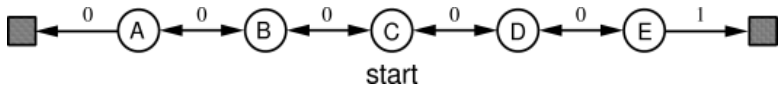
TD(λ) الگوریتم

ورودی: $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m, z$

خروجی: \mathbf{w}

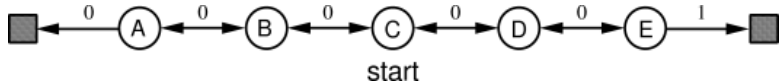
- 1: $\mathbf{w} \leftarrow \mathbf{w}$ ▷ دنباله‌ی مشاهدات- نتیجه
- 2: $P_t \leftarrow P(\mathbf{w}, \mathbf{x}_1)$ ▷ بردار وزن‌ها
- 3: $\mathbf{S} \leftarrow \nabla_{\mathbf{w}} P_t$ ▷ بردار وزن‌ها را با مقادیر دلخواه مقداردهی اولیه کن
- 4: for all $\mathbf{x}_t, t = 2, \dots, m$ do ▷ \mathbf{x}_1 مشاهده شد
- 5: $P_{t-1} \leftarrow P_t$ ▷ \mathbf{S} همان جمع وزن‌دار گرادیان‌هاست
- 6: $P_t \leftarrow P(\mathbf{w}, \mathbf{x}_t)$ ▷ برای مشاهدات \mathbf{x}_2 تا \mathbf{x}_m
- 7: $\Delta \mathbf{w}_{t-1} \leftarrow \alpha(P_t - P_{t-1})\mathbf{S}$ ▷ میزان افزایش بردار وزن‌ها برای مشاهده‌ی قبلی
- 8: $\mathbf{S} \leftarrow \nabla_{\mathbf{w}} P_t + \lambda \mathbf{S}$
- 9: $\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}_{t-1}$ ▷ بروزرسانی بردار وزن‌ها
- 10: end for
- 11: $\Delta \mathbf{w}_{t-1} \leftarrow \alpha(z - P_t)\mathbf{S}$
- 12: $\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}_{t-1}$ ▷ بروزرسانی نهایی بردار وزن‌ها

مثال قدم برداشتن تصادفی (Random Walk)



- تمامی اپیزودها از خانه‌ی مرکز، C، شروع می‌شوند. با احتمال‌های برابر حرکت به سمت راست، یا چپ انجام می‌شود. با رسیدن به هریک از خانه‌هایی که با مربع نمایش داده شده‌اند اپیزود پایان می‌یابد.
- هدف یافتن احتمال این‌که از هرخانه به مربع انتهایی سمت راست برسیم
- این احتمال برای مربع سمت راست مساوی ۱، و برای مربع سمت چپ مساوی صفر است.
- دو نمونه از دنباله‌ی مشاهده-نتیجه: CDCBA0 و CDE1

مثال قدم برداشتن تصادفی (Random Walk) (ادامه)



- برای پیاده‌سازی از $TD(\lambda)$ خطی استفاده شده؛ یعنی

$$P_t = \mathbf{w}^T \mathbf{x}_t \quad \bullet$$

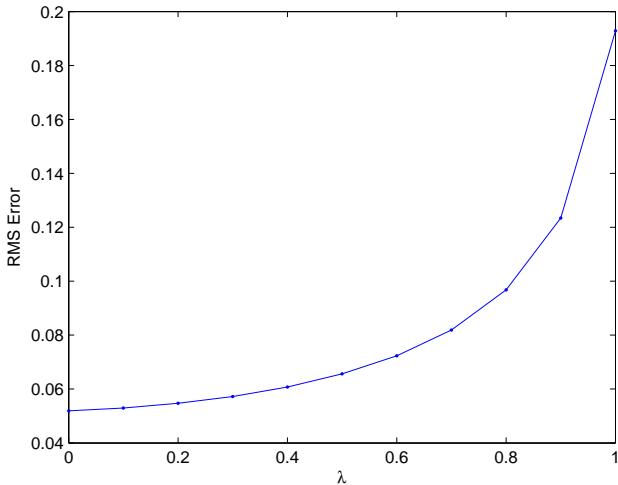
$$\nabla_w P_t = \mathbf{x}_t \quad \bullet$$

