

سؤال ①:

الف) عبارت اول: $E_{s \sim D, a \sim \mu(a|s)} \{Q(s, a)\}$ - $E_{s \sim D, a \sim \hat{\mu}_B(a|s)} \{Q(s, a)\}$

است. مسئله اصلی در off-line RL این است که Distributional Shift وجود دارد و در واقع Q-value ها overestimate می شوند، جمله اول که داریم

$E_{s \sim D, a \sim \mu(a|s)} \{Q(s, a)\}$ سؤال این است که عبارت می که overestimate

می شوند را کاهش دهیم اما (s, a) هایی وجود دارند که داخل دیتای اصلی وجود دارند و در واقع overestimate نشده اند. در اینجا عبارت دوم که داریم $E_{s \sim D, a \sim \hat{\mu}_B(a|s)} \{Q(s, a)\}$ - $E_{s \sim D, a \sim \mu(a|s)} \{Q(s, a)\}$ سؤال این است که

عبارتی که قبلاً دیدیم و می بینیم overestimate نشده اند را کنترل کنیم به این صورت است می که قبلاً دیدیم را overestimate نمی کنیم.

(ب) ابتدا بهینه سازی داخل را به صورت \max_{μ} را حل می کنیم

Objective: $\max_{\mu} E_{s \sim D, a \sim \mu(a|s)} \{Q(s, a)\} + E_{s \sim D} \{H(\mu(a|s))\}$

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

حال می توان این مسئله را برای هر s به صورت جداگانه حل کرد و در نهایت داریم که مسئله بهینه سازی ما به شکل زیر است:

$\max_{\mu} \sum_a \mu(a|s) Q(s, a) - \sum_a \mu(a|s) \log(\mu(a|s))$

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

حال با استفاده از لاگرانژ داریم که:

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

$\frac{\partial L}{\partial \mu(a|s)} = Q(s, a) - 1 - \log \mu(a|s) + \lambda = 0 \Rightarrow \mu(a|s) = e^{-1+\lambda} e^{Q(s, a)}$

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

$\mu(a|s) = \frac{e^{Q(s, a)}}{\sum_{a'} e^{Q(s, a')}} \Rightarrow \mu(a|s) = \frac{e^{Q(s, a)}}{\sum_{a'} e^{Q(s, a')}} \Rightarrow$

$$CQL = \min_Q \left\{ \alpha \mathbb{E}_{S \sim D} \left\{ \sum_a \mu(a|s) Q(s,a) \right\} + \mathbb{E}_{S \sim D} \left\{ H(\mu(\cdot|s)) \right\} \right\}$$

$$CQL = \min_Q \left\{ \alpha \mathbb{E}_{S \sim D} \left\{ \mathbb{E}_{a \sim \mu(a|s)} \{Q(s,a)\} + \mathbb{E}_{S \sim D} \left\{ H(\mu(\cdot|s)) \right\} \right\} \right. \\ \left. - \alpha \mathbb{E}_{S \sim D} \left\{ \mathbb{E}_{a \sim \hat{\pi}_B(a|s)} \{Q(s,a)\} \right\} + \frac{1}{2} \mathbb{E}_{S,a,s' \sim D} \left\{ (Q - \hat{B}^{\pi_k} \hat{Q}^k)^2 \right\} \right\} =$$

$$\Rightarrow \textcircled{I} = \min_Q \left\{ \alpha \mathbb{E}_{S \sim D} \left\{ \alpha \mathbb{E}_{a \sim \mu(a|s)} \{Q(s,a)\} + H(\mu(\cdot|s)) \right\} \right\} =$$

$$= \min_Q \left\{ \mathbb{E}_{S \sim D} \left\{ \alpha \sum_a \mu(a|s) Q(s,a) - \sum_a \mu(a|s) \log(\mu(a|s)) \right\} \right\} =$$

$$= \min_Q \left\{ \mathbb{E}_{S \sim D} \left\{ \alpha \sum_a \mu(a|s) \left\{ Q(s,a) - Q(s,a) + \sum_{a'} \log \left\{ \sum_{a'} \exp(Q(s,a')) \right\} \right\} \right\} \right\}$$

$$= \min_Q \left\{ \mathbb{E}_{S \sim D} \left\{ \alpha \sum_a \log \alpha \log \left(\sum_a \exp(Q(s,a)) \right) \right\} \right\} \Rightarrow$$

$$\Rightarrow CQL = \min_Q \left\{ \mathbb{E}_{S \sim D} \left\{ \log \sum_a \exp(Q(s,a)) - \mathbb{E}_{a \sim \hat{\pi}_B(a|s)} \{Q(s,a)\} \right\} \right. \\ \left. + \frac{1}{2} \mathbb{E}_{S,a,s' \sim D} \left\{ (Q - \hat{B}^{\pi_k} \hat{Q}^k)^2 \right\} \right\}$$

$$\mu^* = \arg \max_{\mu} \left\{ \mathbb{E}_{S \sim D, a \sim \mu(a|s)} \{Q(s,a)\} + R(\mu) \right\} \quad (2)$$

$$CQL = \min_Q \left\{ \alpha \mathbb{E}_{S \sim D, a \sim \mu^*(a|s)} \{Q(s,a)\} + \frac{1}{2} \mathbb{E}_{S,a,s' \sim D} \left\{ (Q(s,a) - \hat{B}^{\pi_k} \hat{Q}^k(s,a))^2 \right\} \right. \\ \left. + R(\mu^*) \right\}$$

