

Полагаться или не полагаться? Оценка вмешательств для адекватного использования больших языковых моделей

Дата: 2025-03-08 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2412.15584>

Рейтинг: 70

Адаптивность: 75

Ключевые выводы:

Основная цель исследования - оценить эффективность различных вмешательств (интервенций) для достижения надлежащего уровня доверия к LLM. Главные результаты показывают, что хотя интервенции могут значительно снизить чрезмерное доверие к LLM, они обычно не улучшают надлежащее доверие в целом. Интервенции, как правило, снижают общий уровень доверия, уменьшая чрезмерное доверие за счет полезного доверия.

Объяснение метода:

Исследование предлагает практически применимые стратегии для улучшения взаимодействия с LLM, особенно "предупреждения о надежности" и "неявные ответы". Предоставляет важную концептуальную основу для понимания баланса между чрезмерным и недостаточным доверием. Некоторые методы технически сложны для реализации обычными пользователями, а результаты показывают неоднозначную эффективность в разных контекстах задач.

Ключевые аспекты исследования: 1. **Сравнение интервенций для улучшения надлежащего доверия к LLM:** Исследование оценивает три типа интервенций, влияющих на то, насколько пользователи доверяют советам LLM: предупреждения о надежности (reliance disclaimer), выделение неопределенностей (uncertainty highlighting) и неявные ответы (implicit answer).

Экспериментальная методология: Проведен онлайн-эксперимент с 400 участниками, выполнявшими два типа задач: логические рассуждения и количественную оценку. Участники сначала отвечали самостоятельно, затем получали совет от LLM (с одной из трех интервенций или в контрольной группе) и отвечали повторно.

Метрики для оценки доверия: Разработаны метрики для оценки уровня доверия пользователей к LLM, включая "относительное доверие к LLM" (RLR),

"относительное доверие к себе" (RSR) и "коэффициент надлежащего доверия" (ARR).

Калибровка уверенности: Исследуется, как различные интервенции влияют на уверенность пользователей в своих ответах и насколько эта уверенность соответствует фактической точности.

Компромисс между чрезмерным и недостаточным доверием: Выявлено, что интервенции, уменьшающие чрезмерное доверие к LLM, часто увеличивают недостаточное доверие, что указывает на сложность разработки сбалансированных решений.

Дополнение:

Применимость методов в стандартном чате

Не все методы из исследования требуют дообучения или API. Вот что можно применить в стандартном чате:

Предупреждения о надежности (reliance disclaimer) - полностью применимы без каких-либо технических модификаций. Пользователь может: Добавить в промпт просьбу, чтобы модель заканчивала ответы фразой о необходимости проверки информации. Самостоятельно критически оценивать ответы LLM, помня о возможных ошибках.

Неявные ответы (implicit answer) - также полностью применимы в стандартном чате:

Пользователь может запрашивать "объясни рассуждения, но не давай прямого ответа". Можно просить модель показать ход решения без указания финального ответа. Эффективно для образовательных целей и сложных решений.

Концепция "выделения неопределенностей" - хотя технически нельзя увидеть подсветку токенов, можно:

Просить модель указать, в каких частях ответа она менее уверена. Запрашивать оценку уверенности для разных частей объяснения. Просить модель отметить предположения и факты отдельно. ### Ожидаемые результаты от применения концепций:

Улучшение надлежащего доверия - более осознанное взаимодействие с LLM, где пользователь опирается на модель, когда это полезно, и на себя, когда совет модели сомнителен.

Повышение когнитивной вовлеченности - особенно метод "неявных ответов" заставляет пользователя активно участвовать в процессе решения.

Лучшая калибровка уверенности - понимание, что уверенность в ответах должна

соответствовать их фактической точности.

Снижение чрезмерного доверия - особенно важно для критических задач, где ошибки могут иметь серьезные последствия.

Важно отметить, что исследование показывает компромисс между чрезмерным и недостаточным доверием - снижение одного часто приводит к увеличению другого. Поэтому выбор метода должен зависеть от контекста задачи и потенциальных рисков ошибок.

