$$\begin{bmatrix} \sum_{(x,y+,y-)\in D} \log \sigma\left( r_\varphi(x,y+) - r_\varphi(x,y-) \right) \end{bmatrix} \left.\begin{matrix} \end{matrix}\right\} \begin{matrix} O.F. \\ \text{for the} \\ RM \end{matrix}$$

$$(x, y_+, y_-)$$

$RM \quad \varphi \leftarrow \varphi_*^* $

maximize $\downarrow \varphi^*$

To train the policy

$$\theta^* \leftarrow \arg\max \, E_{\pi_\theta(y|x)} \left[ \underbrace{r_{\varphi^*}(x,y)}_{\uparrow} - \underbrace{\beta KL\left( \pi_\theta(y|x) \| \pi_{ref}(y|x) \right)}_{\text{Regularized Reward}} \right]$$

RM

PM

Ref.

## DPO

- No RM (Explicit)
- NO RL
- Policy model will act as a RM

$$r^* \leftarrow \arg\max_{r} \sum_{(x,y+,y-)\in D} \log \sigma\left( r(x,y+) - r(x,y-) \right)$$

Optimal Policy model
$$\Rightarrow \pi^* \leftarrow \arg\max_{\pi} E_{\pi(y|x)} r^*(x,y) - \beta \cdot KL\left( \pi^*(y|x) \| \pi_{ref}(y|x) \right) - $$
$$\text{s.t.} \quad \sum_{y\in Y} \pi^*(y|x) = 1$$

$$\nabla_{\pi(y_0|x)} \mathcal{L}(\pi, \lambda) = 0$$

$$\mathcal{L}(\pi, \lambda) = \sum_{y\in Y} \pi^*(y|x) r^*(x,y) - \sum_{y\in Y} \pi^*(y|x) \log \frac{\pi^*(y|x)}{\pi_{ref}(y|x)} + \lambda\left( \sum_{y\in Y} \pi^*(y|x) - 1 \right)$$

$$\nabla_{\pi(y_0|x)} \mathcal{L} = r^*(x,y) - \left[ 1 + \log \frac{\pi^*(y|x)}{\pi_{ref}(y|x)} \right] + \lambda = 0$$

$$\Rightarrow r^*(x,y) + \lambda - 1 = \log \frac{\pi^*(y|x)}{\pi_{ref}(y|x)}$$

$$\therefore r^*(x,y) + \bar{\lambda} = \log \frac{\pi^*}{\pi_{ref}} \Leftarrow \text{(1)} \qquad \text{where } \bar{\lambda} = \lambda - 1$$

$$\Rightarrow e^{(r^*(x,y) + \bar{\lambda})} = \frac{\pi^*}{\pi_{ref}}$$

$$\Rightarrow \pi^*(y|n) = \pi_{ref} \cdot \exp(r^*(n,y) + \bar{\lambda}) \quad \text{---} \textcircled{1}$$

This is the optimal policy in trm of optimal reward,

$$\sum_{y \in Y} \pi^*(y|n) = 1 \quad \Rightarrow \sum \pi_{ref}(y|n) \exp(r^*(n,y) + \bar{\lambda}) = 1$$

$$\Rightarrow \exp(\bar{\lambda}) = \frac{1}{\sum \pi_{ref}(y|n) \exp(r^\circ(n,y))} = \frac{1}{Z}$$

For $\textcircled{1}$, $\boxed{\pi^*(y|n) = \pi_{ref} \frac{\exp(r^*(n,y))}{Z}} \quad \text{---}\textcircled{1}$ policy model in trm of RM

We will write the reward in trm of the policy.

fm Eq $\textcircled{1}$ $\quad r^*(x,y) + \bar{\lambda} = \log \pi^* / \pi_{ref}$ $\qquad \bar{\lambda} = -\log Z \approx \log \frac{1}{Z}$

$$\Rightarrow r^*(n,y) = \log \pi^* / \pi_{ref} - \bar{\lambda}$$

$$\rightarrow \boxed{r^*(n,y) = \log \frac{\pi^*}{\pi_{ref}} - \log Z} \qquad \text{reward in trm of policy}$$

the parametric policy & reward

$$r_\theta(x,y) = \beta \log \frac{\pi_\theta(y|x)}{\pi_{ref}(y|x)} - \log Z_x(\theta) \quad =\!=\!\textcircled{3}$$