→ \min_Q نه قیمت قبل من توان بازی هر S این را به بازی کرد

$$CQL = \min_Q \left\{ \mathbb{E}_{a \sim \mu^*(a|s)} \{Q(s,a)\} + \frac{1}{2} \mathbb{E}_{a \sim \hat{\pi}_B(a|s)} \left\{ \left| Q(s,a) - \hat{B}^{\pi^k} Q(s,a) \right|^2 \right\} \right\}$$

A

حال می توان با استفاده از این عبارت نسبت به $Q(s,a)$ مشتق گرفت:

$$\frac{\partial CQL}{\partial Q} = 0 \Rightarrow \frac{\partial}{\partial Q(s,a)} \left\{ \alpha \sum_a \mu^*(a|s) Q(s,a) + \frac{1}{2} \sum_a \hat{\pi}_B(a|s) \left(Q(s,a) - \hat{B}^{\pi^k} Q(s,a) \right)^2 \right\} = 0$$

$$\Rightarrow \alpha \mu^*(a|s) + \hat{\pi}_B(a|s) \left(Q(s,a) - \hat{B}^{\pi^k} Q(s,a) \right) = 0 \Rightarrow$$

$$\Rightarrow Q(s,a) = \hat{B}^{\pi^k} Q(s,a) - \alpha \frac{\mu^*(a|s)}{\hat{\pi}_B(a|s)}$$

fixed-point : $\hat{Q}^{\pi}(s,a) = \hat{B}^{\pi^{\pi}} \hat{Q}^{\pi}(s,a) - \alpha \frac{\mu^*(a|s)}{\hat{\pi}_B(a|s)}$

$$\forall Q, s, a \in D : -C_\delta(s,a) \leq \hat{B}^{\pi^k} Q(s,a) - B^{\pi} Q(s,a) \leq C_\delta(s,a)$$

$$\forall Q : B^{\pi} Q = R + \gamma P^{\pi} Q$$

همین داریم:

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

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

$$\Rightarrow \hat{Q}^{\pi}(s,a) \leq (I - \gamma P^{\pi})^{-1} R - \alpha \left\{ (I - \gamma P^{\pi})^{-1} \left(\frac{\mu}{\hat{\pi}_B} \right) \right\}(s,a) + \left\{ (I - \gamma P^{\pi})^{-1} C_\delta \right\}(s,a)$$

$$\Rightarrow \hat{Q}^{\pi}(s,a) \leq Q^{\pi}(s,a) - \alpha \left\{ (I - \gamma P^{\pi})^{-1} \left(\frac{\mu}{\hat{\pi}_B} \right) \right\}(s,a) + \left\{ (I - \gamma P^{\pi})^{-1} C_\delta \right\}(s,a)$$

الف) در Inverse Reinforcement Learning (LRL) هدف به دست آوردن تابع reward است با استفاده از دیتاهای expert. در این صورت با توجه به اینکه یک تابع به دست می آوریم در مقابل دیتاهای که expert آنها را ندیده است، به صورت جنرال به عمل می کند و همچنین در مقابل دیتاهای که نویز داشته باشند بهتر عمل می کند.

در مقابل Behavioral Cloning (BC) به صورت مستقیم رفتار expert را تقلید می کند. با توجه به اینکه BC صرفاً تقلید می کند. در مقابل دیتاهای نویزی ممکن است action بهتری بعد و در زمان های که اکت دیده نشده، نمی تواند خوب عمل کند زیرا دادهای برای تقلید ندارد.

به علاوه با توجه به اینکه LRL یک تابع reward به دست می آورد، تقریباً می تواند رفتار expert را عکس تحلیل کند. که در بعضی از شرایط این اتفاق می تواند منجر به این شود که LRL از یک رفتار sup-optimal یک رفتار بهتر پیدا کند.

