a) ۱ استفاده از $\log \sigma\left(r_\theta(x,y_w) - r_\theta(x,y_L)\right)$ باعث می‌شود به هنگام

زیاد شدن اختلاف امتیاز $r_\theta(x,y_w), r_\theta(x,y_L)$ گرادیان به سمت صفر شدن

میرود زیرا $\log \sigma(\cdot)$ در مقادیر زیاد حرارت یابنش صغر می‌کند و به همین علت از زیاد شدن

اختلاف امتیاز این دو جلوگیری می‌شود.

b)

$$\text{objective}(\phi) = E_{(x,y)\sim D_{\pi_\phi^{RL}}}\left[r_\theta(x,y) - \beta \log\left(\frac{\pi_\phi^{RL}(y|x)}{\pi^{SFT}(y|x)}\right)\right]$$

$$+ \gamma\, E_{x\sim D_{\text{pretrain}}}\left[\log\left(\pi_\phi^{RL}(x)\right)\right]$$

$r_\theta(x,y) \longrightarrow$ جهت حداکثر شدن Reward برای هر جفت $(x,y)$ به نوعی alignment مورد نظر.

$\log\left(\frac{\pi_\phi^{RL}(y|x)}{\pi^{SFT}(y|x)}\right) \longrightarrow$ برای جلوگیری از فراموشی زبان توسط مدل و تولید توکن‌های یک جمله مهم حداکثر کردن $\pi_\phi^{RL}(y|x)$ به نوعی باعث می‌شود فنم از توزیع اولیه $\pi^{SFT}(y|x)$ optimization فاصله نگیرد.

$\log\left(\pi_\phi^{RL}(x)\right) \longrightarrow$ برای جلوگیری از خراب شدن کارایی مدل در تسک‌های NLP که در مدل pretraining وجود دارند.

c) با توجه به اینکه $(x,y)\sim D_{\pi_\phi^{RL}}$ می‌آید و به نوعی دیتا static نیست، وابسته به $\pi_\phi^{RL}$ است و از $r_\theta(x,y)$ به همان مستقل گیری بر حسب $\phi$ صفر نشود.

$$l_\theta = \mathbb{E}_{\pi_\theta}[G_t]$$

$$\nabla_\theta l_\theta = \nabla_\theta \mathbb{E}_{\pi_\theta(\tau),\tau}[G_t] = \nabla_\theta \int \pi_\theta(\tau) G_\tau \, d\tau$$

$$= \int \nabla_\theta \pi_\theta(\tau) G_\tau \, d\tau = \int \pi_\theta(\tau) \frac{1}{\pi_\theta(\tau)} \nabla_\theta \pi_\theta(\tau) G_\tau \, d\tau$$

$$= \int \pi_\theta(\tau) \nabla_\theta \log \pi_\theta(\tau) G_\tau \, d\tau$$

$$= \mathbb{E}_{\pi_\theta,\tau}\left[ \nabla_\theta \log \pi_\theta(\tau) G_\tau \right]$$

$$= \mathbb{E}_{\pi_\theta,\tau}\left[ \left( \sum_t \nabla_\theta \log \pi_\theta(a_t \mid s_t) G_t \right) \right]$$

برای محاسبه‌ی $KL$ نیاز هست که توزیع $\pi_\theta^{RL}(y|m)$ ، $\pi^{SFT}(y|m)$ را داشته باشیم. اما در عمل ما این مقادیر را از $(y_i)$ هم نقاط نداریم. به همین علت باید بیاییم تخمین‌زنند‌ه‌ی این $KLD$ را تخمین زد. (.) عضو گیری از $\pi^{RL}(y|m)$).

$$D_{KL}(q \| P) \approx \frac{1}{N} \sum_{i=1}^{N} \log \frac{q(m_i)}{P(m_i)} \quad , \quad m_i \sim q(m)$$

این تخمین‌زنند‌ه unbiased هست اما زیادی variance دارد.

$$r = \frac{P(m)}{q(m)} \implies k_1 = -\log r$$

این تخمین به طور معمول مقادیر منفی لاست که به صورت شهودی منطقی نیست.

این تخمین به طور متوسط فضای عوضی‌ها منفی می‌کند. در حالی که $KL$ منفی نمی‌شود. $\sigma^2_r$ variance دارد. به جای $(.)$ اما تخمین‌زنند‌ه زیر مقادیر منفی نمی‌دهد و $\sigma^2_r$ variance کمتری دارد.

$$D_{KL}(q \| P) \approx \frac{1}{N} \sum \frac{1}{2} \left[ \log q(n_i) - \log P(n_i) \right]^2, \quad n_i \sim q(n)$$

این تخمین biased هست اما دارای Variance کمتری است.

$$K_2 = \frac{1}{2} (\log r)^2$$

$$E_q[K_2] = E_q\left[ \frac{1}{2} \log^2 r \right] \quad \underset{f = \frac{1}{2} \log^2 n}{\overset{f\text{-divergence}}{\rule{3cm}{0.4pt}}} \quad f(n) = -\log n \quad , \quad f''(1) = 1$$