- به جای بردارهای مشاهده از بردارهای پایه‌ی \mathbb{R}^5 استفاده شده؛ یعنی

$$\mathbf{x}_C = (0, 0, 1, 0, 0)^T \quad \bullet$$

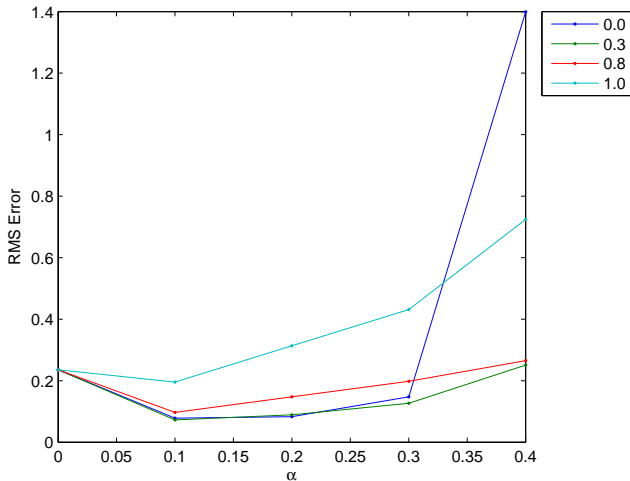
$$\mathbf{x}_E = (0, 0, 0, 0, 1)^T \quad \bullet$$

تأثیر λ بر دقت یادگیری PredictionRandomWalk.m را اجرا کنید

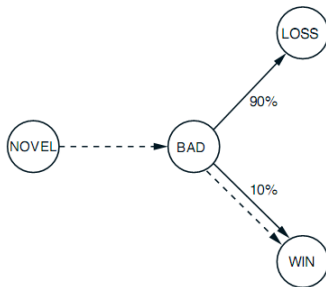


تأثیر α و λ بر دقت یادگیری

PredictionRandomWalkAlphaEffect.m را اجرا کنید



مثال - محیط‌های مارکوف



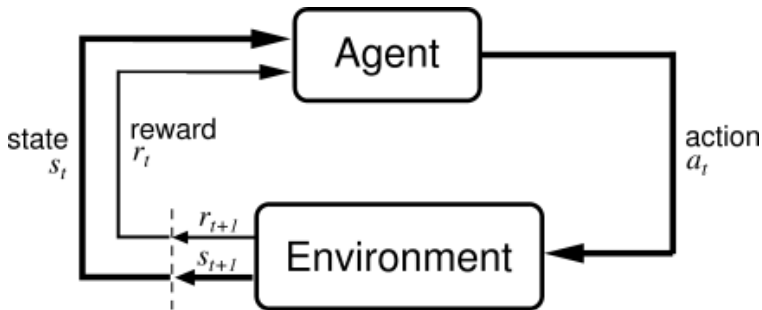
- مزیت روش‌های تفاضل‌های زمانی نسبت به یادگیری با نظارت تنها در سرعت نیست
- در برخی موقعیت‌ها روش‌های تفاضل زمانی جواب صحیح‌تری نسبت به روش‌های یادگیری با نظارت ارائه می‌دهند
- مثلاً در محیط‌های مارکوف

مقدمه‌ای بر یادگیری تقویتی

- زیر شاخه‌ای از یادگیری ماشین، آغاز در نظریه‌ی کنترل
- یادگیری از طریق تعامل با محیط
- اجزای تشکیل دهنده‌ی مسئله‌ی یادگیری تقویتی:
 - عامل
 - محیط
 - کنش
 - پاداش (سیگنال تقویتی)
- هدف عامل: در هر وضعیت، طوری رفتار کند که مجموع پاداش‌های مورد انتظار دریافتی از محیط، در دراز مدّت بیشینه شود

