## RLHF objective:

$$\max_{\pi} \; \underset{x,y}{\mathbb{E}} \left[ \underbrace{r(x,y)}_{frozen} - \beta D_{kL} \left( \pi(y|x) \; || \; \pi_{ref}(y|x) \right) \right)$$

non-differentiable → (pointing to $\pi(y|x)$)

current aligned LLM → (pointing to $\pi(y|x)$)

SFT (instruction-tuned LLM) → (pointing to $\pi_{ref}(y|x)$)

$(x,y)$
data used for SFT is
different than that used for RLHF,
but come from same distribution

- why do we need RL?

## DPO (direct preference optimization):

→ no explicit reward model
→ not going to sample outputs $y|x$ from the model
        ↳ "rollouts"
→ "preference tuning"

$$\max_{\pi} \; \underset{x,y}{\mathbb{E}} \left[ r(x,y) - \beta \log \frac{\pi(y|x)}{\pi_{ref}(y|x)} \right]$$

$$= \min_{\pi} \; \mathbb{E}_{x,y} \left[ \log \frac{\pi(y|x)}{\pi_{ref}(y|x)} - \frac{1}{\beta} r(x,y) \right]$$

let's introduce a new policy $\pi^*$ that incorporates the reward term as well as $\pi_{ref}$

$$\pi^*(y|x) = \frac{\pi_{ref}(y|x) \exp\left(\frac{1}{\beta} r(x,y)\right)}{\underbrace{\sum_{y} \pi_{ref}(y|x) \exp\left(\frac{1}{\beta} r(x,y)\right)}}$$

$Z(x)$, normalizer / partition function

Substitute $Z(x)$ into our objective:

$$\min_{\pi} \; \mathbb{E}_{x,y} \; \log \left[ \frac{\pi(y|x)}{\frac{1}{Z(x)} \pi_{ref}(y|x) \exp\left(\frac{1}{\beta} r(x,y)\right)} - \log Z \right]$$

$$= \min_{\pi} \; \mathbb{E}_{x,y} \; \log \left[ \frac{\pi(y|x)}{\pi^*(y|x)} \right] - \log Z$$

$\hookrightarrow$ KL div

$$= \min_{\pi} \; \mathbb{E}_{x} \; D_{KL}\left( \pi(y|x) \, || \, \pi^*(y|x) \right) - \log Z$$

KL div. is minimized at $0$ when $\pi(y|x) = \pi^*(y|x)$

$$\pi(y|x) = \pi^*(y|x) = \frac{\pi_{ref}(y|x)\exp(\frac{1}{\beta}r(x,y))}{Z(x)}$$

$\underbrace{\phantom{\pi}}_{\text{optimal policy}}$

Solve the above for $r(x,y)$

$$r(x,y) = \boxed{\beta \log \frac{\pi^*(y|x)}{\pi_{ref}(y|x)} + \beta \log Z}$$

Bradley-Terry pref model:

$$p(y_w > y_L | x) = \frac{\exp(r(x,y_w))}{\exp(r(x,y_w)) + \exp(r(x,y_L))}$$

Substitute reward function:

$$p(y_w > y_L | x) = \frac{1}{1 + \exp\left(\beta \log \frac{\pi^*(y_L|x)}{\pi_{ref}(y_L|x)} - \beta\log\frac{\pi^*(y_w|x)}{\pi_{ref}(y_w|x)}\right)}$$

Convert to loss fn (neg. log likelihood)

$$L_{DPO}(\pi_\theta, \pi_{ref}) = -\underset{x,y_w,y_L}{\mathbb{E}} \log \sigma\left(\beta \log \frac{\pi_\theta(y_w|x)}{\pi_{ref}(y_w|x)} - \beta \log \frac{\pi_\theta(y_L|x)}{\pi_{ref}(y_L|x)}\right)$$

$\downarrow$
aligned model we are training

nice properties of DPO!
  — no explicit reward model
  — no need for rollouts from the policy

$$\underset{X}{\text{O}} \; X, Y_w, Y_L$$

SFT
LLM

$\xrightarrow{\hspace{3cm}}$

finetune on
pref. judgments
using DPO loss above

DPO
LLM