

Забывание, вызванное отрицанием, в больших языковых моделях

Дата: 2025-02-26 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2502.19211>

Рейтинг: 65

Адаптивность: 70

Ключевые выводы:

Исследование изучает, проявляют ли большие языковые модели (LLM) эффект забывания, вызванного отрицанием (NIF) - когнитивное явление, наблюдаемое у людей, когда отрицание неверных атрибутов объекта или события приводит к худшему запоминанию информации по сравнению с утверждением правильных атрибутов. Результаты показали, что ChatGPT-3.5 демонстрирует значимый эффект NIF, GPT-4o-mini показывает маргинально значимый эффект, а LLaMA-3-70B не проявляет данного эффекта.

Объяснение метода:

Исследование демонстрирует, что некоторые LLM (особенно ChatGPT-3.5) хуже запоминают информацию, представленную в отрицательной форме. Это позволяет пользователям оптимизировать взаимодействие с LLM, предпочитая утвердительные формулировки для лучшего сохранения информации. Результаты различаются между моделями, что важно учитывать при выборе LLM для конкретных задач.

Ключевые аспекты исследования: 1. **Феномен негативно-индуцированного забывания (NIF):** Исследование изучает, проявляют ли языковые модели (LLM) эффект негативно-индуцированного забывания, который наблюдается у людей - когда отрицание неверной информации приводит к худшему запоминанию, чем подтверждение верной информации.

Методология тестирования: Авторы адаптировали экспериментальную структуру Zang et al. (2023) для тестирования ChatGPT-3.5, GPT-4o-mini и LLaMA-3-70B, используя задачи на верификацию и свободное воспроизведение информации из рассказа.

Результаты по моделям: ChatGPT-3.5 продемонстрировал значимый эффект NIF, GPT-4o-mini показал маргинально значимый эффект, а LLaMA-3-70B не проявил данного эффекта, демонстрируя очень высокую точность воспроизведения.

Сравнение с человеческими когнитивными смещениями: Исследование

расширяет понимание того, как LLM могут воспроизводить когнитивные смещения, характерные для людей, без явного программирования таких эффектов.

Дополнение: Для работы методов данного исследования не требуется дообучение или специальный API. Исследователи использовали стандартный интерфейс чата с моделями (ChatGPT-3.5, GPT-4o-mini, LLaMA-3-70B), и основные концепции можно применить в обычном диалоге с LLM.

Концепции и подходы, применимые в стандартном чате:

Предпочтение утвердительных формулировок: Вместо "не делай X" можно использовать "делай Y". Например, вместо "не используй сложные термины" лучше сказать "используй простые, понятные слова".

Повторение ключевой информации в утвердительной форме: Если необходимо использовать отрицание, можно дополнительно повторить ту же информацию в утвердительной форме для лучшего запоминания.

Учет различий между моделями: Более новые модели (GPT-4, LLaMA-3) могут лучше справляться с запоминанием информации в контексте отрицаний, что можно учитывать при выборе модели.

Проверка усвоения информации: После предоставления инструкций с отрицаниями можно попросить модель повторить ключевые моменты для проверки их запоминания.

Применяя эти концепции, пользователи могут ожидать следующие результаты: - Более точное следование инструкциям - Снижение вероятности "забывания" важной информации - Улучшение последовательности и связности в длительных диалогах - Более эффективное управление контекстом взаимодействия с LLM

Это исследование особенно ценно тем, что выявляет когнитивное ограничение, которое может влиять на повседневное взаимодействие с LLM, и предлагает простой способ его преодоления через адаптацию формулировок.

Анализ практической применимости: 1. **Феномен негативно-индуцированного забывания (NIF):** - Прямая применимость: Средняя. Пользователи могут учитывать, что информация, представленная в отрицательной форме, может хуже запоминаться моделью, и формулировать запросы соответствующим образом. - Концептуальная ценность: Высокая. Помогает понять, что LLM могут иметь "слабости памяти" при работе с отрицаниями, что важно при формулировании инструкций. - Потенциал для адаптации: Значительный. Пользователи могут переформулировать отрицательные утверждения в положительные для повышения вероятности их сохранения моделью.

Методология тестирования: Прямая применимость: Низкая. Экспериментальная структура скорее представляет академический интерес. Концептуальная ценность: Средняя. Демонстрирует подход к тестированию когнитивных аспектов LLM.

Потенциал для адаптации: Низкий. Обычным пользователям трудно воспроизвести такие эксперименты.

Результаты по моделям:

Прямая применимость: Средняя. Пользователи могут учитывать разницу в обработке отрицаний между моделями при выборе LLM для конкретных задач. Концептуальная ценность: Высокая. Показывает, что более новые модели (GPT-4o-mini, LLaMA-3) лучше справляются с запоминанием отрицаемой информации. Потенциал для адаптации: Средний. Можно выбирать модель в зависимости от специфики задачи.

Сравнение с человеческими когнитивными смещениями:

Прямая применимость: Средняя. Понимание, что LLM имеют "человеческие" когнитивные смещения, помогает предсказывать их поведение. Концептуальная ценность: Высокая. Дает фундаментальное понимание ограничений LLM. Потенциал для адаптации: Значительный. Пользователи могут адаптировать свои стратегии взаимодействия с LLM, учитывая эти смещения.

Prompt:

Использование знаний о забывании, вызванном отрицанием (NIF) в промптах для GPT ## Ключевое понимание эффекта NIF Исследование показало, что некоторые языковые модели (особенно ChatGPT-3.5 и в меньшей степени GPT-4o-mini) демонстрируют эффект забывания, вызванного отрицанием - они хуже запоминают информацию, представленную в форме отрицания, чем в утвердительной форме.

Пример промпта с учетом эффекта NIF

Неоптимальный промпт: [=====] Создай инструкцию по безопасности для химической лаборатории. Обязательно укажи, что нельзя смешивать хлор и аммиак, не следует хранить легковоспламеняющиеся вещества рядом с источниками тепла, и не забудь упомянуть, что нельзя есть в лаборатории. [=====]

Оптимизированный промпт: [=====] Создай инструкцию по безопасности для химической лаборатории. Обязательно укажи следующие правила: 1. Храни хлор и аммиак отдельно друг от друга 2. Размещай легковоспламеняющиеся вещества вдали от источников тепла 3. Принимай пищу только в специально отведенных местах вне лаборатории

Для каждого правила добавь краткое объяснение, почему это важно, и представь информацию в утвердительной форме для лучшего запоминания. [=====]

Объяснение применения знаний из исследования

Замена отрицаний на утверждения: Вместо "нельзя смешивать X и Y" → "храни X и Y отдельно"

Позитивное переформулирование: Вместо "не ешь в лаборатории" → "принимай пищу в отведенных местах"

Структурирование информации: Пронумерованный список делает утверждения более заметными и легче запоминаемыми

Запрос на утвердительные формулировки: Явное указание модели представлять информацию в утвердительной форме

Запрос объяснений: Просьба объяснить причины правил усиливает связи между концепциями в "памяти" модели

Такой подход особенно важен при работе с ChatGPT-3.5, где эффект NIF наиболее выражен, и может быть полезен для взаимодействия с другими моделями для повышения точности запоминания критически важной информации.