مقدمه‌ای بر یادگیری تقویتی (ادامه)

تعامل عامل با محیط در یادگیری تقویتی



مدل بیشینه کردن پاداش در دراز مدت

مدل افق نامحدود تخفیف یافته

هدف بیشینه کردن امید ریاضی زیر است:

$$E\left(\sum_{t=0}^{\infty} \gamma^t r_t\right) \quad \text{و} \quad 0 \leq \gamma < 1$$

پاداش‌های آینده، با ضریب تخفیف γ به طور هندسی تخفیف داده می‌شوند

- ضریب γ موجب می‌شود که پاداش آنی، وزن بیشتری نسبت به پاداش‌های مورد انتظار در آینده داشته باشد، و در عین حال پاداش‌های آینده نیز نادیده گرفته نشوند

مدل پیشینه‌کردن پاداش در دراز مدت (ادامه)

- هرچه γ به ۱ نزدیک‌تر باشد، الگوریتم یادگیری آینده‌نگرتر می‌شود و هرچه γ به صفر نزدیک‌تر باشد الگوریتم یادگیری حریص‌تر می‌شود
- کاربرد دیگر ضریب γ ، کران‌دار کردن سری فوق است (کاربرد در اثبات قضایای همگرایی)
- انتخاب $\gamma = 1$ ، غیرمجاز است مگر در سناریوهای اپیزودیک با تعداد مراحل محدود

استفاده از روش تفاضل‌های زمانی در یادگیری تقویتی

- یک عامل در هر لحظه‌ی t ، زوج مرتب (\mathbf{x}_t, r_t) را دریافت می‌کند
- بردار مشاهدات عامل در لحظه‌ی t است \mathbf{x}_t
- r_t یک عدد حقیقی معرف پاداش عامل در لحظه‌ی t
- بردار مشاهدات، \mathbf{x}_t می‌تواند صرفاً حالت محیط باشد، یا علاوه بر آن شامل کنش عامل در لحظه‌ی t نیز باشد:

$$\mathbf{x}_t = \mathbf{s}_t \quad \text{یا} \quad \mathbf{x}_t = \langle \mathbf{s}_t, a_t \rangle$$

- هدف: با استفاده از روش تفاضل‌های زمانی، با دریافت دنباله‌های (\mathbf{x}_t, r_t) ، برای $t = 0, 1, \dots$ ، در هر لحظه‌ی t ، پیش‌بینی P_t از کمیت زیر را انجام دهیم:

$$R_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots = \sum_{i=1}^{\infty} \gamma^{i-1} r_{t+i}$$

استفاده از روش تفاضل‌های زمانی در یادگیری تقویتی

- یک عامل در هر لحظه‌ی t ، زوج مرتب (\mathbf{x}_t, r_t) را دریافت می‌کند
- بردار مشاهدات عامل در لحظه‌ی t است \mathbf{x}_t
- r_t یک عدد حقیقی معرف پاداش عامل در لحظه‌ی t
- بردار مشاهدات، \mathbf{x}_t می‌تواند صرفاً حالت محیط باشد، یا علاوه بر آن شامل کنش عامل در لحظه‌ی t نیز باشد:

$$\mathbf{x}_t = \mathbf{s}_t \quad \text{یا} \quad \mathbf{x}_t = \langle \mathbf{s}_t, a_t \rangle$$

- هدف: با استفاده از روش تفاضل‌های زمانی، با دریافت دنباله‌های (\mathbf{x}_t, r_t) ، برای $t = 0, 1, \dots$ ، در هر لحظه‌ی t ، پیش‌بینی P_t از کمیت زیر را انجام دهیم:

$$R_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots = \sum_{i=1}^{\infty} \gamma^{i-1} r_{t+i}$$