ب) در روش MaxEnt LRL یک روش از LRL است که در آن هدف به دست آوردن یک ریوارد فانکشن است که بتواند رفتار expert را توجیه کند طبقاً می توان چندین تابع مختلف را گذاشت و روش MaxEnt LRL فانکشن را انتخاب می کند که در آن بین همه distribution های مختلف بهترین آنترپی را داشته باشد. بدین صورت policy که یاد گرفته می شود، حداقل بایس را دارد و بهترین عدم قطعیت بدین صورت یک روش بسیار خوب که policy پیدا می شود که generalized ترین است.

$$P(\pi|\pi) = \frac{1}{Z(\pi)} \exp \left\{ \sum_{t=0}^T R(s_t, a_t) \right\}$$

$$\text{Objective} : \max_{\theta} \left\{ \sum_{\pi \in \mathcal{D}} P(\pi|\theta) \log P(\pi|\theta) - \lambda \left(\mathbb{E}_{\pi \in \mathcal{D}} \left\{ \Phi(s, a) \right\} - \mathbb{E}_{\pi \in \mathcal{E}} \left\{ Q(s, a) \right\} \right) \right\}$$

$$H(\pi) = - \sum_{a,s} \pi(a|s) \log \pi(a|s)$$

(7)

حاکم کردن آنشوری باعث می شود که policy غیر قطعی (uncertain) باقی بماند که باعث می شود محیط را بهتر explore کنیم و جلوگیری کنیم از برای (s,a) هایی که expert دیدات جلوگیری می کند که overfit نشوند.

$$D: \text{discriminator} \quad J(\pi_E, L) = \mathbb{E}_{\pi_E} \{ \log(1 - D(s, a)) \}$$

$$J(\pi, L) = \mathbb{E}_{\pi} \{ \log(D(s, a)) \} - \lambda H(\pi)$$

این تابع جلوگیری از تلاش می کند که فاصله بین پالیسی expert با پالیسی است و یا بیشتر کند. بازی min-max به این صورت کار می کند که ابتدا با maximize کردن نسبت به D تلاش می کند فاصله بین policy اکسپرت با policy یادگیری شده را زیاد کند و سپس با minimize کردن نسبت به π تلاش می کند که بهترین جواب را در بهترین محیط پیدا کند.

(8) MaxEnt ZRL: مسأله از نظر پیچیدگی objective و حساب نسبت به ویژگی های انتخاب شده دارد. در بعضی اوقات نیاز به تعداد زیادی sample دارد. در روش GAIL با توجه به اینکه یادگیرنده adversarial است باعث می شود که پیچیدگی حل معادله اصلی کاهش پیدا کند. و نسبت به ویژگی ها خیلی مختلف است. (که می تواند یک neural net باشد)

همچنین عبارت آنشوری داخل GAIL باعث می شود که overfitting, exploration بهتر شوند.

(الف) روش‌های که تنها از مدل محیط استفاده می‌کنند:

محکم‌کردن مدل به نسبت پیشی به مدل است. اگر محیط بیش از اندازه ساده سازی شده باشد یا مشکل داشته باشد، خروفتی مدل یک $sub\text{-optimal policy}$ می‌شود یا حتی $policy$ آبی را همچنین یک مسئله به نام $reward\text{-degradation}$ وجود دارد که ممکن است $policy$ ها را نتوان به دنیای واقعی تعمیم داد. همچنین این مقودها نمی‌توانند انتقادات غیره ای را پیش بینی کنند و در شرایطی می‌توانند $overfit$ شوند.

روش‌هایی که تنها با محیط تعامل می‌کنند:

تعامل با محیط اغلب هزینه زیادی دارد و مدت زیادی صرف تعامل می‌شود. همچنین در بعضی $action$ ها ممکن است که تعامل با محیط آبی باشد. خوبی روش‌های ترکیبی:

این روش‌ها می‌توانند از خوبی‌های هر دو روش بهره ببرند. کمک این صورت می‌توان $generalization$ بیشتری داشته باشیم.

هزینه کمتری برای بدست آوردن $policy$ ها به دامنه (از لحاظ زمان و منابع) و همچنین در جاهایی که آبی از $interaction$ با محیط بهره‌مند است.

(ب)

روش‌های MBPO از دو روش مدل محیط و تعامل با محیط بهره‌مند می‌برند.

این روش بدین صورت کار می‌کند که ابتدا با استفاده از تعامل با محیط یک سری داده ایجاد می‌کنند و سپس با یادگیری مدل تعدادی داده جدید ایجاد می‌کنند و $dataset$ جمع شده از این دو مدل ایجاد می‌کنند. بدین صورت مدل می‌تواند نمونه‌های مختلف بیشتری را بیند.

انتخاب طول $roll-out$ یک انتخاب حیاتی است زیرا

طول کم $roll-out$ به معنی این است که اگر مدل را $minimize$ می‌کنند اما $diversity$ کمتری دارد.

طول زیاد $roll-out$ مشکل $diversity$ را رفع می‌کند اما به خاطر تجمع خطاهای باعث می‌شود که خطای بیشتری داشته باشد.

در نتیجه باید یک انتخاب درست بین برای طول $roll-out$ مشخص کرد که $trade-off$ خوبی بین خطا و $diversity$ داشته باشد.

