

大模型面经之llm在任务型对话系统中的评价指标

瓦力算法学研所 2024年05月26日 15:34 安徽

任务型对话系统 (Task-oriented Dialogue Systems, TODS) 已经广泛应用于客户服务、个人助理、智能客服等领域。这些系统的核心目标是帮助用户完成特定的任务，如预订酒店、查询信息或执行交易。然而，要确保这些系统的有效性和用户满意度，就需要一套科学的评价指标来衡量它们的性能。以下是一些常用的评价指标：

1. BLEU (Bilingual Evaluation Understudy)

BLEU是一种广泛用于机器翻译和文本生成任务的评价指标，它通过比较机器生成的文本和一组参考文本之间的重叠程度来评估生成文本的质量。在任务型对话系统中，BLEU可以用来评估系统生成的回复与人类生成的标准回复之间的相似度。BLEU分数越高，表示系统生成的回复越接近人类生成的回复。

2. 任务完成率 (Task Success Rate)

任务完成率是衡量任务型对话系统性能的关键指标之一。它指的是系统成功帮助用户完成指定任务的比例。例如，在预订酒店的任务中，如果系统能够引导用户完成预订流程并确认预订信息，则认为任务成功完成。

3. 对话轮次 (Dialogue Turns)

对话轮次是指完成一个任务所需的对话回合数。理想情况下，对话轮次越少，说明系统的效率越高。然而，对话轮次的减少并不总是意味着更好的用户体验，因为过少的轮次可能导致用户需求没有被充分理解。

4. 意图识别准确率 (Intent Recognition Accuracy)

在任务型对话系统中，意图识别是指系统理解用户输入的意图并将其映射到相应的任务或操作的能力。意图识别准确率是衡量系统在这一过程中表现的指标，高准确率意味着系统能够更好地理解用户的需求。

5. 实体识别准确率 (Entity Recognition Accuracy)

实体识别是指系统从用户输入中提取关键信息（如时间、地点、人名等）的能力。实体识别准确率是衡量系统在这一过程中表现的指标，这对于执行特定任务（如预订、查询等）至关重要。

6. 响应时间 (Response Time)

响应时间是指系统接收到用户输入后生成回复所需的时间。快速的响应时间可以提高用户体验，尤其是在需要即时反馈的场景中。

7. 人工评价

人工评价（如mos分等）是最直接的评价方式，它依赖于人类评估者对对话系统的输出进行主观评价。人工评价可以包括多个维度，如流畅性、相关性、准确性、满意度等。人工评价的优点是能够综合考虑多种因素，但缺点是成本较高，且可能受到评估者主观性的影响。

(1) . 评价维度

人工评价通常包括多个维度，这些维度可以是：

- **流畅性**：对话的自然程度，是否像人类之间的对话。
- **相关性**：系统回复是否与用户的意图和上下文紧密相关。
- **准确性**：系统提供的信息是否正确无误。

- **完整性**：系统是否提供了用户所需的所有信息。
- **一致性**：系统在对话过程中是否保持一致性，不出现矛盾。
- **满意度**：用户对对话体验的整体满意程度。

(2) . 评价方法

人工评价可以采取不同的方法进行：

- **直接评价**：评估者直接与对话系统进行交互，然后根据体验给出评价。
- **旁观评价**：评估者观察用户与系统的交互过程，然后给出评价。
- **回放评价**：评估者观看或听取用户与系统的交互记录，然后进行评价。

(3) . 评价标准

人工评价通常需要一套明确的标准或指南，以确保评价的一致性和可重复性。这些标准可能包括：

- **评分量表**：如1到5或1到10的评分系统。
- **检查清单**：列出需要评估的具体项目。
- **评价指南**：提供评价过程中的指导和建议。