

400104964

نم
ایلمین ریست

(الف) در یادگیری، تغییر در $distribution$ شل است. $state-action$ ها

شوند و مقدار تخصیص زده شده برای دیگر $action$ ها متوسط مدل بسیار بیشتر از مقدار واقعی $action$ ها و این

عمل $propagate$ شود به این مقدار با عبارت $E[Q]$ و Q ای که نزدیک

است را کاهش دهیم و از این شل $action$ های دیگر را $action$ های دیگر را افزایش دهیم و این $action$ ها

به $action$ ها و دیگر $action$ ها را افزایش دهد و در نهایت $action$ ها به $action$ های دیگر $action$ های دیگر

عبارت $E[Q]$ را از این $action$ ها $action$ های دیگر را $action$ های دیگر را $action$ های دیگر را $action$ های دیگر را

$distribution$ shift ندارند.

(ب) ابتدا بهینه سازی داخلی زیر را حل می کنیم.

$$\max_p E_{s \sim D} [E_{a \sim p(s,a)} [Q(s,a)] - H(p(\cdot|s))]$$

$$s.t. \forall s: \sum_a p(a|s) = 1$$

می توانیم برای هر s آن را به عنوان $action$ حل کنیم پس آن را بهینه سازی به نرم زیر تبدیل می شود:

$$\max_p E_{a \sim p(s,a)} [Q(s,a)] + H(p(\cdot|s))$$

$$s.t. \sum_a p(a|s) = 1$$

ساده را به سبب تاثیرش بر این عمل می‌بینیم.

$$L(\mu, \lambda) = \sum_a p(a|s) (Q(s, a) - \log p(a|s)) + \lambda \left(\sum_a p(a|s) - 1 \right)$$

$$\frac{\partial L(\mu, \lambda)}{\partial p(a|s)} = \sum_a \alpha Q(s, a) - \log p(a|s) - 1 + \lambda = 0$$

$$\Rightarrow p(a|s) = \frac{e^{-1+\lambda} \sum_a \alpha Q(s, a)}{e^{-1+\lambda} \sum_a \alpha Q(s, a)} = 1$$

$$\Rightarrow \frac{e^{-1+\lambda}}{\sum_a \alpha Q(s, a)} = 1 \Rightarrow p(a|s) = \frac{e^{\alpha Q(s, a)}}{\sum_a e^{\alpha Q(s, a)}}$$

در نتیجه، ساده‌ترین به شکل زیر تبدیل می‌شود:

$$\min_Q \alpha E_{s \sim D} \left[\sum_a p(a|s) (Q(s, a) - Q(s, a) + \log \sum_a e^{\alpha Q(s, a)}) \right] - E_{a \sim \pi_B} [Q(s, a)] + \frac{1}{2} E_{s, a, s'} [(Q - \hat{\pi}_K Q^K)^2]$$

$$\Rightarrow \min_Q \alpha E_{s \sim D} [\log \sum_a e^{\alpha Q(s, a)}] - E_{a \sim \pi_B} [Q(s, a)] + \frac{1}{2} E_{s, a, s'} [(Q - \hat{\pi}_K Q^K)^2]$$

ساده به سبب این ارزش می‌بینیم، ساده‌ترین را حل کنیم و نرم‌ترین را ساده‌ترین این کار مانع از غرور

زبان در هر دو Q های نزدیک می‌شود و این را می‌توانیم به سبب این کار مانع از غرور

conservative policy

(ج) سانه به شکل زیر تبدیل می شود:

$$p^* = \arg \min_p E_{S \sim D, a \sim p(a|S)} [Q(s, a)] + R(p)$$

$$\min_Q \alpha E_{S \sim D, a \sim p^*(a|S)} [Q(s, a)] + \frac{1}{2} E_{S, a, S' \sim D} [(Q - \hat{B}^{\pi_K} \hat{Q}^K)^2] + R(p^*)$$

برای هر S سانه را جدا حل می کنیم:

$$\min_Q \alpha E_{a \sim p^*(a|S)} [Q(s, a)] + \frac{1}{2} \int_a \hat{\pi}_B(a|S) (Q(s, a) - \hat{B}^{\pi_K} \hat{Q}^K)^2$$

