

Generative Verifiers: Reward Modeling as Next-Token Prediction

Abstract

验证器或奖励模型通常用于增强大型语言模型（LLMs）的推理性能

一种常见的方法是 Best-of-N 方法 其中由验证模型对 LLM 生成的 N 个候选解决方案进行排序，然后选择最好的一个

虽然基于 LLM 的验证器通常被训练为判别式验证器来对解决方案进行评分 它们没有利用预训练的 LLMs 的文本生成功能

联合验证和答案生成

为了克服这一限制 本文使用常用的下一个 token 预测来训练验证器

它们与指令微调无缝集成

与传统的验证器相比，此类生成验证器（GenRM）可以受益于 LLMs 的多个优势

支持思维链推理

可以通过多数投票来利用额外的推理时间以实现更好的评分

Best-of-N 求解算法和小学数学推理任务问题时 GenRM 作为验证器性能比判别式验证器性能提高 16–64%

效果

GenRM 在数据集大小、模型容量和推理时间计算方面具有良好的扩展性

1 Introduction

虽然大型语言模型（LLMs）表现出卓越的能力，但它们经常自信地犯逻辑和事实错误

解决这个问题的常见策略是 Best-of-N

LLM 为给定问题生成 N 个候选解决方案

以及一个学习奖励模型，称为“verifier”

对这些解决方案进行排名并选择最合适的一个

策略的有效性取决于验证器的准确性

基于 LLM 的验证模型通常被训练为判别式奖励模型（RM）

传统推理领域

为候选解决方案分配分数值

用于将它们分类为正确或错误

问题

这种评分方法并没有利用 LLMs 根本设计目的文本生成功能

没有利用生成式 LLMs 的固有优势

instruction tuning

思维链推理

利用额外的推理阶段的计算来提高性能

LLM-as-a-Judge 方法

通过提示现有的生成式 LLMs

也提供了上述优势

它通常不如经过基于 LLMs 训练的验证模型

称之为 GenRM

建议使用下一个 token 预测来训练验证模型

利用 LLMs 的文本生成功能

图 2：在 GSM8K 测试中使用生成验证器的示例

在这里，解决方案是不正确的，因为它忽略了问题中的“each”一词

判别式 RM 无法识别解决方案中的这种微妙错误

GenRM -CoT 验证器可以可靠地检测到该错误

因为 GenRM -CoT 接受了基于合成思维链推理的下一个标记预测的训练，使其能够明确地推理解决方案

为了生成解决方案的分数值

验证模型现在使用提示词 “Is the answer correct?”

并将分数值表示为单个文本 token 预测的概率（例如“是”或“否”）

假设在训练期间的数据有思维过程可用

它可以在使用“是”或“否”标记预测正确性之前生成思维过程来训练它的推理性能

GenRM 天然支持 CoT 推理

图 3：生成验证器的图示

给定一个问题和一个候选答案

GenRM 通过 SFT 直接微调 LLM 来回答问题“答案是否正确（是/否）？”

推理过程中，通过提取“yes”标记的概率来获得验证模型分数

GenRM-CoT 微调 LLM，以在产生最终的 Yes/No token 之前生成验证链（CoT）基本原理

推理时，我们对多个 CoT 思维过程进行抽样，并使用多数投票来计算“是”的平均概率，使 GenRM-CoT 能够利用额外的推理计算来更好地验证

图 1：在 Best-of-N 性能方面，将 GenRM 与标准验证方法在多个推理任务上进行比较

效果

通过将奖励模型重新设计为下一个 token 预测

GenRM 在提高 Best-of-N 性能方面优于其他方法

水平虚线对应于直接生成性能（Best-of-1）

GenRM 利用 LLMs 的生成功能，使经过微调的验证器模型能够在推理时有效地利用思维链和多数投票来识别推理错误

2 Preliminaries

Next-token 预测

预训练和微调 LLMs 的典型方法

监督微调（SFT）目的是减少模型预测的下一个 token 与给定目标 token 之间的交叉熵损失

给定一个数据集 $D = \{(x, y)\}$

输入上下文的 x 和目标响应 y

$$\mathcal{L}_{\text{SFT}}(\theta, D) = -\mathbb{E}_{(x, y) \sim D} \left[\sum_{t=1}^{|y|} \log p_{\theta}(y_t | x, y_{<t}) \right]$$

判别式验证器

正确答案和错误答案的数据集上微调 LLM 作为分类器

验证模型直接分配一个 0-1 的分数估计解决方案 y 对于某个问题 x 来说是正确的概率

判别式验证模型不利用 LLMs 的文本生成功能

3 GenRM

概述

基于 LLM 的判别式验证器不利用预训练的 LLMs 的文本生成功能

建议使用标准的 next-token 预测来训练验证器同时这个验证器可以生成文本

GenRM 使用 LLM 在各个 token 上的概率分布来表示答案的正确性，而不是单独预测一个数值

这使 GenRM 的生成能力保持完整

因为验证决策一个 token

同时还实现了 LLMs “免费”的多项优势

答案生成和验证

思维链推理

3.1 直接验证器（GenRM）

最简单的形式中

数据集 $D_{\text{Direct}} = \{(x, y), I, \text{'Yes'}\} \cup \{(x, y'), I, \text{'No'}\}$ ， $I = \text{'Is the answer correct (Yes/No)?'}$

