

# Рисование панд: Бенчмарк для LLM в генерации кода для построения графиков

Дата: 2025-02-26 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2412.02764>

Рейтинг: 75

Адаптивность: 85

## Ключевые выводы:

Исследование представляет новый бенчмарк PandasPlotBench для оценки способности языковых моделей (LLM) генерировать код для визуализации табличных данных на основе естественно-языковых инструкций. Основные результаты показывают, что современные LLM хорошо справляются с генерацией кода для популярных библиотек визуализации (Matplotlib, Seaborn), но испытывают трудности с менее распространенными (Plotly), а также что сокращение длины задания минимально влияет на качество генерируемых визуализаций.

## Объяснение метода:

Исследование предоставляет конкретные рекомендации по использованию LLM для визуализации данных, применимые для широкой аудитории. Выводы о влиянии длины инструкций, эффективности моделей и библиотек имеют прямую практическую ценность. Ограничения включают фокус на Python и потенциальное устаревание сравнительных результатов с выходом новых версий моделей.

## Ключевые аспекты исследования: 1. **PandasPlotBench** - исследование представляет новый бенчмарк для оценки способности языковых моделей генерировать код визуализации данных на основе естественно-языковых инструкций. Содержит 175 уникальных задач с фокусом на визуализацию данных из Pandas DataFrame.

**Оценка различных LLM** - бенчмарк тестирует различные модели (GPT-4o, Claude 3, Gemini, Llama и др.) на способность генерировать код визуализации, обнаруживая существенные различия в их производительности.

**Влияние длины задачи** - исследование показывает, что сокращение длины инструкции минимально влияет на качество генерируемых визуализаций, что важно для пользовательского опыта.

**Сравнение библиотек визуализации** - бенчмарк оценивает эффективность моделей при работе с разными библиотеками (Matplotlib, Seaborn, Plotly), выявляя

значительные различия в знакомстве моделей с этими библиотеками.

**Методология оценки** - разработана двойная система оценки: на основе визуального сравнения с эталоном и на основе соответствия задаче, что обеспечивает комплексную оценку качества генерации.

## Дополнение: Исследование не требует дообучения или API для применения его методов в стандартном чате. Ученые использовали расширенные техники (такие как API различных моделей) для проведения бенчмарка, но ключевые концепции и подходы применимы в стандартном чате с LLM.

Основные применимые концепции:

**Структурированное описание данных** - исследование показывает, что включение `df.head(5)` вместе с информацией о типах столбцов критически важно для качественной генерации кода. Пользователи могут применять этот подход, предоставляя структурированное описание своих данных в чате.

**Краткость инструкций** - исследование демонстрирует, что даже одно предложение может быть достаточным для качественной визуализации при наличии хорошего описания данных. Пользователи могут формулировать краткие запросы, экономя время и токены.

**Выбор подходящей библиотеки** - понимание, что модели лучше справляются с популярными библиотеками (Matplotlib, Seaborn), позволяет пользователям делать оптимальный выбор инструментов.

**Разделение задачи и стилизации** - исследование показывает эффективность разделения инструкций на описание задачи и стилизации. Пользователи могут применять этот подход, сначала запрашивая базовую визуализацию, а затем отдельно уточняя стилизацию.

Применяя эти концепции в стандартном чате, пользователи могут ожидать следующих результатов: - Более точную генерацию кода для визуализации данных - Сокращение времени на формулировку запросов - Улучшение качества визуализаций - Минимизацию ошибок в сгенерированном коде

Эти подходы не требуют специальных API или дообучения моделей и могут быть непосредственно использованы в любом чате с современными LLM.

## Анализ практической применимости: **Аспект 1: PandasPlotBench - Прямая применимость:** Высокая для аналитиков данных и исследователей. Пользователи могут понять, какие модели лучше справляются с задачами визуализации, и использовать эти знания при выборе инструментов. - **Концептуальная ценность:** Показывает, что современные LLM способны эффективно генерировать код для визуализации данных, что меняет подход к анализу данных. - **Потенциал для адаптации:** Методология бенчмарка может быть адаптирована для тестирования других задач, связанных с генерацией кода.