به صورت پویا

2) نگاهمانور که در سمت های قبل توضیح دادیم. اوس MBPO گام اول
متنی به interaction با هم می آید. که به صورت پیوسته در صورت نیاز با مدل
interact کند و مدل خود را update کند. در اینجا حالتی که داده ها offline
باشند، MBPO نمی تواند با هم interact کند، خطاهای خود را تصحیح کند.
مشکلات ۱۰۵۱، ۱۰۵۲، ۱۰۵۳، ۱۰۵۴، ۱۰۵۵، ۱۰۵۶، ۱۰۵۷، ۱۰۵۸، ۱۰۵۹، ۱۰۶۰، ۱۰۶۱، ۱۰۶۲، ۱۰۶۳، ۱۰۶۴، ۱۰۶۵، ۱۰۶۶، ۱۰۶۷، ۱۰۶۸، ۱۰۶۹، ۱۰۷۰، ۱۰۷۱، ۱۰۷۲، ۱۰۷۳، ۱۰۷۴، ۱۰۷۵، ۱۰۷۶، ۱۰۷۷، ۱۰۷۸، ۱۰۷۹، ۱۰۸۰، ۱۰۸۱، ۱۰۸۲، ۱۰۸۳، ۱۰۸۴، ۱۰۸۵، ۱۰۸۶، ۱۰۸۷، ۱۰۸۸، ۱۰۸۹، ۱۰۹۰، ۱۰۹۱، ۱۰۹۲، ۱۰۹۳، ۱۰۹۴، ۱۰۹۵، ۱۰۹۶، ۱۰۹۷، ۱۰۹۸، ۱۰۹۹، ۱۱۰۰، ۱۱۰۱، ۱۱۰۲، ۱۱۰۳، ۱۱۰۴، ۱۱۰۵، ۱۱۰۶، ۱۱۰۷، ۱۱۰۸، ۱۱۰۹، ۱۱۱۰، ۱۱۱۱، ۱۱۱۲، ۱۱۱۳، ۱۱۱۴، ۱۱۱۵، ۱۱۱۶، ۱۱۱۷، ۱۱۱۸، ۱۱۱۹، ۱۱۲۰، ۱۱۲۱، ۱۱۲۲، ۱۱۲۳، ۱۱۲۴، ۱۱۲۵، ۱۱۲۶، ۱۱۲۷، ۱۱۲۸، ۱۱۲۹، ۱۱۳۰، ۱۱۳۱، ۱۱۳۲، ۱۱۳۳، ۱۱۳۴، ۱۱۳۵، ۱۱۳۶، ۱۱۳۷، ۱۱۳۸، ۱۱۳۹، ۱۱۴۰، ۱۱۴۱، ۱۱۴۲، ۱۱۴۳، ۱۱۴۴، ۱۱۴۵، ۱۱۴۶، ۱۱۴۷، ۱۱۴۸، ۱۱۴۹، ۱۱۵۰، ۱۱۵۱، ۱۱۵۲، ۱۱۵۳، ۱۱۵۴، ۱۱۵۵، ۱۱۵۶، ۱۱۵۷، ۱۱۵۸، ۱۱۵۹، ۱۱۶۰، ۱۱۶۱، ۱۱۶۲، ۱۱۶۳، ۱۱۶۴، ۱۱۶۵، ۱۱۶۶، ۱۱۶۷، ۱۱۶۸، ۱۱۶۹، ۱۱۷۰، ۱۱۷۱، ۱۱۷۲، ۱۱۷۳، ۱۱۷۴، ۱۱۷۵، ۱۱۷۶، ۱۱۷۷، ۱۱۷۸، ۱۱۷۹، ۱۱۸۰، ۱۱۸۱، ۱۱۸۲، ۱۱۸۳، ۱۱۸۴، ۱۱۸۵، ۱۱۸۶، ۱۱۸۷، ۱۱۸۸، ۱۱۸۹، ۱۱۹۰، ۱۱۹۱، ۱۱۹۲، ۱۱۹۳، ۱۱۹۴، ۱۱۹۵، ۱۱۹۶، ۱۱۹۷، ۱۱۹۸، ۱۱۹۹، ۱۲۰۰، ۱۲۰۱، ۱۲۰۲، ۱۲۰۳، ۱۲۰۴، ۱۲۰۵، ۱۲۰۶، ۱۲۰۷، ۱۲۰۸، ۱۲۰۹، ۱۲۱۰، ۱۲۱۱، ۱۲۱۲، ۱۲۱۳، ۱۲۱۴، ۱۲۱۵، ۱۲۱۶، ۱۲۱۷، ۱۲۱۸، ۱۲۱۹، ۱۲۲۰، ۱۲۲۱، ۱۲۲۲، ۱۲۲۳، ۱۲۲۴، ۱۲۲۵، ۱۲۲۶