

NLP_双关语

摘要

- LLM生成幽默句子的能力仍然存在局限 (ChatGPT的幽默生成能力仅限于创作25条不重复的笑话。)
- 工作简介
 - 我们提出一种多阶段课程偏好学习框架，同步优化双关结构偏好与幽默偏好。具体而言，我们改进了直接偏好优化 (DPO) 算法以解决多目标对齐难题。
- 新数据集
 - 为推进该领域发展，我们构建了包含2.1k条双关语及对应标注的中文双关数据集

简介

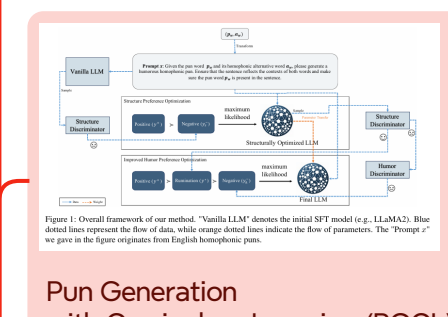
- 双关语
 - 双关语生成旨在根据给定的双关词与替代词组合 (如表所示) 生成具有幽默效果的语句，其核心目标包含双重维度：双关结构与幽默效果
- 为从新维度提升LLM能力，偏好学习/对齐已成为当前主流方法 (Casper等人，2023)
 - 然而，在双关生成任务中直接应用直接偏好优化 (DPO) (Rafailov等人，2023) 等偏好学习方法效率较低——这源于我们需在较小规模数据集上同步实现双关结构偏好与幽默偏好的双重对齐。
- 文章提出的解决办法：多阶段课程学习框架
 - 第一阶段优先学习双关结构偏好
 - 第二阶段进阶学习幽默偏好
 - 为解决灾难性的遗忘问题：在该阶段引入一种改进的幽默偏好对齐算法
 - 与标准DPO方法 (仅从正负样本对中学习) 不同，我们通过为每个训练样本添加第三元——即使用相同输入提示从第一阶段模型获得的额外生成输出——将其构建为三元组样本
 - 提出一种新损失函数，既能适配这种多目标对齐任务，又可缓解灾难性遗忘效应。
- 数据集
 - 英语双关语数据集: English pun dataset from SemEval 2017 task 7 (SemEval) (Miller et al., 2017)
 - 包含1298个同形双关语
 - 1098 个用双关语和替代词注释的谐音双关语
 - 中文数据集 ChinesePun
 - Table 2. Statistics of the ChinesePun dataset
 - Table 1. Statistics of the ChinesePun dataset
 - 我们提出了一种多阶段课程学习框架，结合创新的三元组偏好学习方法，显著提升了大语言模型生成幽默双关语的能力。
 - 我们发布了首个面向中文双关语生成任务的数据集，旨在推动中文幽默理解与生成领域的研究进展。
 - 我们在中英文双关数据集上进行了全面实验，结果充分证明了所提出方法的优越性。

相关工作

- 双关语模型上
 - 然而，上述所有方法都是专门为小规模模型设计的。据我们所知，我们的工作首次增强 LLM 的双关语生成能力。
- 偏好学习上
 - 直接偏好优化 (DPO) (Rafailov等人，2023) 通过配对偏好实现模型对齐，显著提升效率与稳定性。
 - 基于DPO和DSJC (Zhao等人，2022) 改进的统计拒绝采样优化 (RSO) (Liu等人，2023) 采用拒绝采样实现更有效的优化
 - Xu等人 (2023) 则提出渐进式对比训练课程，通过从简单到复杂的数据对逐步提升模型性能。
 - 本文工作创新
 - 我们创新性地提出多阶段课程学习框架和三元组偏好学习损失函数，专门针对双关语生成任务中的多目标偏好对齐稳定性进行优化

问题定义

- 同音异义
 - 我们系统的输入由双关词构成，包含一个双关词 (pw, 如weak) 及其对应替代词 (aw, 如 week)
- 同形异义
 - 在同形异义双关的情形下，当双关词在特定语境中具有逻辑连贯的双重含义时，我们遵循Tian等人 (2022) 的方法，采用相同表征形式 (即设 $DW = \Delta W$)
- 合格的生成结果需满足两个条件：符合双关结构 (如包含双关词) 且具有幽默性。
 - lift weights only on Saturday and Sunday because Monday to Friday are weak days
 - 虚假日和工作日
- 蓝色虚线表示数据流
- 黄色虚线表示参数流