مشتق سانه را برابر 0 قرار می دهیم و داریم:

$$\alpha p^*(a|S) + \hat{\pi}(a|S) (\hat{Q}^{K+1}(s, a) - \hat{B}^{\pi_K} \hat{Q}^K(s, a)) = 0$$

$$\Rightarrow \hat{Q}^{K+1}(s, a) = \hat{B}^{\pi_K} \hat{Q}^K(s, a) - \alpha \frac{p^*(a|S)}{\hat{\pi}_B(a|S)}$$

(د) با توجه به سنجش مبتدیان و دران نشانه شده داریم (fixed point):

$$\hat{Q}^{\pi}(s, a) = \hat{B}^{\pi} \hat{Q}^{\pi}(s, a) - \alpha \frac{p^*(a|S) - C(s, a) \hat{B}^{\pi} \hat{Q}^{\pi}(s, a) - R(s, a)}{\hat{\pi}_B(a|S)}$$

$$\hat{Q}^{\pi}(s, a) \leq \hat{B}^{\pi} \hat{Q}^{\pi}(s, a) + C_{\delta}(s, a) - \alpha \frac{p^*(a|S)}{\hat{\pi}_B(a|S)} = R(s, a) + \gamma \hat{P}^{\pi} \hat{Q}^{\pi}(s, a) + C_{\delta}(s, a) - \alpha \frac{p^*(a|S)}{\hat{\pi}_B(a|S)}$$

$$\Rightarrow (I - \gamma \hat{P}^{\pi}) \hat{Q}^{\pi}(s, a) \leq R(s, a) + C_{\delta}(s, a) - \alpha \frac{p^*(a|S)}{\hat{\pi}_B(a|S)}$$

$$\Rightarrow \hat{Q}^{\pi}(s, a) \leq \underbrace{(I - \gamma \hat{P}^{\pi})^{-1} R(s, a)}_{\hat{Q}^{\pi}(s, a)} + \underbrace{E_{\delta}}_{\hat{Q}^{\pi}(s, a)} (I - \gamma \hat{P}^{\pi})^{-1} C_{\delta}(s, a) - (I - \gamma \hat{P}^{\pi})^{-1} \alpha \frac{p^*(a|S)}{\hat{\pi}_B(a|S)}$$

$$\Rightarrow \hat{Q}^{\pi}(s, a) \leq \hat{Q}^{\pi}(s, a) + (I - \gamma \hat{P}^{\pi})^{-1} C_{\delta}(s, a) - \alpha (I - \gamma \hat{P}^{\pi})^{-1} \frac{p^*(a|S)}{\hat{\pi}_B(a|S)}$$

(2) الف) BC سعی می کند با لایه ها راه صورت تصمیم از آنسوی فضای مشخص یاد بگیرد که این باعث می شود که توانایی generalization آن کاهش یابد و بسیار به نسبت داده ها حساس شود. IRL سعی می کند ابتدا برآورد مناسبی را به دست آورد که این سبب می شود تا قدرت تعمیم ما بالاتر رود و نسبت به نویزها مقاوم تر شود. این باعث می شود آموختن های IRL انعطاف پذیرتر در سازگار تر شوند و انبوه های نسبت رفتارهای expert را درک کند نه اینکه صرفاً صرفاً را تقلید کند. همچنین IRL می تواند با لایه های sub optimal را بپذیرد.

ب) در این مثال ما یک نرم افزار داریم که می توانیم آن را با یک سیستم دیگر به انجام کارهای دیگر این باعث می شود تا از بین تمام با لایه های ممکن که می تواند رفتار expert را توجیه کند این را برگزینیم. بهترین و به جامع بهترین رنדר من را دارند. در این حالت ممکن می شود که با لایه ما تعداد اعیان رنדר است و همچنین چنین رفتار expert است که این از overfit شدن جلوگیری می کند. همچنین exploration را تقویت می کند.