، ۱۲۲۷، ۱۲۲۸، ۱۲۲۹، ۱۲۳۰، ۱۲۳۱، ۱۲۳۲، ۱۲۳۳، ۱۲۳۴، ۱۲۳۵، ۱۲۳۶، ۱۲۳۷، ۱۲۳۸، ۱۲۳۹، ۱۲۴۰، ۱۲۴۱، ۱۲۴۲، ۱۲۴۳، ۱۲۴۴، ۱۲۴۵، ۱۲۴۶، ۱۲۴۷، ۱۲۴۸، ۱۲۴۹، ۱۲۵۰، ۱۲۵۱، ۱۲۵۲، ۱۲۵۳، ۱۲۵۴، ۱۲۵۵، ۱۲۵۶، ۱۲۵۷، ۱۲۵۸، ۱۲۵۹، ۱۲۶۰، ۱۲۶۱، ۱۲۶۲، ۱۲۶۳، ۱۲۶۴، ۱۲۶۵، ۱۲۶۶، ۱۲۶۷، ۱۲۶۸، ۱۲۶۹، ۱۲۷۰، ۱۲۷۱، ۱۲۷۲، ۱۲۷۳، ۱۲۷۴، ۱۲۷۵، ۱۲۷۶، ۱۲۷۷، ۱۲۷۸، ۱۲۷۹، ۱۲۸۰، ۱۲۸۱، ۱۲۸۲، ۱۲۸۳، ۱۲۸۴، ۱۲۸۵، ۱۲۸۶، ۱۲۸۷، ۱۲۸۸، ۱۲۸۹، ۱۲۹۰، ۱۲۹۱، ۱۲۹۲، ۱۲۹۳، ۱۲۹۴، ۱۲۹۵، ۱۲۹۶، ۱۲۹۷، ۱۲۹۸، ۱۲۹۹، ۱۳۰۰، ۱۳۰۱، ۱۳۰۲، ۱۳۰۳، ۱۳۰۴، ۱۳۰۵، ۱۳۰۶، ۱۳۰۷، ۱۳۰۸، ۱۳۰۹، ۱۳۱۰، ۱۳۱۱، ۱۳۱۲، ۱۳۱۳، ۱۳۱۴، ۱۳۱۵، ۱۳۱۶، ۱۳۱۷، ۱۳۱۸، ۱۳۱۹، ۱۳۲۰، ۱۳۲۱، ۱۳۲۲، ۱۳۲۳، ۱۳۲۴، ۱۳۲۵، ۱۳۲۶، ۱۳۲۷، ۱۳۲۸، ۱۳۲۹، ۱۳۳۰، ۱۳۳۱، ۱۳۳۲، ۱۳۳۳، ۱۳۳۴، ۱۳۳۵، ۱۳۳۶، ۱۳۳۷، ۱۳۳۸، ۱۳۳۹، ۱۳۴۰، ۱۳۴۱، ۱۳۴۲، ۱۳۴۳، ۱۳۴۴، ۱۳۴۵، ۱۳۴۶، ۱۳۴۷، ۱۳۴۸، ۱۳۴۹، ۱۳۵۰، ۱۳۵۱، ۱۳۵۲، ۱۳۵۳، ۱۳۵۴، ۱۳۵۵، ۱۳۵۶، ۱۳۵۷، ۱۳۵۸، ۱۳۵۹، ۱۳۶۰، ۱۳۶۱، ۱۳۶۲، ۱۳۶۳، ۱۳۶۴، ۱۳۶۵، ۱۳۶۶، ۱۳۶۷، ۱۳۶۸، ۱۳۶۹، ۱۳۷۰، ۱۳۷۱، ۱۳۷۲، ۱۳۷۳، ۱۳۷۴، ۱۳۷۵، ۱۳۷۶، ۱۳۷۷، ۱۳۷۸، ۱۳۷۹، ۱۳۸۰، ۱۳۸۱، ۱۳۸۲، ۱۳۸۳، ۱۳۸۴، ۱۳۸۵، ۱۳۸۶، ۱۳۸۷، ۱۳۸۸، ۱۳۸۹، ۱۳۹۰، ۱۳۹۱، ۱۳۹۲، ۱۳۹۳، ۱۳۹۴، ۱۳۹۵، ۱۳۹۶، ۱۳۹۷، ۱۳۹۸، ۱۳۹۹، ۱۴۰۰، ۱۴۰۱، ۱۴۰۲، ۱۴۰۳، ۱۴۰۴، ۱۴۰۵، ۱۴۰۶، ۱۴۰۷، ۱۴۰۸، ۱۴۰۹، ۱۴۱۰، ۱۴۱۱، ۱۴۱۲، ۱۴۱۳، ۱۴۱۴، ۱۴۱۵، ۱۴۱۶، ۱۴۱۷، ۱۴۱۸، ۱۴۱۹، ۱۴۲۰، ۱۴۲۱، ۱۴۲۲، ۱۴۲۳، ۱۴۲۴، ۱۴۲۵، ۱۴۲۶، ۱۴۲۷، ۱۴۲۸، ۱۴۲۹، ۱۴۳۰، ۱۴۳۱، ۱۴۳۲، ۱۴۳۳، ۱۴۳۴، ۱۴۳۵، ۱۴۳۶، ۱۴۳۷، ۱۴۳۸، ۱۴۳۹، ۱۴۴۰، ۱۴۴۱، ۱۴۴۲، ۱۴۴۳، ۱۴۴۴، ۱۴۴۵، ۱۴۴۶، ۱۴۴۷، ۱۴۴۸، ۱۴۴۹، ۱۴۵۰،

