

Насколько надежны чат-боты как аннотаторы текста? Иногда

Дата: 2025-02-25 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2311.05769>

Рейтинг: 68

Адаптивность: 75

Ключевые выводы:

Основная цель исследования - систематически оценить эффективность открытых (open-source) языковых моделей (LLMs) по сравнению с ChatGPT и стандартными подходами к классификации с помощью машинного обучения для задач аннотации текста. Главные результаты показали, что производительность ChatGPT и открытых моделей значительно варьируется и часто непредсказуема, при этом супервизорный классификатор DistilBERT обычно превосходит обе группы моделей.

Объяснение метода:

Исследование предоставляет ценные практические знания о выборе моделей для аннотирования текста, демонстрируя, что традиционные методы с учителем часто превосходят LLM. Общая методология (zero/few-shot, типы промптов) применима широкой аудиторией, но полная реализация рекомендаций требует технических навыков для обучения моделей с учителем, что снижает доступность для некоторых пользователей.

Ключевые аспекты исследования: 1. **Сравнительная оценка моделей для аннотирования текста:** Исследование систематически сравнивает эффективность различных моделей для задач аннотирования текста: ChatGPT (закрытая модель), открытые LLM и классические модели машинного обучения с учителем (supervised).

Методология эксперимента: Авторы тестируют модели на двух бинарных задачах классификации твитов: определение политического контента и определение наличия примеров людей ("exemplars"). Применяются различные подходы: zero-shot, few-shot, с использованием как общих, так и специально разработанных промптов.

Результаты производительности: В большинстве случаев модель DistilBERT с учителем превосходит как ChatGPT, так и открытые LLM. GPT-4 показывает хорошие результаты только в некоторых задачах, а производительность моделей значительно варьируется в зависимости от задачи.

Проблемы открытой науки: Исследование подчеркивает проблемы использования закрытых моделей (непрозрачность, высокая стоимость, проблемы с защитой

данных) и оценивает, существует ли компромисс между точностью классификации и принципами открытой науки.

Практические рекомендации: Авторы рекомендуют осторожно подходить к использованию ChatGPT для аннотирования текста и предлагают вместо этого использовать размеченные людьми данные для обучения моделей с учителем.

Дополнение: Исследование не требует дообучения или API для применения основных методов и концепций. Хотя авторы использовали API для доступа к ChatGPT, многие подходы можно адаптировать для стандартного чата с LLM.

Основные концепции, которые можно применить в стандартном чате:

Few-shot обучение: Предоставление примеров в промпте значительно улучшает качество классификации. Пользователи могут включать 3-5 примеров текстов с правильными метками перед основным запросом.

Специализированные промпты: Исследование показывает, что специально разработанные промпты (с определениями категорий) обычно работают лучше, чем общие. Пользователи могут включать четкие определения категорий в свои запросы.

Понимание ограничений: Осознание того, что производительность LLM может значительно варьироваться в зависимости от задачи, помогает формировать реалистичные ожидания и проверять результаты.

Итеративное улучшение: Пользователи могут экспериментировать с различными формулировками промптов и количеством примеров для оптимизации результатов.

Применяя эти концепции в стандартном чате, пользователи могут получить более точные и надежные результаты классификации текста, хотя и не на уровне специализированных моделей с учителем.

Анализ практической применимости: 1. **Сравнительная оценка моделей для аннотирования текста** - Прямая применимость: Высокая. Исследование дает пользователям четкое представление о том, какие подходы лучше работают для аннотирования текста. Обычные пользователи могут использовать эти знания при выборе инструментов для решения своих задач. - Концептуальная ценность: Высокая. Исследование демонстрирует, что LLM не всегда являются лучшим решением для конкретных задач, и что традиционные методы могут быть более эффективными. - Потенциал для адаптации: Средний. Результаты могут быть адаптированы для других типов задач классификации текста.

Методология эксперимента Прямая применимость: Средняя. Пользователи могут использовать описанные промпты и подходы (zero-shot, few-shot) для своих задач аннотирования. Концептуальная ценность: Высокая. Демонстрирует важность правильного выбора промпта и предоставления примеров для повышения качества работы LLM. Потенциал для адаптации: Высокий. Методология может быть

адаптирована для различных задач классификации текста.

Результаты производительности

Прямая применимость: Высокая. Пользователи получают конкретные метрики эффективности различных моделей, что помогает принять обоснованное решение о выборе инструмента. Концептуальная ценность: Высокая. Результаты показывают, что производительность LLM может быть непредсказуемой и сильно зависит от конкретной задачи. Потенциал для адаптации: Средний. Результаты специфичны для конкретных задач, но общие выводы о вариативности производительности применимы широко.

Проблемы открытой науки

Прямая применимость: Средняя. Помогает пользователям осознать ограничения закрытых моделей с точки зрения воспроизводимости результатов. Концептуальная ценность: Высокая. Поднимает важные вопросы о прозрачности, воспроизводимости и стоимости использования закрытых моделей. Потенциал для адаптации: Низкий. Эти проблемы являются системными и не могут быть легко адаптированы пользователями.

Практические рекомендации

Прямая применимость: Высокая. Конкретные рекомендации по выбору подхода к аннотированию текста. Концептуальная ценность: Высокая. Помогает пользователям принимать более обоснованные решения о выборе инструментов. Потенциал для адаптации: Средний. Рекомендации могут быть адаптированы для различных задач классификации текста.

Prompt:

Использование знаний из исследования о чат-ботах как аннотаторах текста ##
Ключевые уроки исследования

Исследование показывает, что: - Эффективность LLM для аннотации текста сильно варьируется - Специализированные промпты работают лучше общих - Few-shot подход превосходит zero-shot - DistilBERT обычно превосходит генеративные модели

Пример улучшенного промпта

Вот пример промпта, который использует знания из исследования:

[=====] Я хочу, чтобы ты выполнил задачу классификации текста, определив, содержит ли следующий твит политический контент.

Вот несколько примеров для понимания задачи: 1. "Новая налоговая политика администрации вызвала споры в Конгрессе" - ПОЛИТИЧЕСКИЙ 2. "Сегодня

прекрасная погода для пикника в парке" - НЕПОЛИТИЧЕСКИЙ 3. "Президент подписал указ о защите окружающей среды" - ПОЛИТИЧЕСКИЙ

При анализе используй следующие критерии: - Упоминаются ли политические фигуры, партии или институты - Обсуждаются ли законы, политика или государственное управление - Содержится ли политическая риторика или идеология

Твит для анализа: [ТВИТ]

Дай ответ в формате "ПОЛИТИЧЕСКИЙ" или "НЕПОЛИТИЧЕСКИЙ", а затем кратко объясни свое решение. [=====]

Почему это работает лучше

Использует few-shot подход - включает примеры для обучения модели, что улучшает точность согласно исследованию

Применяет специализированный промпт - содержит конкретные критерии для задачи классификации, что дает лучшие результаты, чем общие инструкции

Структурирует ответ - запрашивает конкретный формат ответа, что снижает неоднозначность

Включает объяснение - просит модель объяснить свое решение, что позволяет оценить качество рассуждений

Дополнительные рекомендации

- Для критических задач лучше использовать супервизорные модели вроде DistilBERT
- Тестируйте разные версии промптов на небольшой выборке перед полномасштабным применением
- Для сложных задач классификации предоставляйте больше разнообразных примеров
- Учитывайте, что даже с оптимальным промptom результаты могут быть непредсказуемыми