Анализ практической применимости: 1. **Предупреждения о надежности (reliance disclaimer)** - **Прямая применимость**: Высокая. Простое добавление предупреждения "Не забудьте проверить эту информацию" к ответам LLM может быть легко реализовано любым пользователем в обычном чате. Это наиболее эффективная интервенция, улучшающая надлежащее доверие без значительного увеличения времени использования. - **Концептуальная ценность**: Средняя. Помогает пользователям понять, что ответы LLM не всегда точны и требуют критической оценки. - **Потенциал для адаптации**: Высокий. Можно легко внедрить в любые взаимодействия с LLM, включая персональные запросы и рабочие задачи.

Выделение неопределенностей (uncertainty highlighting) **Прямая применимость**: Низкая для обычных пользователей, так как требует доступа к вероятностям токенов LLM. **Концептуальная ценность**: Средняя. Показывает, что LLM может быть менее уверенной в определенных частях своего ответа. **Потенциал для адаптации**: Средний. Хотя технические детали сложны для реализации, концепция может быть адаптирована в виде просьбы к LLM указать уровень уверенности в разных частях ответа.

Неявные ответы (implicit answer)

Прямая применимость: Высокая. Пользователи могут просить LLM объяснить рассуждения без предоставления прямого ответа. **Концептуальная ценность**: Высокая. Заставляет пользователей активно участвовать в процессе решения, а не просто принимать готовый ответ. **Потенциал для адаптации**: Высокий. Применимо в образовательных, исследовательских и профессиональных контекстах, где важно понимание процесса.

Метрики оценки доверия

Прямая применимость: Средняя. Помогают осознать важность как чрезмерного, так и недостаточного доверия к LLM. **Концептуальная ценность**: Высокая. Предоставляют структуру для понимания взаимодействия с LLM. **Потенциал для адаптации**: Средний. Могут быть адаптированы для самооценки пользователем своего взаимодействия с LLM.

Выводы о калибровке уверенности

Прямая применимость: Высокая. Осознание того, что пользователи часто

становятся более уверенными после неправильных решений о доверии. **Концептуальная ценность:** Высокая. Помогает понять опасность необоснованной уверенности при использовании LLM. **Потенциал для адаптации:** Высокий. Может использоваться для разработки стратегий взаимодействия с LLM в критически важных областях.

Prompt:

Применение знаний из исследования в промптах для GPT ## Ключевые выводы исследования

Исследование показывает, что различные интервенции влияют на то, как пользователи доверяют ответам языковых моделей:

Простые дисклеймеры наиболее эффективны для улучшения надлежащего доверия **Сложные методы** (например, выделение неопределенности) могут вызвать негативное восприятие **Неявные ответы** улучшают самостоятельность, но требуют больше времени на обработку ## Пример промпта с применением выводов исследования

[=====] Ты - эксперт по финансовому планированию. Мне нужен совет по диверсификации инвестиционного портфеля на сумму \$50,000.

Важные правила для твоего ответа: 1. Начни с дисклеймера, что твои рекомендации не заменяют консультацию лицензированного финансового советника. 2. Вместо прямого списка конкретных активов, опиши принципы и стратегии распределения, чтобы я самостоятельно принял(а) окончательное решение. 3. Если в каком-то аспекте есть значительная неопределенность, просто укажи на это простым языком без сложных технических деталей. 4. Задай мне 2-3 вопроса в конце, которые помогут мне критически оценить твои рекомендации. [=====]

Объяснение применения исследования

Дисклеймер в начале промпта следует выводу исследования о том, что простые предупреждения эффективно калибруют доверие пользователя.

Просьба о принципах вместо конкретных решений реализует концепцию "неявного ответа", стимулируя самостоятельное мышление пользователя.

Указание говорить о неопределенности простым языком учитывает вывод о том, что сложные методы выделения неопределенности могут ухудшить восприятие.

Запрос на вопросы для самопроверки помогает пользователю критически оценить информацию, что способствует надлежащему уровню доверия.

Такой подход к составлению промптов помогает достичь баланса между полезностью информации от LLM и предотвращением чрезмерного доверия к ней.