در این روش چند مدل مختلف به روی $offline\ data$ آموزش می‌دهند. به این صورت هر مدل جداگانه برای $next\ state$ تصمیم‌گیری می‌کند. حال من آن‌ها را با هم مقایسه می‌کنم و بر اساس $threshold$ برای قبول کردن یا رد کردن آن‌ها تصمیم می‌گیرم.

$$Q_{\text{new}}(s, a) = \min \{ Q_{\text{old}}(s, a), r + \gamma \max_{a' \in \pi} Q_{\text{target}}(s, a') \}$$

COMBO از این‌ها در روش CQL و Dyna استفاده می‌کنند
CQL : از روش CQL این‌ها را استفاده می‌کنند که با قرار دادن یک پناهی
جلوی Overestimation برای Q-value ها را می‌گیرد و همگانی‌تر می‌کند

Dyna: از این ابزار استفاده می کنند نه این ابزار
با استفاده از
طرح های مجرّبی ایجاد کنند و از آنها استفاده کنند.

COMBO با استفاده از این ایده ها تلاش می کند که با داده های *explorative* انجام دهد و *Q-value* ها را به صورت مقایسه ای آپدیت کند. به این صورت به هدف دو روش Dyna و CQL رسیده است.

(9) هر دو روش COMBO برای Model Rel offline RL هستند.

COMBO

Strengths

(1) از مدل روش مبتنی بر مدل و محافظه کارانه استفاده می کنند که یک بالانس بین exploration و safety پیدا کند.

(2) sample-efficient: استفاده از یک مدل برای ایجاد داده ها باعث می شود که

efficiency بهتر شود.

(3) Robustness: محافظه کارانه بودن این روش باعث می شود که به عملکردی قابل اطمینانتری برسم.

Weaknesses

(1) Model dependency: وابسته به کیفیت مدل است.

(2) Conservatism: بعضی اوقات خیلی محافظه کارانه عمل می کند.

(3) Complexity: ایجاد تعامل بین exploration و Conservatism هزینه زیادی داشته باشد.

Model Rel

Strengths

(1) Explicit Uncertainty: عدم قطعیت را به صورت مستقیم در نظر می گیرد.

(2) Efficiency: استفاده از آنسامبل ها efficient می کند.

(3) Safety: با تغییر MDP مطمئن می شود که policy امن نباشد امن است.

Weaknesses

(1) Model dependency

(2) Computational Intensity: استفاده از آنسامبل زمان زیادی می برد.

(3) Conservatism

در کل هر دو روش برای offline RL مناسب هستند. COMBO یک تعامل بین محافظه کاری و exploration به قرار می دهد که Model Rel کاملاً بر شرایط غیرقطعی متمرکز می کند. هم دوی این روش ها از Model استفاده می کنند که این انحراف باعث وابستگی بالایی آنها به دقت مدل می باشد نه مشکل زا است.

الف) طبق نامساوی Hoeffding داریم که: $S_n = \sum_i X_i, a_i \leq X_i \leq b_i$

$$\hookrightarrow P(|E(X) - S_n| \geq \epsilon) \leq \exp \left\{ - \frac{2\epsilon^2}{\sum_i (b_i - a_i)^2} \right\} \quad X_i \rightarrow \text{i.i.d.}$$

حال اگر داشته باشیم: $a_i = -1, b_i = 1$ و همچنین $\bar{X} = \frac{\sum X_i}{n}$

$$\hookrightarrow P(|E(X) - \bar{X}| \geq \epsilon) \leq \exp \left\{ - \frac{n\epsilon^2}{2} \right\}$$

فرض کنیم که $\epsilon = \sqrt{\frac{2 \log(1/\delta)}{n}}$ داریم که

$$P \left\{ |\bar{X} - E(X)| \geq \sqrt{\frac{2 \log(1/\delta)}{n}} \right\} \leq \delta \Rightarrow P \left(\mu \geq \hat{\mu} + \sqrt{\frac{2 \log(1/\delta)}{n}} \right) \leq \delta$$

در نتیجه داریم که با احتمال $1 - \delta$ $\mu \leq \hat{\mu} + \sqrt{\frac{2 \log(1/\delta)}{n}}$ و.
در نتیجه با احتمال بالا داریم که $\hat{\mu} + \sqrt{\frac{2 \log(1/\delta)}{n}}$ یک حد بالا برای μ است.

ب) بدون کاست از جمله فرض می کنیم که بهترین arm همان arm اول باشد
و $\mu^* = \mu_1$ باشد.