و باعث مقاوم شدن آن نسبت به variation های مختلف دنیا می شود. به افزایش generalization

$$\max_{\theta} \sum_{\tau \sim D} \log p(\tau | \theta) - \lambda \| E_{\pi_{\theta}} [\phi(s, a)] - E_{\pi_E} [\phi(s, a)] \| \quad p(\tau | \pi) = \frac{1}{Z(\theta)} \exp \left(\sum_{t=0}^T R(s_t, a_t) \right)$$

$$J(\pi_E, \pi) = E_{\pi_E} [\log(1 - D(s, a))] \quad J(\pi, \pi) = \lambda H(\pi) - E_{\pi} [\log(D(s, a))] \quad \text{ج)}$$

D تلاش می کند تا مایه بین با لایه expert و با لایه در حال یادگیری را افزایش دهد. در حقیقت تلاش می کند

بهترین را از خود جدا کند. در آن با لایه در حال یادگیری در بهترین حالت خود نسبت به با لایه expert باشد.

در این حال با بیان کردن نسبت به π داریم بالبر مان را در بهترین حالت عصبی جدید در تقسیم به این
 همان نقش adversarial آن است. حال با بیان کردن $H(\pi)$ - همان با بیان کردن $H(\pi)$ در عصبی
 داریم تقسیم به π به بالبر تا جای ممکن میزنیم است و شکل نقش را حل میزنیم.

(د) در این های نقش میزنیم داشته $\max Ent IKL$ با π داشته نیاز به محاسبه میزنیم کرده و محدودیت
 در انتخاب میزنیم $generalization, scalability$ نسبت به داده های دیده شده را دارند. $GAIL$ است
 از شبکه های عصبی عمیق یادگیری میزنیم ها از داده های عام $state-action$ این مشکلات را برطرف میزنیم
 به داده های دست ساز اضافه را حاضر میزنیم این باعث افزایش مقیاس پذیری و تقسیم به π را برطرف میزنیم
 adversarial روش $GAIL$ تقسیم جامع تری از expert را فراهم میزنیم و محدودیت های روش های قبلی را برطرف میزنیم
 یادگیری adversarial هم به $discriminator$ را جدید در تقسیم

3 الف) روش‌هایی که تنها از معیارهای model-based استفاده می‌کنند، اغلب از عدم دقت مدل رنج می‌برند.

در تراننده سیاست‌های optimal type مصرف‌شده از سوی دیگر روش‌هایی که حامله تعاملات دنیای واقعی

مفکند اند. می‌تواند بر هزینه، زمان و ریسک و خطرات باشد به خصوص در محیط‌هایی که پیچیده یا ناامن هستند. روش‌های ترکیبی

مانند MBPO نقاط قوت هر دو رویکرد را با هم ترکیب می‌کند. اما از جنبه‌های ساختاری، برای بهینه‌سازی برای ترکیب داده‌های

آموزش اضافی استفاده می‌کند. sample efficiency را افزایش می‌دهد و همچنین تعاملات دنیای واقعی را برای

تصمیم‌گیری در تعاملات دنیای واقعی به کار می‌آید. به مصرف سیاست‌های مقاوم‌تر و خردتر می‌شود. این رویکرد ترکیبی

معادل ایجاد می‌کند که از طریق مدل‌ها استفاده می‌کند. در حالی که با بازخورد واقعی، قابلیت اطمینان را تضمین می‌کند.

ب) MBPO جنبه‌های بهینه‌سازی بهینه‌سازی را با استفاده از مدل دنیای واقعی برای ترکیب rollout ها ترکیب می‌کند.

که پس از مدل جنبه‌های سازی داده‌های دنیای واقعی ترکیب می‌شود. rollout مدل‌های مختلف برای rollout

ها همان معادل بین بایاس و واریانس است. مدل کمتر سبب کاهش خطاهای تصحیح می‌شود اما ممکن است

نتراننده اشاره‌ها را explore کند که سبب افزایش بایاس و کاهش واریانس می‌شود. مدل بیشتر سبب می‌شود.

آ داده‌های بیشتری را مشاهده کنیم و پس خطاهای تصحیح مدل را افزایش می‌دهد و سبب کاهش دقت مدل می‌شود.