推理时，我们使用“是”标记的概率

作为验证器对答案重新排名的分数 $r_{\text{Direct}}(x, y) = p_{\theta}(\text{Yes} | x, y, I)$

3.2 统一生成和验证

GenRM 将验证答案正确或错误的建模与生成正确答案的有监督微调 SFT 无缝集成

可以通过简单地更改包括验证和生成任务的数据混合比例以更改 SFT 损失

给定一个验证数据集 $D_{\text{Verify}} = \begin{matrix} D_{\text{Direct}} \\ D_{\text{CoT}} \end{matrix}$

生成数据集 带有正确性标记的问题-答案对（可选地带有 CoT 原理） D_{Correct}

$$\mathcal{L}_{\text{GenRM}}(\theta, D_{\text{Verify}}) = \mathcal{L}_{\text{SFT}}(\theta, D_{\text{Verify}}) + \lambda \mathcal{L}_{\text{SFT}}(\theta, D_{\text{Correct}})$$

其中 $\lambda > 0$ 是一个超参数 控制验证与生成任务的比例

统一的训练可以提高验证器和生成器的性能 两任务具有一定的相关性

具体地说

由于验证通常涉及复杂的推理 生成式验证器 可以从 CoT 推理中受益

对解决方案的正确性做出决定之前生成中间推理步骤

这可能会识别直接验证器遗漏的细微推理错误 图 3 底部

数据集 D_{CoT} 上最小化 SFT 损失 $\mathcal{L}_{\text{GenRM}}$

训练 CoT 验证器

problem-solution 对

输入

提示词 $I_{\text{CoT}} = \text{'Let's verify step by step.'}$

理由（思维过程） V_{CoT} 可以是人类生成的，也可以是 LLM 生成的 探索了这两种情况

targets

最终的问题 $I = \text{'Is the answer correct (Yes/No)?'}$

最终答案 'Yes' or 'No' token

$D_{\text{CoT}} = \{(x, y), I_{\text{CoT}}, (V_{\text{CoT}}, I, \text{'Yes'})\} \cup \{(x, y'), I_{\text{CoT}}, (V_{\text{CoT}}, I, \text{'No'})\}$ $I_{\text{CoT}} = \text{'Let's verify step by step.'}$

3.3 思维链验证器（GenRM-CoT）

从 GenRM -CoT 生成一个 CoT 基本原理 然后使用“yes”的概率来分配正确性分数

使用“是”的概率来分配正确性分数： $r_{\text{CoT}}(x, y) = p_{\theta}(\text{Yes} | x, y, I_{\text{CoT}}, V_{\text{CoT}}, I)$ ，其中 $V_{\text{CoT}} \sim p_{\theta}(\cdot | x, y, I_{\text{CoT}})$

与直接使用相比 上述 CoT 奖励还取决于

$I_{\text{CoT}} = \text{'Let's verify step by step.'}$

自生成的 V_{CoT}

生成式验证器可能会使用不同的推理路径

并对同一 problem-solution 对产生不同的正确概率

我们希望边缘化中间推理路径，以选择最一致的正确性答案

并将这些理由的 CoT 验证器分数平均化

使用前 k 个理由进行多数投票 $r_{\text{Majority}}(x, y) = \frac{1}{K} \sum_{t=1}^K p_{\theta}(\text{Yes} | x, y, I_{\text{CoT}}, V_{\text{CoT}}^{(t)}, I)$ ，其中 $V_{\text{CoT}}^{(t)} \sim p_{\theta}(\cdot | x, y, I_{\text{CoT}})$

CoT 验证器的单个推理理由可能存在推理错误 进行平均来减轻此类错误的影响

GenRM-CoT 可以利用额外的推理时间来提高其准确性 这是判别式验证器无法做到的

默认根据 32 票的多数投票报告 GenRM-CoT 性能

CoT Verifier 的综合验证理由生成

探索了在 GSM8K 上使用模型生成的原理

如果简单地使用给定问题-解决方案对和 'Let's verify step by step' 提示

并根据它们是否正确答案的正确性来过滤生成的理由 准确率为 50%

生成的理由质量很差

除了要验证的问题和答案之外，我们还提供了一个参考提示词 使 LLM 指出所提供的答案中的任何推理错误

4. Experiments

Problem: Carol spends 4 hours writing a song, half that much time recording it, and 90 minutes editing it. What percentage of her total work time did she spend editing?

Solution: Half of 4 hours is 2 hours. So the time spent recording is 2 hours. Add the editing time of 90 minutes and the total time for the entire project is 4 + 2 + 1.5 = 7.5 hours. That means she spent (90 / 7.5) * 100 = 12% of the total time editing. The answer is 12.

Discriminative RM correctness score: 0.8414

GenRM-CoT: Let's verify step by step.

Step 3: That means she spent (90 / 7.5) * 100 = 12% of the total time editing. Expected Answer: Incorrect. The editing time should be in hours, so it should be (1.5 / 7.5) * 100 = 20%.

Verification: Is the answer correct (Yes/No)? No

GenRM-CoT (Majority Voting) score: 0.2498

图 4：GenRM-CoT 捕捉到判别式验证器无法捕捉到的细微错误的示例

候选解决方案没有将 90 分钟转换为 1.5 小时，然后除以 7.5

判别式验证器无法检测到这个错误

我们提出的 GenRM-CoT 模型能够使用分步生成验证来识别这个错误

参考资料