استفاده از روش تفاضل‌های زمانی در یادگیری تقویتی (ادامه)

- در روش تفاضل‌های زمانی، مقادیر بردار وزن‌ها یادگرفته می‌شوند.
- در صورت استفاده از روش تفاضل‌های زمانی خطی، طول بردار وزن‌ها برابر است با طول بردار مشاهدات.
- در یادگیری تقویتی در صورتی که بردار مشاهدات صرفاً شامل وضعیت محیط باشد، به بردار وزن‌ها می‌گویند تابع ارزش و آن را با V نمایش می‌دهند
- در صورتی که بردار مشاهدات علاوه بر وضعیت محیط، شامل کنش نیز باشد، آن را با Q نمایش می‌دهند

استفاده از روش تفاضل‌های زمانی در یادگیری تقویتی (ادامه)

- در هر لحظه محیط دارای وضعیت s است که s یک عضو از فضای حالات S می‌باشد
- با فرض گسسته و محدود بودن فضای حالات، به ازای هر حالت در فضای حالات، یک درایه در بردار وزن‌ها در نظر بگیریم. در این صورت بدست آوردن ارزش هر حالت به سادگی جستجو در یک جدول خواهد بود.

$$|V| = n(S)$$

- بردار مشاهده را طوری در نظر می‌گیریم که فقط مشخص کند، کدام حالت در فضای حالات مشاهده شده است، بنابراین هر بردار مشاهده، یک پایه برای فضای $\mathbb{R}^n(S)$ است

- با استفاده از TD(0) خطی:

$$\Delta V_t = \alpha(P_{t+1} - P_t) \nabla_V P_t$$

$$P_t = V^T \mathbf{x}_t \Rightarrow \nabla_V P_t = \mathbf{x}_t$$

- اگر وضعیت محیط در لحظه‌ی t ، همان وضعیت شماره‌ی s ام باشد در فضای حالت S ، تنها درایه‌ی s ام بردار \mathbf{x}_t ، مساوی یک است. بنابراین:

$$V_{t+1}(s) \leftarrow V_t(s) + \alpha(P_{t+1} - P_t)$$

- که در آن منظور از $V(s)$ ، درایه‌ی s ام بردار V است؛ و منظور از V_t ، نسخه‌ای از بردار V است که در زمان t ، در دسترس بود.

- اگر $\gamma = 0$ ؛ آن‌گاه $P_{t+1} = r_{t+1}$ ، و همچنین داریم:

$$P_t = P_t(\mathbf{x}_t) = V^T \mathbf{x}_t = V(s)$$

- تابع پیش‌بینی همان بردار ارزش‌ها است، و با یادگیری آن به طور خودبه‌خود تابع پیش‌بینی را نیز یاد گرفته‌ایم

$$V_{t+1}(s) \leftarrow V_t(s) + \alpha(r_{t+1} - V_t(s)) \quad (7)$$

- s وضعیت محیط در لحظه‌ی t است
- این الگوریتم در هر لحظه‌ی t ، مقادیر بردار ارزش‌ها در گام زمانی قبلی را بروز می‌کند

• اگر $\gamma \neq 0$ ؛ آن‌گاه

$$\begin{aligned} P_{t+1} = R_t &= r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots \\ &= r_{t+1} + \gamma (r_{t+2} + \gamma r_{t+3} + \dots) \\ &= r_{t+1} + \gamma R_{t+1} \end{aligned}$$

- اگر وضعیت بعد از s را با s' نشان دهیم، آن‌گاه می‌دانیم که $V_t(s')$ تخمینی برای R_{t+1} است در لحظه‌ی t
- بنابراین می‌توانیم R_t را با $r_{t+1} + \gamma V_t(s')$ تخمین بزنیم

$$V_{t+1}(s) \leftarrow V_t(s) + \alpha (r_{t+1} + \gamma V_t(s') - V_t(s)) \quad (8)$$