MBPO به طور مویا خود را تصحیح می‌دهد تا از مزایای هر دو به صورت همزمان استفاده کند. سه مرحله 3 مرحله

(ج) MBPO عملی تر از حد صحت های افلاخ استفاده شود زیرا نیاز به حداقل با صحت واقعی دارد و در یادگیری افلاخ

ما این دسترس را فراهم این باعث می شود فضای مدل ما به صورت نرینه ای افزایش یابد و منفرجه عمل در صحت

سیاست شود و نیاز به شلات QOD را حاصل شد

(د) دلیل از روش های عملی برای اندازه گیری عدم تقصیت به این صورت است که از مدل های اساسی استفاده

کنیم. در این حالت چند مدل روی یک دنیا ترین می شود و نتایج آنها با هم ترکیب و مقایسه می شود و این بین

این مدل ها به ما میزان عدم تقصیت را می دهد پس سیاست ها از ورود به مناطق با عدم تقصیت بالا منفرجه می شود

و عمل در این مورد قابل اعتدالتری تقصیر می شود

(ه) COMBO استراتژی های محافظه کارانه را با ایده های Dyna و QAL ترکیب می کند تا جانش های یادگیری افلاخ

را حل کند روش Dyna از مدل برای شبیه سازی تجربه ها برای ادیت کردن یادگیری استفاده می کند و بار روش QAL

که به صورت Conservative یادگیری استرس انجام می دهد ادغام می شود (اندازه های) که تقصیر یادگیری من شده و سیاست

بسیار کمتر از حد می شود. این ترکیب اطمینان می دهد که مدل محافظه کارانه در من مانده و over estimate و بر حاد

امن با داده ها خاص عملی می شود و تقریباً شبیه سازی شده را برای بهبود استقامت در عمل سیاست به طار می شود

این روش ترکیب به طار مؤثر تعادل بین التاب و محافظه کاری را برای بهبود نتایج یادگیری فراهم می کند

$$Q_{new}(s,a) = \min_i Q_{old}(s,a) + r + \gamma E[Q_{target}(s,a)]$$

sam

(د) مدل‌های معادله‌ای مانند COMBO، MOREL، اینتر و «robustness» را با بررسی از آن‌های رگرسیونی

جبریم کردن است. استی‌های غیر تقعر اولویت بندی می‌کند و برای این کار از روش‌های معادله‌ای استفاده می‌کند.

COMBO از ایندیکس پائین معادله‌ای استفاده می‌کند تا از آن‌های Self معین شود در حالی که MOREL

به صورت خاص خاصه با عدم تعقیب بالا را جبریم می‌کند.

COMBO به نقاط مرتب: (۱) الیورین معادل ^{معادله‌ای} sample efficient (۲) جبریمی (۳) over estimation

و لورالین کار در سبب جبریمی از آن‌های رگرسیونی به سبب جبریمی ضعیف می‌شوند.

نقاط ضعف: (۱) احسان معادله‌ای بودن بیش از حد به سبب جبریمی می‌شود (۲) بی‌صدایی پیاده‌سازی

MOREL به نقاط مرتب: (۱) جبریمی کردن عدم تعقیب به صورت مستقیم و صریح (۲) تقنین ایندیکس بالا

نقاط ضعف: (۱) جبریمی معادله‌ای کارانه نه عملی است به جواب ^{optimal} sub در مناطق کسری می‌رسد

(۲) زیاده‌روی در جفا استفاده از آن اصل (۳) عملی است در صورت پیدا کردن مناطق با برابری بالا را از دست بدهد

اثبات ها با استفاده از نابرابری های هوفدینگ و Tor Lattimore و Csaba Szepesvári اثبات ها را می بینیم

$$S_n = \sum_{i=1}^n X_i$$

$$a_i \leq X_i \leq b_i$$

$$\bar{X}_n = \frac{S_n}{n}$$

الف) اثبات برای Hoeffding

4

$$P(E[S_n] - S_n \geq t) \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right) \quad X_i \rightarrow i.i.d$$

$$P(E[X] - \bar{X}_n \geq \frac{t}{n}) \leq \exp\left(-\frac{2t^2}{4n}\right)$$

در صورتی که $a_i = -1$ و $b_i = 1$ می شود داریم:

$$\Rightarrow P(E[X] - \bar{X}_n \geq \epsilon) \leq \exp\left(-\frac{n\epsilon^2}{2}\right) \Rightarrow P(\mu \geq \hat{\mu} + \epsilon) \leq \exp\left(-\frac{n\epsilon^2}{2}\right)$$

$$P(\mu \geq \hat{\mu} + \epsilon) \leq \exp\left(-\frac{n \frac{2 \log \frac{1}{\delta}}{2}}{2}\right) = \delta$$

