

25.03.04

Reward Model (usually a small language model)

$r_\theta(x, y)$
prompt \swarrow \searrow response

$X \xrightarrow{\text{SFT}} y_1 > y_2 > y_3 > y_4$

MLE,

$$\max \frac{1}{\binom{4}{2}} \sum_{(y_w, y_l)} \log \left[\frac{\exp(r_\theta(x, y_w))}{\exp(r_\theta(x, y_w)) + \exp(r_\theta(x, y_l))} \right]$$

\downarrow \downarrow
winning response losing response

New Model

For DPD:

Reward Modeling with $r_\theta(x, y) = \log \left[\frac{p_{RL}(y|x)}{p^{\text{SFT}}(y|x)} \right]$
SFT Model \swarrow