口 切换模式

关于作者

回答

37

关注他

AI星际漫游

文章

76

vx公众号[AI星际漫游]

登录/注册

关注者

141

● 发私信

# dpo缺点及DPOP解决方案



AI星际漫游

vx公众号[AI星际漫游]

■ 来自专栏·每天学点深度学习 >

11 人赞同了该文章 >

一、 存在的问题: 标准直接偏好优化+ (DPO) 的失败模式

在大型语言模型 $^+$  (LLM) 的对齐中,目标是使模型的行为符合人类的偏好。直接偏好优化 (DPO) 是一种很有前景的技术,它试图直接通过偏好数据(即给定提示  $\boldsymbol{x}$ ,人类更喜欢回答  $\boldsymbol{y_w}$  而不是  $\boldsymbol{y_t}$ )来优化模型,而无需像传统的基于人类反馈的强化学习 (RLHF) 那样显式地训练一个奖励模型。

#### 1. 标准 DPO 的目标和机制

DPO 的核心思想是,存在一个隐式的奖励函数  $r^*(x,y)$ ,人类的偏好是根据这个奖励函数生成的。一个理想的模型  $\pi^*$  会根据这个奖励函数来分配概率:

$$\pi^*(y|x) \propto \pi_{\mathrm{ref}}(y|x) \exp(eta r^*(x,y))$$

#### 其中:

- ・  $\pi_{\mathrm{ref}}(y|x)$  是参考模型的概率(通常是经过监督微调 $^+$  (SFT) 的模型)。
- $\beta$  是一个超参数,控制奖励函数对概率分布的影响程度。

DPO 的目标是找到一个策略  $\pi_{\theta}$  (模型参数为  $\theta$ ),使其能够最好地解释观察到的偏好数据。它通过最大化以下目标函数(或最小化其负数形式的损失函数)来实现这一点:

$$L_{ ext{DPO}}(\pi_{ heta}; \pi_{ ext{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim D} \left[ \log \sigma \left( eta \log rac{\pi_{ heta}(y_w | x)}{\pi_{ ext{ref}}(y_w | x)} - eta \log rac{\pi_{ heta}(y_l | x)}{\pi_{ ext{ref}}(y_l | x)} 
ight) 
ight]$$

# 这里:

- ・D 是人类偏好数据集,包含  $(x,y_w,y_l)$  的三元组。
- ・ $\pmb{y_w}$  是偏好的 (winning) 回答, $\pmb{y_l}$  是不偏好的 (losing) 回答。
- $\cdot$   $\sigma(z)=1/(1+e^{-z})$  是 sigmoid 函数。
- $m{\hat{r}}_{ heta}(y|x)$  可以看作是模型  $m{\pi}_{ heta}$  相对于参考模型  $m{\pi}_{ ext{ref}}$  对回答  $m{y}$  给出的隐式奖励 $m{\hat{r}}_{ heta}(x,y)$ 。

所以, DPO 损失可以更简洁地写为:

$$L_{ ext{DPO}} = -\mathbb{E}_{(x,y_w,y_l)\sim D}\left[\log\sigma(\hat{r}_{ heta}(x,y_w) - \hat{r}_{ heta}(x,y_l))
ight]$$

DPO 试图最大化偏好回答  $y_w$  的隐式奖励与不偏好回答  $y_l$  的隐式奖励之间的差异。

# 2. DPO 的失败模式: 降低偏好样本的似然性

尽管标准的 DPO 损失函数旨在提高偏好回答  $y_w$  相对于不偏好回答  $y_l$  的"相对概率"(即  $\frac{\pi_{\theta}(y_w|x)}{\pi_{\theta}(y_l|x)}$  增大),但它可能无意中导致模型降低对偏好回答  $y_w$  本身的绝对似然性  $\log \pi_{\theta}(y_w|x)$ 。

# 理论原因和机制:

・ 关注相对差异而非绝对值:
 \hat{r}\_{\text{theta}(x, y\_l) 的}
 ▲ 赞同 11
 ▼ ● 添加评论
 ④ 分享
 ● 喜欢
 ◆ 收藏
 ● 申请转载