در جای ϵ قرار دهیم $\sqrt{\frac{2 \log \frac{1}{\delta}}{n}}$ داریم:

$$\Rightarrow P(\mu \geq \hat{\mu} + \sqrt{\frac{2 \log \frac{1}{\delta}}{n}}) \leq \delta$$

$$\mu \leq \hat{\mu} + \sqrt{\frac{2 \log \frac{1}{\delta}}{n}}$$

در نتیجه با احتمال حداقل $1 - \delta$ داریم:

و در آن بالا برای بیان این می توانیم هر عمل است

ب) فرض می کنیم K بازه و n تمام زمان داریم و بین n شش از شش ها میانه در نظر می گیریم که بازوی اول همیشه است

یعنی $\mu = \mu^*$ فرض می کنیم دانش ما در t این اجرای بازی را با $X_{t,i}$ نشان می دهیم و همچنین می توانیم

تقریب می دانش حاصل از انجام این بازی m بار اجرای آن را با $\hat{\mu}_{i,m}$ نشان می دهیم که برابر است با $\frac{\sum_{t=1}^m X_{t,i}}{m}$

برای استنباط باید برای $\frac{E[T_i(n)]}{n}$ دست آوریم ابتدا باید دید که G را به شکل زیر تعریف کنیم:

$$G_i = \left\{ \mu_i < \min_{t \in [n]} UCB_i(t, \delta) \right\} \cap \left\{ \hat{\mu}_{i,u_i} < \sqrt{\frac{2 \log \frac{1}{\delta}}{u_i}} + \mu_i \right\}$$

هرگاه G_i اتفاق بیفتد

دو مورد را اثبات می‌کنیم. ① در صورتی که G_j رخ دهد، بازی را حذف می‌کنیم و بازی را انتخاب شده است.

• اثبات: از مرحله آخر انتخاب شده است. فرض می‌کنیم $u_i > T_i(n)$ باشد. در این صورت به دلیل اینکه بازی را

نمی‌توانیم از u_i بازی انتخاب شده است باید $t \in [n]$ وجود داشته باشد که $u_i = T_i(t)$ و $A_t = i$. در نتیجه داریم:

$$UCB_i(t-1, \delta) = \hat{\mu}_i(t-1) + \sqrt{\frac{2 \log \frac{1}{\delta}}{T_i(t-1)}} = \hat{\mu}_i + \sqrt{\frac{2 \log \frac{1}{\delta}}{u_i}} < \mu_i < UCB_i(t-1, \delta)$$

$$\Rightarrow A_t = \arg \max_j UCB_j(t-1, \delta) \neq i \quad \checkmark$$

② G_j احتمال رخ دادن در هر

• اثبات: این به این معنیست که G_j به احتمال رخ دهد. صحت تقریب داریم:

$$G_j^c = \{ \mu_i \geq \min_{t \in [n]} UCB_i(t, \delta) \} \cup \{ \hat{\mu}_i + \sqrt{\frac{2 \log \frac{1}{\delta}}{u_i}} \geq \mu_i \}$$

صحت Union bound

$$P(\mu_i \geq \min_{t \in [n]} UCB_i(t, \delta)) = P(\bigcup_{t \in [n]} \mu_i \geq \hat{\mu}_i + \sqrt{\frac{2 \log \frac{1}{\delta}}{T_i(t)}}) \leq P(\bigcup_{t \in [n]} \mu_i \geq \hat{\mu}_i + \sqrt{\frac{2 \log \frac{1}{\delta}}{t}})$$

$$\leq \sum_{t=1}^n P(\mu_i \geq \hat{\mu}_i + \sqrt{\frac{2 \log \frac{1}{\delta}}{t}}) \leq n \delta$$

$$\frac{1}{2} \Delta_i \geq \sqrt{\frac{2 \log \frac{1}{\delta}}{u_i}} \quad \text{حال فرض می‌کنیم } u_i \text{ را به صورتی انتخاب می‌کنیم که (صحت صورتی) } (\Delta_i \geq \sqrt{\frac{2 \log \frac{1}{\delta}}{u_i}})$$

$$\Rightarrow P(\hat{\mu}_i + \sqrt{\frac{2 \log \frac{1}{\delta}}{u_i}} \geq \mu_i) = P(\hat{\mu}_i - \mu_i \geq \Delta_i - \sqrt{\frac{2 \log \frac{1}{\delta}}{u_i}}) \leq P(\hat{\mu}_i - \mu_i \geq \frac{1}{2} \Delta_i) \stackrel{\text{Hoeffding}}{\leq} \exp\left(-\frac{u_i \Delta_i^2}{8}\right)$$

$$P(G_j^c) \leq n \delta + \exp\left(-\frac{u_i \Delta_i^2}{8}\right) \quad \checkmark$$