**Аспект 2: Оценка различных LLM - Прямая применимость:** Пользователи получают четкое понимание, какие модели лучше подходят для задач визуализации данных (GPT-4o и Claude 3.5 Sonnet показывают наилучшие результаты). - **Концептуальная ценность:** Демонстрирует, что даже большие открытые модели (Llama 70B и 405B) могут конкурировать с проприетарными решениями. - **Потенциал для адаптации:** Результаты помогают пользователям оптимизировать выбор модели в зависимости от задачи и доступных ресурсов.

**Аспект 3: Влияние длины задачи - Прямая применимость:** Пользователи могут формулировать краткие запросы без существенной потери качества результатов, что ускоряет рабочий процесс. - **Концептуальная ценность:** Понимание того, что модели могут эффективно работать с краткими инструкциями, меняет подход к взаимодействию с ними. - **Потенциал для адаптации:** Разработчики интерфейсов могут создавать более удобные системы, не требующие подробных инструкций.

**Аспект 4: Сравнение библиотек визуализации - Прямая применимость:** Пользователи узнают, что модели лучше справляются с популярными библиотеками (Matplotlib, Seaborn) и могут учитывать это при выборе инструментов. - **Концептуальная ценность:** Понимание ограничений моделей при работе с менее распространенными библиотеками (Plotly). - **Потенциал для адаптации:** Пользователи могут адаптировать свои запросы, учитывая эти ограничения, или использовать дополнительные подсказки для работы с менее известными библиотеками.

**Аспект 5: Методология оценки - Прямая применимость:** Двойная система оценки может быть использована пользователями для оценки качества генерируемого кода. - **Концептуальная ценность:** Понимание разницы между визуальным сходством и функциональным соответствием задаче. - **Потенциал для адаптации:** Методология может быть адаптирована для оценки других задач генерации кода.

## Prompt:

Использование знаний из исследования PandasPlotBench в промтах для GPT На основе предоставленного исследования о способностях языковых моделей генерировать код для визуализации данных, можно создать более эффективные промты для GPT. Вот ключевые выводы и их практическое применение:

## Ключевые принципы для эффективных промтов:

**Предпочтение популярным библиотекам:** Используйте Matplotlib и Seaborn вместо Plotly **Оптимальное описание данных:** Включайте первые 5 строк DataFrame и типы столбцов **Краткость инструкций:** Даже короткие инструкции (1 предложение) работают эффективно **Выбор мощных моделей:** GPT-4o и подобные дают лучшие результаты для сложных визуализаций ## Пример эффективного промта:

[=====] Создай код для визуализации данных используя библиотеку Matplotlib.

Мой DataFrame (первые 5 строк): [=====] год продажи рост\_процент регион 0 2018 1200 NaN Север 1 2019 1350 12.5 Север 2 2020 1100 -18.5 Север 3 2021 1450 31.8 Север 4 2022 1600 10.3 Север [=====]

Типы данных: год: int64 продажи: int64 рост\_процент: float64 регион: object

Задание: Построй линейный график продаж по годам с точками и подписями значений, добавь вторую ось Y для процента роста. [=====]

## Почему это работает:

- Библиотека: Явно указана Matplotlib, с которой GPT работает лучше всего (75/89 баллов)
- Данные: Предоставлены первые 5 строк и типы данных, что дает оптимальное понимание структуры
- Лаконичность: Задание сформулировано в одном предложении, но содержит все необходимые детали
- Конкретность: Четко указаны требования к графику (линейный, с точками, подписями, вторая ось Y)

Такой подход к составлению промтов позволяет получить максимально качественный код для визуализации данных, минимизируя вероятность ошибок и неточностей в сгенерированном коде.