

① (۲) استناد از sigmoid به علت سادگی که اگر افتلات winner و loser زیاد باشد، پدید gradient vanishing افتاد

بیفتد و در تقسیم تابع یاداش طوری آسان‌تر می‌باشد که افتلات winner و loser بسیار زیاد باشد.

$$\frac{\partial G(z)}{\partial z} = 1 - G(z) = 0$$

اگر جایی افتلات زیاد winner و loser را تنظیم، ممکن است باعث شود که به جز یک دسته بقیه دسته‌ها یاداش کمی اختصاص داده شود، با وجود اینکه ممکن است در مسائل سیای دسته‌ها، داده‌های بزرگ و بنابراین نباید یاداش بایستی دریافت کند. به همین علت تابع یاداش مرتباً به سمت صفر کاهش داده می‌شود تا بایاس نباشد.

(ب)  $\theta(z)$ : زوای فضاها داده‌های تبدیل شده مدل زبانی که پارامترهای  $\theta$  را دارد، تابع یاداش را می‌کشد. به عبارت دیگر Policy را آسان‌تر می‌کند که تابع یاداش را می‌کشد.

$$-\beta \log \left( \pi_{\theta}^{RL}(a|x) / \pi^{SFT}(a|x) \right) : \text{این نرم حاصل می‌شود از KL به این دلت است که اجازه می‌دهد در مسیر آسان‌تر تبدیل تابع داده زوای}$$

Policy بهینه‌سازی می‌کند و اساسی تابع یاداش در هر گام از توزیع فرجه Policy اولیه فاصله زیادی بگیرد. زوای Policy با داده‌های تبدیل شده در هر گام آسان‌تر می‌شود و اگر اجازه دهیم از زوای Policy اولیه فاصله زیادی بگیرد، ممکن است که داده‌های خوبی تبدیل شده (noisy) که این داده‌ها مرتباً می‌تواند تابع یاداش را می‌کشد اما فرجه‌های مدل را به سمت یک توزیع بد می‌برد. علاوه بر این از آنجایی که آسان‌تر noisy است، می‌توان این فاصله گرفتن به فاصله‌های نزدیک به Policy فرجه‌های منتقل کرد.

$$E_{\pi_{\theta}^{RL}}[ \log \left( \pi_{\theta}^{RL}(a|x) \right) ] : \text{از آنجایی که داریم مدل را با آسان‌تر می‌کنیم، ممکن است این اثر باعث شود که نتوانیم}$$

مدل در انجام تک‌های مدل MDP که منجر به بنابرین با بهبود یافتن گرفتن روی داده‌های دنباله‌ای MDP از گام‌های آخری که Policy آسان‌تر داده شده، انتظاس می‌دهد به آن داده و می‌کشد کردن این ایده باعث می‌شود عملکر مدل روی تک‌های MDP را نیز خوب کند و این را

(ج) زوای ورودی  $\theta(z)$  با استفاده از سیاست به دست می‌آید و بنابراین تابعی از آن است. بنابراین مشتق  $\theta(z)$  نسبت به  $\theta$  ضرورت و مقدار دارد.

$$\begin{aligned} L_{\theta} &= E_{\pi_{\theta}}[G_e] & \nabla_{\theta} L_{\theta} &= \nabla_{\theta} E_{\pi_{\theta}}[G_e] = \nabla_{\theta} \int \pi_{\theta}(z) G_e dz = \int \nabla_{\theta} \pi_{\theta}(z) G_e dz \\ & & &= \int \pi_{\theta}(z) \frac{1}{\pi_{\theta}(z)} \nabla_{\theta} \pi_{\theta}(z) G_e dz = \int \pi_{\theta}(z) \nabla_{\theta} \log \pi_{\theta}(z) G_e dz \\ & & &= E_{\pi_{\theta,e}}[\nabla_{\theta} \log \pi_{\theta}(z) G_e] \end{aligned} \quad (د)$$

(۲) تابع هدف زیر را به صورت یک لاگرانژی بنویسید. از آن جهت  $\lambda$  و  $\beta$  که ضرایب هستند، مشتق بگیرید و جواب بگیرید.