حال صورت ساده را ثابت می‌کنیم. برای $E[T_i(n)]$ داریم:

$$E[T_i(n)] \leq \overbrace{E[I(G_i)T_i(n)]}^{\leq u_i} + \overbrace{E[I(G_i^c)T_i(n)]}^{\leq np(G_i^c)} \leq u_i + np(G_i^c) \leq u_i + n(n\delta + \exp(-\frac{u_i \Delta_i^2}{8}))$$

اگر u_i را برابر $\left\lceil \frac{16 \log \frac{1}{\delta}}{\Delta_i^2} \right\rceil$ قرار دهیم (دسته بندی در * صریح می‌کند) داریم:

$$E[T_i(n)] \leq \left\lceil \frac{16 \log n}{\Delta_i^2} \right\rceil + n(n\delta + \delta) \stackrel{\delta = \frac{1}{n^2}}{=} \frac{16 \log n}{\Delta_i^2} + 1 + 2 = \frac{16 \log n}{\Delta_i^2} + 3 \quad \checkmark$$

$$R_n = \sum_{i=1}^K \Delta_i E[T_i(n)]$$

(ج) ابتدا اثبات می‌کنیم که R_n از رابطه رویه درجه دست می‌آید:

$$R_n = np^* E\left[\sum_{t=1}^n X_t\right]$$

* اثبات: طبق تعریف regret داریم:

در هر زمان t داریم $\sum_{i=1}^K I[A_t = i] = 1$ درستی:

$$R_n = np^* \sum_{t=1}^n \sum_{i=1}^K E[I[A_t = i] X_t] = \sum_{i=1}^K \sum_{t=1}^n E[(p^* - X_t) I[A_t = i]]$$

$$= \sum_{i=1}^K \sum_{t=1}^n E[E[(p^* - X_t) I[A_t = i] | A_t]]$$

طبق قانون امید ریاضی داریم:

$$\Rightarrow = \sum_{i=1}^K \sum_{t=1}^n E[I[A_t = i] E[(p^* - X_t) | A_t = i]] = \sum_{i=1}^K \sum_{t=1}^n \Delta_i E[I[A_t = i]] = \sum_{i=1}^K \Delta_i \sum_{t=1}^n E[I[A_t = i]]$$

$$\Rightarrow = \sum_{i=1}^K \Delta_i E[T_i(n)] \quad \checkmark$$

$$E[T_i(n)] \leq \frac{16 \log n}{\Delta_i^2} + 3$$

از طرفی طبق تعریف Δ_i داریم:

$$R_n = \sum_{i=1}^K \Delta_i E[T_i(n)] \leq \sum_{i=1}^K 3 \Delta_i + \sum_{i: \Delta_i > 0} \frac{16 \log n}{\Delta_i} \quad \checkmark$$

درستی:

$$R_n = \sum_{i=1}^K \Delta_i E[T_i(n)]$$

(د) طبقه‌بندی داده‌ها:

این جمع را به دو بخش برای داده‌های درجه دوم و زیرمجموعه‌ای تقسیم می‌کنیم.

$$R_n = \sum_{i: \Delta_i < \Delta} \Delta_i E[T_i(n)] + \sum_{i: \Delta_i \geq \Delta} \Delta_i E[T_i(n)] \leq n\Delta + \sum_{i: \Delta_i \geq \Delta} \left(3\Delta_i + \frac{16 \log n}{\Delta_i} \right)$$

$$\leq n\Delta + \sum_{i=1}^K 3\Delta_i + \frac{16K \log n}{\Delta} = \sqrt{16Kn \log n} + \sqrt{16Kn \log n} + \sum_{i=1}^K 3\Delta_i = 2\sqrt{Kn \log n} + 3 \sum_{i=1}^K \Delta_i$$