قبل از ادامه دادن اثبات notation را بصورت زیر تعریف می کنیم:

یادش متوجه شد. در تکاپی اجلی \rightarrow i.i.d variables $\rightarrow (X_{ti})_{t \in [n], i \in [k]} \rightarrow$ arm i

$$\hat{\mu}_{is} = \frac{1}{s} \sum_{u=1}^s X_{ui} \quad R_n = \sum_{i=1}^k \Delta_i (E \{ T_i(n) \})$$

این اثبات بر این صورت است که یک bound برای $E(T_i(n))$ پیدا می کنیم. V_i suboptimal. G_i را به این صورت تعریف می کنیم:

$$G_i = \left\{ \mu_1 \leq \min_{t \in [n]} VCB_1(t, i) \right\} \cap \left\{ \hat{\mu}_{i, n_i} + \sqrt{\frac{2}{n_i} \log\left(\frac{1}{\delta}\right)} \leq \mu_1 \right\}$$

که $u_i \in [n]$ یک ثابت است که در ادامه نشان می دهیم.

حال ۲ چیز را نشان می دهیم:

۱) اگر G_i اتفاق بیفتد، arm i مالیم u_i بار انتخاب شده.

۲) G_i^c (Complement) به احتمال کم اتفاق می افتد.

می‌باشد که داریم $T_i(n) \ll n$

(IV)

$$E(T_i(n)) = E\{ \sum_{G_i \in \mathcal{G}_i} T_i(n) \} + E\{ \sum_{G_i \notin \mathcal{G}_i} T_i(n) \} \leq \sum_{G_i \in \mathcal{G}_i} P(G_i) + \sum_{G_i \notin \mathcal{G}_i} P(G_i)$$

برای اثبات قضیه ① از بهر حال خلف استفاده می‌کنیم.

فرض کنیم که $T_i(n) > u_i$ باشد. پس α -arm بیش از u_i بار در n بار بازی شده پس یک دور داشته $t \in [n]$ وجود دارد که $T_i(t-1) = u_i$ و $A_t = i$ باشد. از تعریف G_i داریم که:

$$UCB_i(t-\delta) = \hat{\mu}_i(t-1) + \sqrt{\frac{2 \log(1/\delta)}{T_i(t-1)}} \rightarrow UCB_i(t-\delta)$$

$$= \hat{\mu}_{i,u_i} + \sqrt{\frac{2 \log(1/\delta)}{u_i}} \quad \{ \mu_i, UCB_i(t-\delta) \}$$

در نتیجه داریم که $A_t = \arg \max_j UCB_j(t-\delta) \neq i$ که متناقض است. اثبات قضیه ②:

برای اثبات قضیه ② طبق تعریف داریم

$$G_i^c = \{ \mu_i \mid \exists t \in [n] \text{ such that } UCB_i(t, \delta) > \hat{\mu}_{i,u_i} + \sqrt{\frac{2 \log(1/\delta)}{u_i}} \} \quad \text{III}$$

لین تعریف $UCB_i(t, \delta)$ داریم که

$$\{ \mu_i \mid \exists t \in [n] \text{ such that } UCB_i(t, \delta) > \hat{\mu}_{i,s} + \sqrt{\frac{2 \log(1/\delta)}{s}} \} = \bigcup_{s \in [n]} \{ \mu_i \mid \exists t \in [n] \text{ such that } UCB_i(t, \delta) > \hat{\mu}_{i,s} + \sqrt{\frac{2 \log(1/\delta)}{s}} \}$$

طبق قانون Union bound داریم که

$$P(\mu_i \mid \exists t \in [n] \text{ such that } UCB_i(t, \delta) > \hat{\mu}_{i,s} + \sqrt{\frac{2 \log(1/\delta)}{s}}) \leq P(\bigcup_{s \in [n]} \{ \mu_i \mid \exists t \in [n] \text{ such that } UCB_i(t, \delta) > \hat{\mu}_{i,s} + \sqrt{\frac{2 \log(1/\delta)}{s}} \})$$

$$\leq \sum_{s=1}^n P(\mu_i \mid \exists t \in [n] \text{ such that } UCB_i(t, \delta) > \hat{\mu}_{i,s} + \sqrt{\frac{2 \log(1/\delta)}{s}}) \leq n \delta$$

Moehfolding (II)