<u>Training the reward fine:</u>

$$(x, y+, y-) \quad \arg\max \sum_{(x,y+,y-) \in D} \log \sigma(r_\theta(x,y+) - r_\theta(x,y-)) \Longleftarrow$$

For Eq $\textcircled{3}$

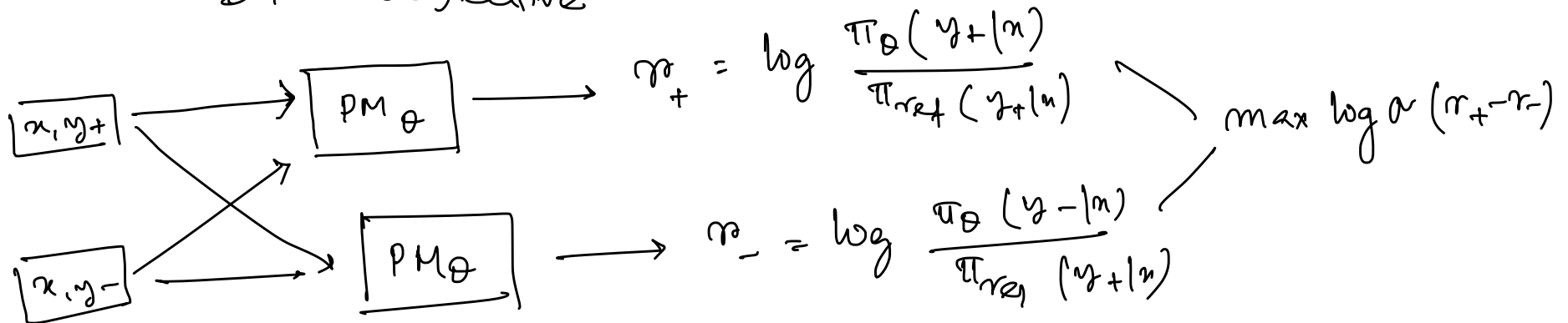$$r_\theta(x,y+) = \beta \log \frac{\pi_\theta(y+|x)}{\pi_{ref}(y+|x)} - \log Z_x(\theta)$$

$$r_\theta(x,y-) = \beta \log \frac{\pi_\theta(y-|x)}{\pi_{ref}(y-|x)} - \log Z_x(\theta)$$

$$\log \sigma\left(\left[\beta \log \frac{\pi_\theta(y+|x)}{\pi_{ref}(y+|x)} - \log \cancel{Z_x(\theta)}\right] - \left[\beta \log \frac{\pi_\theta(y-|x)}{\pi_{ref}(y-|x)} - \log \cancel{Z_x(\theta)}\right]\right)$$

$$= \quad \log \sigma\left(\beta\left[\log \frac{\pi_\theta(y+|x)}{\pi_{ref}(y+|x)} - \log \frac{\pi_\theta(y-|x)}{\pi_{ref}(y-|x)}\right]\right)$$

$$=$$

$$= \log \frac{\exp\left(\beta \log \frac{\pi_\theta(y_+|x)}{\pi_{ref}(y_+|x)}\right)}{\exp\left(\beta \log \frac{\pi_\theta(y_+|x)}{\pi_{ref}(y_+|x)}\right) + \exp\left(\beta \log \frac{\pi_\theta(y_-|x)}{\pi_{ref}(y_-|x)}\right)}$$

DPO Objective



$$\frac{\pi_\theta(y_+|x)}{\pi_{ref}(y_+|x)} \approx 0.8 \qquad y_+ \to \qquad y_+ \sim .1$$

$$\boxed{\frac{.9}{.8} \approx \frac{.1+}{.1}}$$

$$\beta \to (0.3 - 0.003)$$

DPO is prone to generating a biased-policy that favours.
<u>out-of-distribution responses</u>