$$L = \int (r_\theta(x, y) - \beta \log \frac{\pi_\theta(y|x)}{\pi_{ref}(y|x)}) \pi_\theta(y|x) p(x) dx dy + \lambda (1 - \int \pi_\theta(y|x) p(x) dx dy)$$

↓

قید توزیع احتمال بودن  $\pi_\theta$

با این مشتق گرفتن از تابع  $\pi_\theta(y|x)$  مقدار  $x, y$  را ثابت فرض میکنیم.

$$\frac{\partial L}{\partial \pi_\theta(y|x)} = (r_\theta(x, y) - \beta \log \frac{\pi_\theta(y|x)}{\pi_{ref}(y|x)}) p(x) - \beta p(x) - \lambda p(x) = 0$$

$$\Rightarrow \frac{\pi_\theta(y|x)}{\pi_{ref}(y|x)} = \exp\left(\frac{1}{\beta} r_\theta(x, y)\right) = \exp\left(-\frac{\lambda + \beta}{\beta}\right)$$

$$\Rightarrow \pi_\theta(y|x) = \exp\left(-\frac{\lambda + \beta}{\beta}\right) \pi_{ref}(y|x) \exp\left(\frac{1}{\beta} r_\theta(x, y)\right) \quad ; \quad Z(x) = \frac{1}{\exp\left(-\frac{\lambda + \beta}{\beta}\right)}$$

$$\Rightarrow \pi_\theta(y|x) = \frac{1}{Z(x)} \pi_{ref}(y|x) \exp\left(\frac{1}{\beta} r_\theta(x, y)\right)$$

(ب)

$$\frac{\pi_\theta(y|x)}{\pi_{ref}(y|x)} Z(x) = \exp\left(\frac{1}{\beta} r_\theta(x, y)\right) \Rightarrow \log \frac{\pi_\theta(y|x)}{\pi_{ref}(y|x)} + \beta \log Z(x) = \frac{r_\theta(x, y)}{\beta}$$

$$\Rightarrow \beta \log \frac{\pi_\theta(y|x)}{\pi_{ref}(y|x)} + \beta \log Z(x) = r_\theta(x, y) \quad L_R(\theta, D) = -E_{(x, y_w, y_e) \sim D} [\log \delta(r_\theta(x, y_w) - r_\theta(x, y_e))]$$

$$L_{DPO}(\pi_\theta; \pi_{ref}) = -E_{(x, y_w, y_e) \sim D} \left[ \log \delta\left(\beta \log \frac{\pi_\theta(y_w|x)}{\pi_{ref}(y_w|x)} - \beta \log \frac{\pi_\theta(y_e|x)}{\pi_{ref}(y_e|x)}\right) \right]$$

$$\nabla_\theta L_{DPO}(\pi_\theta; \pi_{ref}) = -E_{(x, y_w, y_e) \sim D} \left[ \frac{\partial \log(\delta(A))}{\partial \theta} \right]$$

(ج)

$$A = \beta \left( \log \pi_\theta(y_w|x) - \log \pi_{ref}(y_w|x) \right) - \beta \left( \log \pi_\theta(y_e|x) - \log \pi_{ref}(y_e|x) \right)$$

$$= -E_{(x, y_w, y_e) \sim D} \left[ (1 - \delta(A)) \frac{\partial A}{\partial \theta} \right] = -\beta E_{(x, y_w, y_e) \sim D} \left[ \delta\left(\beta \log \frac{\pi_\theta(y_w|x)}{\pi_{ref}(y_w|x)} - \beta \log \frac{\pi_\theta(y_e|x)}{\pi_{ref}(y_e|x)}\right) (\nabla_\theta \log \pi_\theta(y_w|x) - \nabla_\theta \log \pi_\theta(y_e|x)) \right]$$

$$= -\beta E_{(x, y_w, y_e) \sim D} \left[ \delta(\hat{r}_\theta(x, y_e) - \hat{r}_\theta(x, y_w)) (\nabla_\theta \log \pi_\theta(y_w|x) - \nabla_\theta \log \pi_\theta(y_e|x)) \right]$$