حال که (2) قضیه را ثابت کردیم، فرض کنیم که u_i به صورتی قرار می‌گیرد که انتفا باشد

$$\Delta_i - \sqrt{\frac{2 \log(1/\delta)}{u_i}} \gg c \Delta_i \rightarrow c_i \in (0, 1)$$

$$\mu_i = \mu_i + \Delta_i \Rightarrow \mathbb{P}\left(\hat{\mu}_{i u_i} + \sqrt{\frac{2 \log(1/\delta)}{u_i}} \gg \mu_i\right) =$$

$$= \mathbb{P}\left(\hat{\mu}_{i u_i} - \mu_i \gg \Delta_i - \sqrt{\frac{2 \log(1/\delta)}{u_i}}\right) \leq \mathbb{P}\left(\hat{\mu}_{i u_i} - \mu_i \gg c \Delta_i\right) \leq \exp\left(-\frac{u_i c^2 \Delta_i^2}{2}\right)$$

$$\textcircled{I}, \textcircled{II} \Rightarrow \mathbb{P}(G_i^c) \leq n \delta + \exp\left\{-\frac{u_i c^2 \Delta_i^2}{2}\right\}$$

$$\textcircled{IV} \Rightarrow \mathbb{E}(\bar{T}_i(n)) \leq u_i + n \left(n \delta + \exp\left\{-\frac{u_i c^2 \Delta_i^2}{2}\right\} \right)$$

حال کوچکترین n_i را انتخاب می‌کنیم به صورتی که $\Delta_i - \sqrt{\frac{2 \log(1/\delta)}{u_i}} \gg c \Delta_i$ باشد

$$\hookrightarrow u_i = \left\lceil \frac{2 \log(1/\delta)}{(1-c)^2 \Delta_i^2} \right\rceil$$

حال با توجه به آنکه $\bar{T}_i(n) \leq n$ و با فرض $\delta = \frac{1}{n^2}$

$$\mathbb{E}(\bar{T}_i(n)) \leq u_i + 1 + n \left(1 + n^{-\frac{2c^2}{(1-c)^2}} \right)$$

$$\mathbb{E}(\bar{T}_i(n)) \leq u_i + 1 + n^{1 - \frac{2c^2}{(1-c)^2}} = \left\lceil \frac{2 \log(1/\delta)}{(1-c)^2 \Delta_i^2} \right\rceil + 1 + n^{1 - \frac{2c^2}{(1-c)^2}}$$

حال تنها کافی است $c \in (0, 1)$ را انتخاب کنیم

c اگر خیلی به (1) نزدیک باشد مقدار n به سمت بی‌نهایت می‌رود و همچنین

توان n باید کمتر از صفر باشد. به صورت دلخواه $c = \frac{1}{2}$ قرار می‌دهیم

$$\Rightarrow \mathbb{E}(\bar{T}_i(n)) \leq 3 + \frac{16 \log(n)}{\Delta_i^2}$$

(2) برای اثبات این قضیه، ابتدا اثبات می‌کنیم که

$$R_n = \sum_{i \in A} \Delta_i E\{T_i(n)\}$$

Proof

$$R_n = n \mu^* E\{S_n\} = \sum_{i \in A} \sum_{t=1}^k \sum_{s=1}^n E\{(\mu^* - X_t) I\{A_t = i\}\}$$

$$S_n = \sum_t X_t = \sum_t \sum_{i=1}^k X_t I\{A_t = i\}$$

حال می‌توان نوشت که:

$$E\{(\mu^* - X_t) I\{A_t = i\}\} = I\{A_t = i\} E\{\mu^* - X_t | A_t = i\}$$

$$= I\{A_t = i\} (\mu^* - \mu_{A_t}) = I\{A_t = i\} (\mu^* - \mu_i) =$$

$$= I\{A_t = i\} \Delta_i \Rightarrow$$

$$\Rightarrow E\{(\mu^* - X_t) I\{A_t = i\}\} = E\{I\{A_t = i\} \Delta_i\}$$

$$R_n = \sum_i \Delta_i \sum_t E\{I\{A_t = i\}\} = \sum_{i=1}^k \Delta_i E(T_i(n))$$

حال با بکارگیری این عبارت در چیزی که قبلاً (ب) استفاده کردیم داریم که

$$R_n \leq 3 \sum_{i=1}^k \Delta_i + \sum_{i: \Delta_i > 0} \frac{16 \log(n)}{\Delta_i}$$

$$R_n = \sum_{i=1}^k \Delta_i \mathbb{E}\{\bar{T}_i(n)\}$$

(د) داریم که

$$\Delta = \sqrt{\frac{16k \log(n)}{n}} \quad \text{فرض کنیم}$$

$$R_n = \sum_{i: \Delta_i < \Delta} \Delta_i \mathbb{E}(\bar{T}_i(n)) + \sum_{i: \Delta_i \geq \Delta} \Delta_i \mathbb{E}(\bar{T}_i(n))$$

$$\sum_{i: \Delta_i < \Delta} T_i(n) \approx n$$

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

$$\Rightarrow R_n \leq n\Delta + \sum_{i: \Delta_i \geq \Delta} \left\{ 3\Delta_i + \frac{16 \log(n)}{\Delta_i} \right\}$$

$$\leq n\Delta + \frac{16k \log(n)}{\Delta} + 3 \sum_i \Delta_i$$

حال با جایگزینی Δ داریم

$$R_n \leq 8\sqrt{nk \log(n)} + 3 \sum_i \Delta_i$$

انبار این سؤال با استفاده از کتاب

Bandit Algorithms

Tor Lattimore and Csaba Szepesvári