

Обратите внимание на разрыв уверенности: избыточная уверенность, калибровка и эффекты отвлекающих факторов в больших языковых моделях

Дата: 2025-02-16 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2502.11028>

Рейтинг: 72

Адаптивность: 85

Ключевые выводы:

Исследование направлено на анализ проблемы калибровки уверенности в больших языковых моделях (LLM). Основные результаты показывают, что хотя более крупные модели (например, GPT-4o) в целом лучше калиброваны, они более подвержены отвлечению на неверные варианты ответов, в то время как меньшие модели больше выигрывают от предоставления вариантов ответов, но хуже справляются с оценкой неопределенности.

Объяснение метода:

Исследование выявляет критическую проблему избыточной уверенности LLM и предоставляет практические стратегии улучшения взаимодействия. Показывает, как формулировать запросы с вариантами ответов для повышения точности, особенно для меньших моделей. Объясняет различия в поведении моделей разного размера и влияние типов вопросов на калибровку.

Ключевые аспекты исследования: 1. **Проблема избыточной уверенности в LLM:** Исследование фокусируется на проблеме калибровки уверенности в больших языковых моделях, которые часто демонстрируют избыточную уверенность в неправильных ответах, что может вводить пользователей в заблуждение.

Влияние размера модели и дистракторов: Авторы изучают, как размер модели и наличие вариантов ответов (дистракторов) влияют на точность и калибровку уверенности LLM.

Сравнительный анализ моделей: Исследование анализирует различные модели (GPT-4o, GPT-4-turbo, GPT-4o-mini, Llama 3, Gemma2) и их поведение в задачах с открытыми вопросами и вопросами с множественным выбором.

Типы вопросов и калибровка: Авторы исследуют, как разные типы вопросов (о

датах, числах, людях, местах) влияют на точность и калибровку моделей.

Метрики для оценки калибровки: В исследовании используется метрика ожидаемой ошибки калибровки (ECE) для измерения расхождения между уверенностью модели и фактической точностью.

Дополнение:

Применимость методов в стандартном чате

Исследование не требует дообучения моделей или специального API для применения его основных выводов. Большинство методов и концепций можно непосредственно использовать в стандартном чате с LLM:

Предоставление вариантов ответов: Пользователи могут формулировать запросы в формате множественного выбора, предлагая модели несколько возможных вариантов ответа. Исследование показывает, что это значительно повышает точность, особенно для меньших моделей.

Критическая оценка уверенности: Понимание того, что высокая уверенность модели не гарантирует правильность ответа, позволяет пользователям более критически оценивать ответы и проверять информацию из других источников.

Адаптация по типам вопросов: Пользователи могут быть более осторожными с вопросами о людях и датах, где модели показывают большую избыточную уверенность, и более доверять ответам на вопросы о местах.

Стратегия проверки: Можно задавать один и тот же вопрос в разных форматах (с вариантами ответов и без) и сравнивать результаты для повышения уверенности в правильности.

Применение этих концепций должно привести к: - Повышению точности получаемых ответов - Более реалистичным ожиданиям от взаимодействия с LLM - Уменьшению риска принятия неправильной информации из-за избыточной уверенности модели - Более эффективным стратегиям формулирования запросов

Анализ практической применимости: 1. **Проблема избыточной уверенности в LLM:** - Прямая применимость: Высокая. Пользователи должны знать, что высокая уверенность LLM не всегда означает правильность ответа, и критически оценивать ответы с высокой заявленной уверенностью. - Концептуальная ценность: Очень высокая. Понимание того, что LLM могут быть избыточно уверены в неправильных ответах, критически важно для формирования правильных ожиданий от взаимодействия с моделями. - Потенциал для адаптации: Высокий. Пользователи могут разработать стратегии проверки ответов с высокой заявленной уверенностью, особенно для фактологических вопросов.