آشنایی با چند تعریف

سیاست عامل

نگاشتی که مشخص می‌کند، عامل در هر وضعیت، ممکن است چه کنشی را انجام دهد را سیاست^۱ عامل می‌گویند، و آن را با π نشان می‌دهند

Policy

سیاست اپسیلون-حریصانه

سیاستی که در آن عامل به احتمال ϵ کنش تصادفی را انجام می‌دهد، و در باقی مواقع کنشی را انجام می‌دهد که در تابع ارزش بیشترین ارزش را بخود اختصاص داده. هرچه اپسیلون بیشتر باشد عامل زمان بیشتری را صرف تجربه کردن محیط ناشناخته اطراف خود می‌کند.^۲

ϵ -greedy

آشنایی با چند تعریف

سیاست عامل

نگاشتی که مشخص می‌کند، عامل در هر وضعیت، ممکن است چه کنشی را انجام دهد را سیاست^۱ عامل می‌گویند، و آن را با π نشان می‌دهند

Policy

سیاست اپسیلون-حریصانه

سیاستی که در آن عامل به احتمال ϵ کنش تصادفی را انجام می‌دهد، و در باقی مواقع کنشی را انجام می‌دهد که در تابع ارزش بیشترین ارزش را بخود اختصاص داده. هرچه اپسیلون بیشتر باشد عامل زمان بیشتری را صرف تجربه کردن محیط ناشناخته اطراف خود می‌کند.^۲

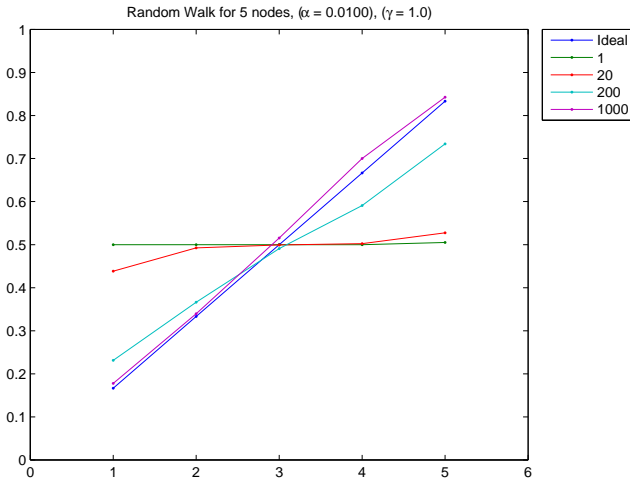
ϵ -greedy

استفاده از روش TD(0) خطی در یادگیری تقویتی

- 1: Initialize V arbitrarily ▷ بردار ارزش‌ها را به‌طور دلخواه مقداردهی اولیه کن.
- 2: **for all** episode **do** ▷ برای هر اپیزود
- 3: Initialize s
- 4: **for all** step of episode **do**
- 5: $a \leftarrow$ action given by π for s
- 6: Take action a ; ▷ کنش بدست آمده از طریق π را انجام بده
- 7: observe r , and next state s' ▷ پاداش و وضعیت بعدی محیط را مشاهده کن
- 8: $V(s) \leftarrow V(s) + \alpha [r + \gamma V(s') - V(s)]$
- 9: $s \leftarrow s'$
- 10: **end for**
- 11: **end for**