详细做法

- 1多阶段“课程”学习 (Curriculum Learning——从易到难学习样本)
 - 1我们分两个阶段使用直接偏好优化 (DPO) 方法，以引导模型优化两个关键偏好：双关语结构和幽默。
 - 2结构偏好对齐
 - (1) 结构偏好优化模块 (上图)，用于在第一阶段增强大语言模型满足双关结构要求的能力
 - 2) 幽默偏好优化模块 (下图)，旨在第二阶段实现更具挑战性的幽默偏好对齐。
 - 对于输入pw和aw，通过模板转换为提示x
 - 1) 直接偏好优化 (DPO) 模块：该模块接收提示x和对应的替代词aw，通过DPO算法优化模型，使其生成更符合双关结构的输出。该模块的输出为模型参数θ。
 - 2) 幽默偏好优化 (HPO) 模块：该模块接收提示x和对应的替代词aw，通过HPO算法优化模型，使其生成更具幽默性的输出。该模块的输出为模型参数θ。
 - 正负样本生成——不满足5条规则其中之一-的样本做负样本
 - DPO损失
 - $\ell_{DPO} = -\log \pi_{\theta}(y^+ | x) + \log \pi_{\theta}(y^- | x)$ (1)
 - $\ell_{DPO} = -\log \pi_{\theta}(y^+ | x) + \log \pi_{\theta}(y^- | x)$ (2)
 - $\ell_{DPO} = -\log \pi_{\theta}(y^+ | x) + \log \pi_{\theta}(y^- | x)$ (3)
 - 使用与第一阶段相同的提示x，和正样本y+并用第一阶段优化后的LM生成候选句子
 - 判别器—— fine-tune RoBERTa-large model
 - 损失函数
 - $\ell_{DPO} = -\log \pi_{\theta}(y^+ | x) + \log \pi_{\theta}(y^- | x)$ (4)
 - $\ell_{DPO} = -\log \pi_{\theta}(y^+ | x) + \log \pi_{\theta}(y^- | x)$ (5)
 - $\ell_{DPO} = -\log \pi_{\theta}(y^+ | x) + \log \pi_{\theta}(y^- | x)$ (6)
 - 参数解释
 - 1. 模型参数: 模型参数θ, 用于生成输出y。
 - 2. 提示: 提示x, 用于生成输出y。
 - 3. 替代词: 替代词aw, 用于生成输出y。
 - 4. 损失函数: 损失函数ℓ, 用于衡量模型输出y与目标y+和y-之间的差异。
 - 5. 判别器: 判别器D, 用于判断输出y是否符合双关结构要求。

方法

- 改进的具备幽默偏好的DPO
 - 由于训练数据有限，多阶段偏好对齐面临灾难性遗忘问题——即第二阶段中结构偏好对齐的效果会逐渐衰退，而难以三元组幽默对齐DPO方法。(triplet humor alignment DPO)
 - (y^+, y_h^-, y^*)
 - 修改幽默偏好三元组
 - 1. 模型参数: 模型参数θ, 用于生成输出y。
 - 2. 提示: 提示x, 用于生成输出y。
 - 3. 替代词: 替代词aw, 用于生成输出y。
 - 4. 损失函数: 损失函数ℓ, 用于衡量模型输出y与目标y+和y-之间的差异。
 - 5. 判别器: 判别器D, 用于判断输出y是否符合双关结构要求。
 - 损失函数
 - $\ell_{DPO} = -\log \pi_{\theta}(y^+ | x) + \log \pi_{\theta}(y^- | x)$ (7)
 - $\ell_{DPO} = -\log \pi_{\theta}(y^+ | x) + \log \pi_{\theta}(y^- | x)$ (8)
 - $\ell_{DPO} = -\log \pi_{\theta}(y^+ | x) + \log \pi_{\theta}(y^- | x)$ (9)
 - 参数解释
 - 1. 模型参数: 模型参数θ, 用于生成输出y。
 - 2. 提示: 提示x, 用于生成输出y。
 - 3. 替代词: 替代词aw, 用于生成输出y。
 - 4. 损失函数: 损失函数ℓ, 用于衡量模型输出y与目标y+和y-之间的差异。
 - 5. 判别器: 判别器D, 用于判断输出y是否符合双关结构要求。