تخمین بی بایاس ولی با variance بالا :

$$D_{KL}(q \| P) \approx \frac{1}{N} \sum_{i=1}^{N} \exp\left\{ \frac{1}{2} \left[ \log q(n_i) - \log P(n_i) \right]^2 - 1 - \left[ \begin{matrix} \log q(n_i) \\ -\log P(x_i) \end{matrix} \right] \right.$$

$$\max_{\sim} \; E_{n \sim D, \, y \sim \pi_\theta(y|m)} \left[ r_\phi(n,y) \right] - \beta D_{KL} \left[ \pi_\theta(y|m) \| \pi_{ref}(y|m) \right]$$

Langrangian :

$$\max_{\sim} \; E_{n \sim D, \, y \sim \pi_\theta(y|m)} \left[ r_\phi(m,y) - \beta D_{nL}\left( \pi_\theta(y|m) \| \pi_{ref}(y|m) \right) \right]$$
$$+ \lambda \left( 1 - \int \pi_\theta(y|m) P(m) \, dm \, dy \right)$$

$$= \int \left( r_\phi(m,y) - \beta \log \frac{\pi_\theta(y|m)}{\pi_{ref}(y|m)} \right) \pi_\theta(y|m) P(m) \, dm \, dy$$
$$+ \lambda \left( 1 - \int \pi_\theta(y|m) P(m) \, dm \, dy \right)$$

$$\frac{\partial L}{\partial \pi_\theta(m|y)} = \left( r_\phi(m,y) - \beta \log \frac{\pi_\theta(y|m)}{\pi_{ref}(y|m)} \right) P(m) - \beta P(m) - \lambda P(m) = ?$$

$$\implies \frac{\pi_\theta(y|m)}{\pi_{ref}(y|m)} = e^{\frac{1}{\beta} r_\phi(m,y) - \frac{\lambda + \beta}{\beta}}$$

$$\implies \pi_\theta(y|m) = e^{-\frac{\lambda + \beta}{\beta}} \pi_{ref}(y|m) \, e^{\frac{1}{\beta} r_\phi(m,y)}$$

$$= \frac{1}{Z(m)} \pi_{ref}(y|m) \, e^{\frac{1}{\beta} r_\phi(m,y)}$$

تعریف می‌کنیم:

$$\pi_\theta(y|n) = \frac{1}{Z(n)} \pi_{ref}(y|n) e^{\frac{1}{\beta} r_\varphi(n,y)}$$

خواص آن را؟ :

$$\log \frac{Z(n) \pi_\theta(y|n)}{\pi_{ref}(y|n)} = \frac{1}{\beta} r_\varphi(n,y)$$

$$\Rightarrow \quad r_\varphi(n,y) = \beta \log Z(n) + \beta \log \frac{\pi_\theta(y|n)}{\pi_{ref}(y|n)} \qquad \text{①}$$

$$L_R(r_\varphi, D) = - E_{(n, y_w, y_L) \sim D} \left[ \log \sigma(r_\varphi(n,y_w) - r_\varphi(n,y_L)) \right] \qquad \text{②}$$

$$\text{①, ②} \Rightarrow L_{DPO}(\pi ; \pi_{ref}) = - E_{n, y_w, y_L \sim D} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(y_w|n)}{\pi_{ref}(y_w|n)} - \beta \log \frac{\pi_\theta(y_L|n)}{\pi_{ref}(y_L|n)} \right) \right]$$

$$\nabla_\theta L_{DPO}(\pi_\theta ; \pi_{ref}) = \nabla_\theta E_{n, y_w, y_L \sim D} \left[ \log \sigma \left( \underbrace{\beta \log \frac{\pi_\theta(y_w|n)}{\pi_{ref}(y_w|n)} - \beta \log \frac{\pi_\theta(y_L|n)}{\pi_{ref}(y_L|n)}}_{u} \right) \right]$$

$$= - E_{n, y_w, y_L \sim D} \left[ \frac{\sigma'(u)}{\sigma(u)} \nabla_\theta \left( \beta \log \frac{\pi_\theta(y_L|n)}{\pi_{ref}(y_L|n)} - \beta \log \frac{\pi_\theta(y_w|n)}{\pi_{ref}(y_w|n)} \right) \right]$$

$$= - E_{n, y_w, y_L \sim D} \left[ \beta \sigma \left( \beta \log \frac{\pi_\theta(y_w|n)}{\pi_{ref}(y_w|n)} - \beta \log \frac{\pi_\theta(y_L|n)}{\pi_{ref}(y_L|n)} \right) \left( \nabla_\theta \log \pi(y_w|n) - \nabla_\theta \log \pi(y_L|n) \right) \right]$$

Reward مربوط به مدل زبانی $\pi_\theta(y|m)$ و مدل زبانی Reference
$\pi_{ref}(y|m)$ یعنی می‌توان بصورت مستقیم بدان نیاز:
Policy را Optimize کرد. Reward function

$$\beta \sigma \left( \beta \log \frac{\pi_\theta(y_w|m)}{\pi_{ref}(y_w|m)} - \beta \log \frac{\pi_\theta(y_L|m)}{\pi_\rho(y_L|m)} \right) \cdot \left[ \begin{array}{c} \\ \end{array} \right]$$

$$\nabla_\theta \log \pi(y_w|m) :$$  جهت افزایش likelihood خروجی مناسب $y_w$

$$-\nabla_\theta \log \pi(y_w|m) :$$  جهت کاهش likelihood خروجی نامناسب $y_L$