Влияние размера модели и дистракторов: Прямая применимость: Средняя. Пользователи могут улучшить точность ответов, предоставляя модели варианты

ответов, особенно для меньших моделей. Концептуальная ценность: Высокая. Понимание того, что разные модели по-разному реагируют на наличие вариантов ответов, помогает выбирать оптимальные стратегии взаимодействия. Потенциал для адаптации: Высокий. Пользователи могут адаптировать свои запросы, предлагая модели несколько вариантов для выбора, особенно при работе с меньшими моделями.

Сравнительный анализ моделей:

Прямая применимость: Средняя. Пользователи могут выбирать модели в зависимости от задачи, учитывая, что большие модели (GPT-4o) лучше калиброваны, но более подвержены дистракции. Концептуальная ценность: Высокая. Понимание различий между моделями разного размера помогает формировать реалистичные ожидания. Потенциал для адаптации: Средний. Обычные пользователи часто ограничены в выборе модели, но могут адаптировать стратегии взаимодействия.

Типы вопросов и калибровка:

Прямая применимость: Высокая. Пользователи могут учитывать, что LLM лучше справляются с определенными типами вопросов (например, о местах), и быть более осторожными с другими (например, о людях). Концептуальная ценность: Высокая. Понимание сильных и слабых сторон моделей в разных типах вопросов помогает формулировать запросы и оценивать ответы. Потенциал для адаптации: Высокий. Пользователи могут адаптировать свои запросы в зависимости от типа информации, которую они ищут.

Метрики для оценки калибровки:

Прямая применимость: Низкая. Обычные пользователи не могут непосредственно измерить ECE. Концептуальная ценность: Средняя. Понимание того, что существуют метрики для оценки калибровки, помогает критически оценивать заявления о способностях LLM. Потенциал для адаптации: Низкий. Технические метрики сложно адаптировать для повседневного использования.

Prompt:

Применение исследования о калибровке уверенности в промптах для GPT ##
Ключевые выводы исследования для использования в промптах

Исследование показывает, что: - Более крупные модели лучше калиброваны, но подвержены отвлечению на неверные варианты - Предоставление структурированных вариантов ответов значительно улучшает точность и калибровку - Разные типы вопросов требуют разных подходов к калибровке уверенности - Даже хорошо калиброванные модели могут проявлять чрезмерную уверенность

Пример промпта с применением этих знаний

[=====] Я задам вопрос, требующий фактической информации. Пожалуйста:

Сначала сформулируй несколько возможных ответов на этот вопрос (минимум 3-4 варианта) Для каждого варианта приведи краткое обоснование, почему он может быть верным Оцени свою уверенность в каждом варианте по шкале от 0 до 100% Если твоя уверенность превышает 70%, дополнительно объясни, на чем основана такая высокая уверенность Выбери окончательный ответ, но если ты не уверен(а) более чем на 60%, явно укажи это Если вопрос касается конкретного человека, уточни о какой именно личности идет речь, чтобы избежать неоднозначности Вопрос: Кто написал роман "Война и мир"? [=====]

Почему этот промпт работает на основе исследования

Структурированные варианты ответов: Промпт требует генерации нескольких вариантов, что согласно исследованию повышает точность (с 35.14% до 73.42% для GPT-4o).

Явная калибровка уверенности: Запрос на оценку уверенности по шкале заставляет модель лучше калибровать свои ответы.

Дополнительное обоснование высокой уверенности: Исследование показало, что модели часто проявляют избыточную уверенность в диапазоне 70-100%, поэтому промпт требует дополнительного обоснования.

Дезамбигуация для вопросов о людях: Исследование выявило, что вопросы о людях наиболее сложны из-за неоднозначности имен, поэтому промпт включает специальный пункт для уточнения личности.

Пороговый уровень уверенности: Установка порога в 60% для явного признания неуверенности помогает избежать избыточной уверенности в пограничных случаях.

Такой подход существенно улучшает калибровку уверенности модели и повышает точность ответов, особенно в задачах, требующих фактической информации.