数据集 (187,315 words in total)

- phonic: 语音双关
- graphic: 同形异义双关
- 我们发现大多数双关语句子只包含一个双关语对。因此，我们的研究将重点完全集中在这种情形。

指标

- 自动评估体系
 - 1. 结构偏好 (Structure Score)
 - 定义: 衡量生成文本中双关词与替代词的对齐程度 (yes/no/2018)
 - 计算: $Structure\ Score = \frac{\sum_{i=1}^N \mathbb{I}(y_i^+ = aw_i)}{N}$
 - 范围: 0.0000 - 1.0000 (0.5000为期望值)
 - 2. 多样性 (Diversity)
 - 定义: 衡量生成文本中双关词与替代词的多样性
 - 计算: $Diversity = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i^+ \neq aw_i)$
 - 范围: 0.0000 - 1.0000 (0.5000为期望值)
 - 3. 平均句长 (Avg Sentence Length)
 - 定义: 衡量生成文本中句子的平均长度
 - 计算: $Avg\ Sentence\ Length = \frac{1}{N} \sum_{i=1}^N |y_i^+|$
 - 范围: 0.0000 - 1.0000 (0.5000为期望值)
- 人工评估
 - 评估数据准备
 - 感知这部分通过机器结果参数下就行了

基准模型

- AmbiPun
 - LLMs LLaMA2-7B and Baichuan2-7B
 - 源自huggingFace Transformers库
- ChatGPT
 - gpt-3.5-turbo版本

实验设置

- 优化器
 - 模型: AdamW (Loshchilov & Hutter, 2017)
 - 学习率: 5e-5
 - 批量大小: 4
 - 训练轮数: 10k (1k/5k, 2017), 验证轮数: 1k/5k
- DPO参数设置
 - 温度系数: 0.5
 - 折扣系数: 0.95
 - 温度 (temperature): 0.05
 - top-p: 0.95
 - top-k: 5

实验结果

- 偏好数据规模
 - 偏好三元组样本量≥10,000组
 - 采样频率: 每步1000
- 幽默偏好数据
 - 1. 数据源: 从公开数据集 (如GPT-4) 中提取的幽默数据
 - 2. 数据清洗: 去除重复、低质量数据
 - 3. 数据标注: 标注幽默类型 (如双关、谐音)
 - 4. 数据验证: 验证数据质量和多样性
- 实验环境
 - 1. 硬件配置: 服务器 (CPU: Intel Xeon, GPU: NVIDIA A100)
 - 2. 软件配置: Python 3.8, PyTorch 1.10, Transformers 4.15
 - 3. 训练环境: 分布式训练 (多GPU)
- 性能突破
 - 多样性 (Dv) ChatGPT显著优于其他模型
 - 模型的dpo版本性能要优于sft指令微调版本
- 实验结果
 - 与基准模型的DPO方法相比，本方案最重要的改进在于采用多阶段课程学习策略——通过“由易到难”的分阶段方式独立优化两个偏好目标 (结构→幽默)，该设计的有效性得到了实验数据的充分验证。

消融实验

- 评估PGCL各组件部分对其的影响——实验与其各变体比较
 - w/o Improved Humor DPO——使用标准的DPO算法而非改进后的
 - w/o Humor——移除幽默对齐阶段 (第二阶段)，仅学习结构偏好
 - 与第二阶段使用标准DPO算法微调的LM相比，PGCL的结构成功率分别高出11.02%和8.56%。这表明，标准DPO算法在进行幽默偏好对齐时，难以保持对结构偏好的良好维持。
- 结果
 - 新的三元组对齐损失能有效指导大语言模型 (LLMs) 在结构偏好与幽默偏好之间保持平衡。

相关前期探索

- DPOs instruction tuning
 - 结论: 相较于指令微调，DPO方法能显著提升大语言模型的性能。
- 多阶段课程学习LM能专注攻克困难任务 (如幽默生成) 而不牺牲整体性能
 - 而DPO-A因双关语生成的固有难度，模型难以直接同步优化，性能提升有限。
- 不同策略在提升模型多目标偏好对齐能力上的性能表现 (基于中文双关数据集和SemEval数据集)