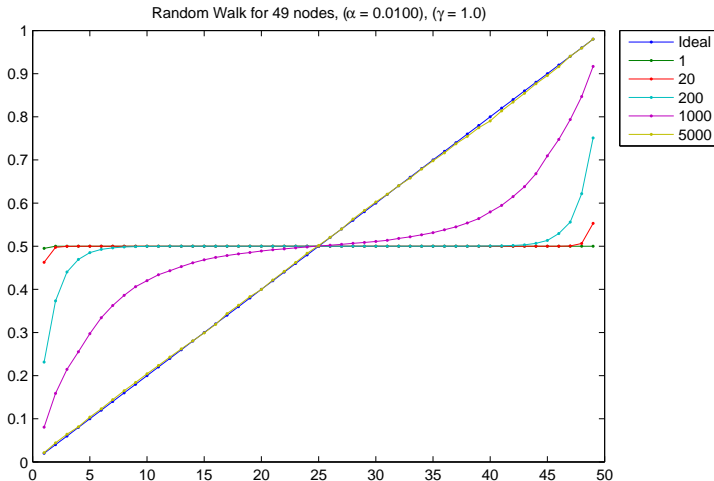
مثال: قدم برداشتن تصادفی Random Walk

RLRandomWalk.m را اجرا کنید



مثال: قدم برداشتن تصادفی Random Walk

RLRandomWalk.m را اجرا کنید



الگوریتم‌های یادگیری مبتنی بر سیاست، و مستقل از سیاست

الگوریتم‌های مبتنی بر سیاست

الگوریتم‌هایی که بهترین پاسخ را برای سیاستی که عامل در پیش گرفته ارائه می‌دهند. اگر عامل سیاست خود را تغییر دهد، پاسخ الگوریتم نیز ممکن است تغییر کند.^۳ مانند الگوریتم یادگیری SARSA.

on-policy

الگوریتم‌های مستقل از سیاست

الگوریتم‌هایی که بهترین پاسخ ممکن را برای محیط ارائه می‌دهند. اگر عامل سیاست خود را تغییر دهد، پاسخ الگوریتم در نهایت تغییر نمی‌کند.^۴ مانند الگوریتم یادگیری Q-Learning

off-policy

الگوریتم‌های یادگیری مبتنی بر سیاست، و مستقل از سیاست

الگوریتم‌های مبتنی بر سیاست

الگوریتم‌هایی که بهترین پاسخ را برای سیاستی که عامل در پیش گرفته ارائه می‌دهند. اگر عامل سیاست خود را تغییر دهد، پاسخ الگوریتم نیز ممکن است تغییر کند.^۳ مانند الگوریتم یادگیری SARSA.

on-policy

الگوریتم‌های مستقل از سیاست

الگوریتم‌هایی که بهترین پاسخ ممکن را برای محیط ارائه می‌دهند. اگر عامل سیاست خود را تغییر دهد، پاسخ الگوریتم در نهایت تغییر نمی‌کند.^۴ مانند الگوریتم یادگیری Q-Learning

off-policy

الگوریتم یادگیری مبتنی بر سیاست SARSA

- 1: Initialize $Q(s, a)$ arbitrarily
- 2: **for all** episode **do**
- 3: Initialize s
- 4: Choose a from s using policy derived from Q (e.g., ϵ -greedy)
- 5: **for all** step of episode **do**
- 6: Take action a , observe r, s'
- 7: Choose a' from s' using policy derived from Q (e.g., ϵ -greedy)
- 8: $Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma Q(s', a') - Q(s, a)]$
- 9: $s \leftarrow s'; a \leftarrow a'$
- 10: **end for**
- 11: **end for**

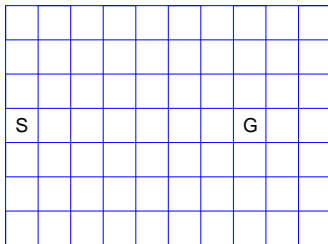
الگوریتم یادگیری مستقل از سیاست Q-Learning

- 1: Initialize $Q(s, a)$ arbitrarily
- 2: **for all** episode **do**
- 3: Initialize s
- 4: **for all** step of episode **do**
- 5: Choose a from s using policy derived from Q (e.g., ϵ -greedy)
- 6: Take action a , observe r, s'
- 7: $Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$
- 8: $s \leftarrow s'$;
- 9: **end for**
- 10: **end for**

