

توضیح محل الگوریتم SAC:

الگوریتم SAC یک الگوریتم model free (با استاندارد از نمونه گیری) و off-policy است. الگوریتمهای RL معمولی تلاش می کنند که راه حلی را به بهترین کردن یادداشتی انجامه و دنبال می کنند. اما در این الگوریتم

دنبال راه حلی (سیاستی هست) که در می کند بهترین یادداشت ممکن را بدست آورد اما هر زمان تا جایی می تواند رفتار غیر قابل

بیشتری و آشنایی خود را حفظ کند. **همدوفایه اساسی: بایداری بیشتر/آشنایی بیشتر**

هدف SAC بهینه کردن مجموع یادداشتی ما به همراه آنتروپی است:

$$J(\pi) = \sum_{t=0}^T E_{(s_t, a_t) \sim p_{\pi}} [r(s_t, a_t) + \alpha H(\pi(a|s_t))]$$

\nearrow ضریب اهمیت
 \uparrow آنتروپی سیاست
 \nearrow یادداشت

حال به اجزای SAC می پردازیم: البته SAC مدلی اینی بدون value:

۱- شبکه Actor (سیاست گذار): این شبکه state را گرفته و پارامترهای یک توزیع (تخمین از mean و std) را برای

انتخاب عمل غرضی می دهد.

۲- شبکه critic (منتقد): این شبکه می خواهد $Q(s, a)$ را البته بصورت soft-max تخمین بزند (x خنثی)،

توجه داریم استاندارد از دو شبکه و انتخاب min بین آنها باعث جلوگیری از overfit می شود.

۳- دو شبکه target-critic:

این شبکه های یکسانی با تناظر از critic هستند برای یادگیری خود actor و critic، این یکی باعث جلوگیری از overfit می شود.

تابع های loss و روند یادگیری:

critic: تلاش می کند خطای MSE بین سیگنالی خود و یک مقدار هدف را کم کند:

$$y = r(s, a) + \gamma \left(\min_{i=1,2} Q_i(s, a) - \alpha \log \pi(a|s) \right)$$

هدف

$$ssE = \sum (y - Q(s, a))^2$$

Actor : تابع هزینه Actor بر اساس هزینه وزن و ابرای KL سن سیاست و توزیع احتمالی (Q-function soft-max) است :

$$L_{actor}(\theta) = E_{a \sim \pi_{\theta}} [\alpha \log \pi_{\theta}(a|s) - \min_{Q_i} Q_i(s, a)]$$

به سمت اعلای رود که $Q(s, a)$ بهترین با آستروپی بستر شود.
- در نتیجه ای، مایباده ردی خود پاراستر θ نیزند پاراستر قابل یادگیری بود.

۵. آنزیمی $\log \pi(a|s)$ - در یاداشت باعث شود که یاداشت درونی برای نگهداری افعال شود. نتیجه این کار

می شود اینده التشاف به مقدار کافی انجام شود و احتمال پیرامندان در

ط. اقتراض بایدار و جلوگیری از overfit: یک سیاست با آسرونی بالا اعطای پذیر است، به جای یادگیری کامل و پیش برای بد موقعیت عامل یاد می گیرد مجموعه ای از احوال خوب وجود دارند.

نبار این تخمین شود که سیاحت کنندگان با اتعاق زیر منظره ای رخ دهد حاصل گزینهای دیگری
برای دانش دارد.

معادله من با soft-max به صورت زیر است :

$$Q_{soft}^*(s, a) = r(s, a) + \gamma E_{s' \sim p(a|s)} [Q_{soft}^*(s', a) - \alpha \log \pi(a|s')]]$$

↓
سیلیت لینک

با اعمال مکرر این ایراتو به این خاطر Q-function تابع نقطه ثابت است (fixed-point) به نقطه ثابت خود مگر می شود. حالا این Q^* چیست؟

ارزشی را به خروج از s و انجام a action و سپس دنبال کردن بهترین سیاست تعاقبی را تعادل بین یادگیری و آنتروپی را برتری کند به دست می آید.

$$\tau^\pi Q(s, a) \triangleq r(s, a) + \gamma E_{s' \sim p} [V_{soft}(s')]]$$

$$V_{soft}(s) = E_{a \sim \pi} [Q(s, a) - \alpha \log \pi(a|s)]$$

این تابع زمانی به نقطه ثابت خود می رسد که اعمال ایراتو دیگر باعث تغییر نمی شود یعنی $\tau^\pi Q = Q$

و طبق نتیجه نقطه ثابت اثبات می شود این مقدار، مقدار مگرایی است.

توجه داریم که یک ثابت نیست انتقابی است ($\gamma < 1$) اثبات می شود به این دلیل ما باید نقطه مگرایی شود.

هدف بدست آوردن مقدار θ به روش است:

$$J_{\text{deal policy}}(\alpha|\theta) = \frac{\exp(\frac{1}{\alpha} Q_{\theta}(s,a))}{2\theta}$$

این 2θ جمع وزن روی \exp \rightarrow

ما است

و برین آن باید در فضای پیکتید اشتغال را حساب کنیم. ضرب این اشتغال به دست آوردن قیمت است
بنابراین از ابزار حاصله و کمترین آن استفاده می‌کنیم

هدف ما این خواهد بود:

$$\min_{\phi} D_{KL} \left(\pi_{\phi}(\alpha|s) \parallel \frac{\exp(\frac{1}{\alpha} Q(s,a))}{2\theta} \right)$$

باید

یا راترمای بدست آمده از ϕ را باید طوری انتخاب کرد که توزیع خرجی آن کمترین فاصله را با مقدار ایده‌آل داشته باشد.

در ریاضیات داریم:

$$D_{KL}(p||Q) = E_{p,p} [\log p - \log Q(s)]$$

پس برین مساله ما خواهیم داشت:

$$D_{KL} = E_{a \sim \pi_{\phi}} \left[\log \pi_{\phi}(\alpha|s) - \log \left(\frac{e^{\frac{1}{\alpha} Q(s,a)}}{2\theta} \right) \right]$$

$$= E_{a \sim \pi_{\phi}} \left[\log \pi_{\phi}(\alpha|s) - \log(e^{\frac{1}{\alpha} Q(s,a)}) + \log 2\theta \right]$$

رابطه ϕ نه اندیش از E بیرون می‌آید و داریم:

$$= E_{a \sim \pi_{\phi}} \left[\log \pi_{\phi}(\alpha|s) - \frac{1}{\alpha} Q_{\phi}(s,a) \right]$$

$$= E_{a \sim \pi_{\phi}} [\alpha \log \pi_{\phi}(\alpha|s) - Q_{\phi}(s,a)]$$

و این میان تابع خطایی است که در Actor داریم یعنی باید می‌نویسیم:

$$E_{a \sim \pi_{\phi}} \left[\alpha \log \pi_{\phi}(\alpha|s) - Q_{\phi}(s,a) \right]$$