证 \log \pi\_{\theta}(y\_w|x) 会增加。

- 低編辑距离问题: 这种失败模式在偏好对 (y\_w, y\_l) 之间文本差异很小(即编辑距离低)时 尤为突出。
  - 假设 y\_w = S\_1 S\_2 \dots S\_k W\_1 \dots W\_m 和 y\_l = S\_1 S\_2 \dots S\_k L\_1 \dots L\_p, 其中 S\_i 是共享的前缀词元, W\_j 和 L\_j 是不同的后续词元。
  - ・ 为了最大化  $y_w$  和  $y_v$  之间的奖励差异,DPO 的梯度可能会集中在使模型增加区分性词元( $w_i$ )的概率,并降低  $y_v$  中对应词元( $w_i$ )的概率。
  - · 然而,DPO 在增加区分词元概率的同时,**可能会减少共享前缀 S\_i 或 y\_w 中后续词元W\_j 的概率**。这可以被理解为一种"拆东墙补西墙"的效应。如果对共享部分的概率惩罚过大,即使区分部分的相对概率增加了,偏好序列 y\_w 的整体(绝对)对数似然性也可能下降。
- 梯度分析: 考虑 DPO 损失对模型输出 logits 的梯度。对于偏好序列 y\_w 中的一个词元 y\_{w,t}, 其梯度的贡献大致为:

 $\label{logit} $$ \frac{\left(\frac{partial L_{\text{DPO}}}{\left(y_{w,t}\right)}\right) \cdot \left(\frac{1 - \sigma(\left(1 - \frac{1}{pi_{\text{text}}}\right)}{1 - \frac{1}{pi_{\text{text}}}(y_{w,t}|x, y_{w,1:t-1})} - \frac{1}{pi_{\text{text}}}(y_{w,t}|x, y_{w,1:t-1})} \right) } $$$ 

#### 其中

 $\label{eq:def:Delta} $$ \Pr_{\hat{r}_{\hat{x}, y_w} - \hat{r}_{\hat{x}, y_l) } $$ $$ \Pr_{\hat{x}, y_l} (x, y_w) - \hat{r}_{\hat{x}, y_l}$ 

对于不偏好序列 y\_l 中的一个词元 y\_{l,t}, 其梯度的贡献大致为:

 $\label{logit} $$ \frac{po}{}{\left(y_{l,t}}\right) \cdot \left(\int_{t_{l,t}} \po(t_{l,t})} \po(t_{l,t})} \left(\int_{t_{l,t}} \po(t_{l,t})} \po(t_{l,t}) \right) \left(\int_{t_{l,t}} \po(t_{l,t}) \po(t_{l,t})} - \frac{1}{\left(\int_{t_{l,t}} \po(t_{l,t})} \po(t_{l,t}) \po(t_{l,t})} \right) \left(\int_{t_{l,t}} \po(t_{l,t}) \po(t_{l,t}) \po(t_{l,t})} \po(t_{l,t}) \po(t_{l,t}) \po(t_{l,t})} \po(t_{l,t}) \po(t_{l,t}) \po(t_{l,t}) \po(t_{l,t})} \po(t_{l,t}) \po(t_{l,t}) \po(t_{l,t}) \po(t_{l,t})} \po(t_{l,t}) \po(t_{l,t}) \po(t_{l,t}) \po(t_{l,t}) \po(t_{l,t}) \po(t_{l,t})} \po(t_{l,t}) \po(t_{l,t})$ 

虽然这些梯度试图增加  $y_w$  的相对奖励并减少  $y_w$  的相对奖励,但当  $y_w$  和  $y_w$  高度相似时,这些更新可能以复杂的方式相互作用,导致  $y_w$  的整体概率下降。DPO 可能只关注少数区分性词元的 logits。

#### 二、解决方案: DPO-Positive+ (DPOP)

为了解决标准 DPO 方法中存在的降低偏好样本似然性的问题,作者设计了一种新的损失函数和训练程序,称为 DPO-Positive (DPOP)。

# 1. DPOP 的目标

DPOP 的核心目标是**在保持 DPO 区分偏好和不偏好回答能力的同时,避免或减轻对偏好回答 y\_w 绝对似然性的不必要惩罚,甚至促进其提升**。 它通过在损失函数中引入一个新的惩罚项来实现这一目标,该惩罚项专门用于激励模型维持偏好样本的高对数似然性。

#### 2. DPOP 损失函数

DPOP 完整损失函数如下:

 $L_{\text{DPOP}}(\pi_{\hat{p}_{\text{theta};pi_{\text{theta$ 

# 这里:

- · \pi\_{\theta} 是当前正在优化的模型。
- · \pi\_{\text{ref}} 是参考模型(通常是经过 SFT 的模型)。
- (x, y\_w, y\_l) 分别是提示、偏好的回答和不偏好的回答。
- \beta 是控制 KL 正则化强度的超参数,与 DPO 中的定义一致。
- \lambda > 0 是一个新的超参数,用于控制新增惩罚项的权重。
- ・ \sigma 是 sigmoid 函数。
- 核心的新增惩罚项是 -\lan {\pi\_{\theta}(y\_w|x)}\right



#### 3. DPOP 如何解决失败模式

DPOP 通过其新增的惩罚项来解决 DPO 的失败模式:

- 激励维持偏好样本的似然性: 惩罚项 \text{max}\left(0, \log\frac{\pi\_{\text{ref}}(y\_w|x)} {\pi\_{\text{ref}}(y\_w|x)}\right) 的作用机制如下:
  - 当 \pi\_{\theta}(y\_w|x) \ge \pi\_{\text{ref}}(y\_w|x) 时(即当前模型对偏好回答 y\_w 的似然性不低于参考模型时), \log\frac{\pi\_{\text{ref}}(y\_w|x)}{\pi\_{\text{max}}(0, \dots) 项为 0, 惩罚项不起作用。
  - 当 \pi\_{\theta}(y\_w|x) < \pi\_{\text{ref}}(y\_w|x) 时(即当前模型降低了偏好回答 y\_w 的似然性), \log\frac{\pi\_{\text{ref}}(y\_w|x)}{\pi\_{\text{ref}}(y\_w|x)}} > 0, 此时惩罚项变为一个正值,通过损失函数的最小化过程,模型将被激励去提高 \pi\_{\theta}(y\_w|x),使其至少不低于(或努力恢复到) \pi\_{\text{ref}}(y\_w|x) 的水平。

展开后,对于词元  $t_k$ (假设其在词汇表中的索引为 i),其梯度关于第 j 个 logit \theta\_j 的 分量为: = \begin{cases} \lambda(1-s\_j^{{y\_w^{<k},x}}) + s\_j^{{y\_l^{<k},x}} - s\_j^{{y\_w^{<k},x}} & \text{if } i=j \ -(\lambda+1)s\_j^{{y\_w^{<k},x}} + s\_j^{{y\_l^{<k},x}} & \text{if } i=j \ -(\lambda+1)s\_j^{{y\_w^{<k},x}} + s\_j^{{y\_l^{<k},x}} & \text{if } i \ ine j \end{cases}\\ \text{if } s\_j^{{context},x}\\ \text{if } i \ ine j \end{cases}\\ \text{if } s\_j^{{context},x}\\ \text{if } i \ ine j \end{cases}\\ \text{if } s\_j^{{context},x}\\ \text{if } \ ine j \end{cases}\\ \text{if } s\_j^{{context},x}\\ \text{if } \ ine j \end{cases}\\ \text{if } s\_j^{{context},x}\\ \text{if } s\_j^{{context},x

• 确保偏好样本的高似然性:通过这种优化压力,模型不能再仅仅通过大幅降低不偏好样本的似然性(相比于降低偏好样本似然性的程度)来最小化损失。它必须同时确保偏好样本的似然性相对于参考模型保持在较高水平。

# 4. DPOP 中的隐式奖励参数化

DPOP 损失函数保留了在 Bradley-Terry 模型下拟合偏好数据的特性。 其隐式奖励参数化如下:

- 对于不偏好的回答 y\_I: \hat{r}\_{\text{DPOP}}(x, y\_I) = \beta \log\frac{\pi\_{\text{ref}}(y\_I|x)}{\pi\_{\text{ref}}(y\_I|x)} \lambda (这与标准 DPO 对 y\_I 的奖励形式相同)
- · 对于偏好的回答 y\_w:

# 5. 与对比损失<sup>+</sup>的联系 (Connection to Contrastive Loss)

尽管 DPOP 的主要动机是解决 DPO 的失败模式,但它也与对比学习中的对比损失 (contrastive loss) 有一定的联系。 DPOP 的损失函数形式可以被看作是带有特定间隔 (margin) 的对比损失,其中间隔 m = \log\frac{1}{\pi\_{\text{ref}}(y\_w|x)}。 标准 DPO 类似于缺少了对比损失中"相似点对"项和间隔项的公式,而 DPOP 通过引入新的惩罚项(可以类比为相似点对的项和间隔)弥补了这一不足。

通过这些机制,DPOP 旨在更稳健地从偏好数据中学习,避免标准 DPO 可能出现的性能退化,并在多种情况下展现出更优越的性能。

编辑于 2025-05-14 00:26·浙

# DPO 深度学习 (Deep Learning)



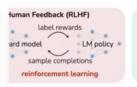
### 推荐阅读



### DPO和实现代码

1. 提出 DPO 的背景RLHF 的缺陷:基于人类反馈的强化学习(RLHF)过程复杂且不稳定,需要拟合奖励模型并使用强化学习微调大型无监督语言模型(LM)。DPO 的目的:通过利用奖励函数与最优策...

Evan Lee



从原理到实战技巧,详解DPO 的优缺点

黑木

# (详解) 言模型!

Direct Pre

引言在上了DPO设作为PPC于大模型齐。其中一个问题