مثال: Grid World

شمای کلی محیط

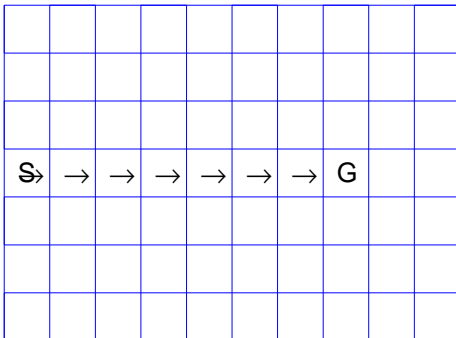
- هدف: آغاز حرکت از S و رسیدن به G
- اعمال ممکن: حرکت در چهار جهت اصلی
- الگوریتم‌های SARSA و Q-Learning یک پاسخ را ارائه می‌دهند



مثال: Grid World

GridWorldsDemo('sg-small-sarsa')

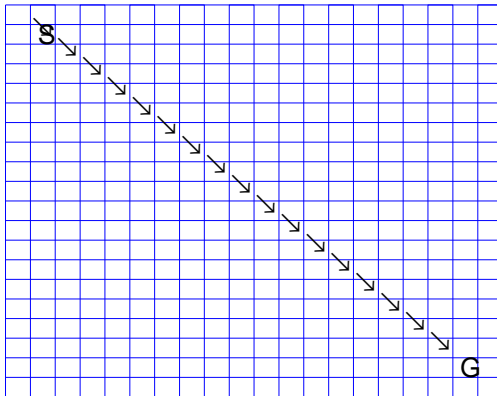
Simple grid-world SARSA – episode 700 – ($\epsilon: 0.100$), ($\alpha = 0.1000$), ($\gamma = 0.9$)



مثال: Grid World

GridWorldsDemo('sg-big-king-sarsa') را اجرا کنید

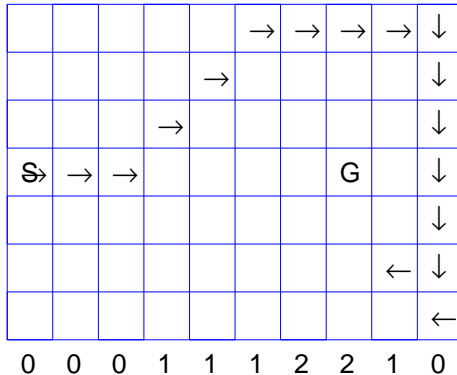
Simple grid-world SARSA – episode 4000 – ($\epsilon: 0.100$), ($\alpha = 0.1000$), ($\gamma = 0.9$)



مثال: Windy Grid World

GridWorldsDemo('wg-small-sarsa') را اجرا کنید

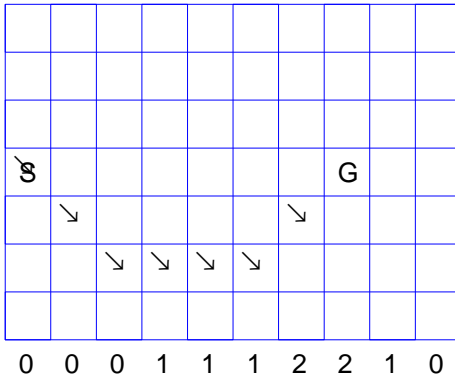
Windy grid-world SARSA – episode 1000 – ($\epsilon: 0.100$), ($\alpha = 0.1000$), ($\gamma = 0.9$)



مثال: Windy Grid World

GridWorldsDemo('wg-small-king-sarsa')

Windy grid-world SARSA – episode 2000 – ($\epsilon: 0.100$), ($\alpha = 0.1000$), ($\gamma = 0.9$)



مثال: Cliff Walking

GridWorldsDemo('cliff-small-sarsa')

Cliff Walking SARSA – episode 15000 – ($\epsilon: 0.100$), ($\alpha = 0.0100$), ($\gamma = 0.9$)

→	→	→	→	→	→	→	→	→	↓
↑									↓
\$	C	C	C	C	C	C	C	C	G

با تشکر از شما

?