

EPIC: Эффективная подсказка для синтеза данных с несбалансированными классами в классификации табличных данных с использованием больших языковых моделей

Дата: 2025-01-13 00:00:00

Ссылка на исследование: <https://arxiv.org/pdf/2404.12404>

Рейтинг: 82

Адаптивность: 90

Ключевые выводы:

Исследование направлено на изучение эффективности использования больших языковых моделей (LLM) для генерации синтетических табличных данных, особенно для решения проблемы несбалансированности классов. Основным результатом - разработка метода EPIC, который использует сбалансированные группированные образцы данных и уникальное отображение переменных для создания качественных синтетических данных, значительно улучшающих производительность классификации машинного обучения.

Объяснение метода:

EPIC предлагает готовые шаблоны промптов для генерации качественных табличных данных с решением проблемы несбалансированных классов. Метод не требует дообучения, работает с различными LLM, включая открытые модели, и демонстрирует превосходные результаты на реальных данных. Подход понятен и доступен пользователям с базовым пониманием работы с данными.

Ключевые аспекты исследования: 1. **EPIC (Effective Prompting for Imbalanced-Class Data Synthesis)** - метод использования языковых моделей (LLM) для генерации синтетических табличных данных с акцентом на решение проблемы несбалансированных классов.

Структурированный подход к промптам - исследование определяет оптимальные компоненты промптов для генерации качественных табличных данных: CSV-формат, группировка примеров по классам, сбалансированная выборка классов и уникальное отображение переменных.

Эффективность без дообучения - метод использует in-context learning (обучение в

контексте) без необходимости дополнительного обучения моделей, что делает его доступным для широкого круга пользователей.

Превосходство над существующими методами - EPIC демонстрирует лучшие результаты на шести реальных наборах данных по сравнению с современными методами, особенно в улучшении качества классификации для миноритарных классов.

Итеративная генерация данных - подход предполагает многократную генерацию синтетических образцов с сохранением распределений и корреляций признаков, аналогичных исходным данным.

Дополнение: Для работы методов исследования EPIC **не требуется** дообучение или специальный API. Авторы используют in-context learning (обучение в контексте) существующих LLM без изменения их параметров. Хотя в исследовании для экспериментов используется API моделей GPT-3.5, Mistral и Llama2, сами методы полностью применимы в стандартном чате с LLM.

Основные концепции и подходы, которые можно применить в стандартном чате:

Структурированное CSV-форматирование - представление табличных данных в виде простых CSV-строк с указанием имен признаков в начале. Это позволяет эффективно передавать структурированную информацию модели.

Балансировка классов в примерах - вместо сохранения исходного распределения классов, включение равного количества примеров для каждого класса. Это помогает модели лучше понимать миноритарные классы.

Группировка примеров по классам - организация примеров в группы по целевому классу, что помогает модели лучше улавливать особенности каждого класса.

Повторение структуры - использование нескольких наборов примеров с одинаковой структурой, что помогает модели распознавать паттерны.

Триггер для генерации - размещение заголовка с именами признаков в конце промпта для запуска генерации новых данных.

Итеративная генерация - многократное применение метода с новыми примерами для создания разнообразных синтетических данных, охватывающих все распределение.

Применяя эти концепции в стандартном чате, пользователи могут:

- Генерировать синтетические примеры для несбалансированных данных
- Создавать дополнительные данные для миноритарных классов
- Улучшать качество классификационных моделей
- Получать данные с сохранением корреляций между признаками
- Расширять наборы данных для тестирования и обучения

Уникальное отображение переменных также можно реализовать в стандартном

чате, хотя это потребует предварительной обработки, которую можно выполнить вручную для небольших наборов данных.

Анализ практической применимости: 1. **Структурированный подход к промптам:** - Прямая применимость: Высокая. Пользователи могут применять предложенные структуры промптов для улучшения генерации табличных данных в любой LLM. Четкое описание форматирования CSV, группировки по классам и балансировки выборки делает метод доступным. - Концептуальная ценность: Высокая. Исследование объясняет, почему определенные компоненты промптов эффективны, что помогает пользователям понимать принципы взаимодействия с LLM. - Потенциал для адаптации: Очень высокий. Методы можно адаптировать для различных задач, включая генерацию данных в других форматах или для других целей.

Уникальное отображение переменных: Прямая применимость: Средняя. Техника полезна для улучшения качества генерации, но требует предварительной обработки данных. Концептуальная ценность: Высокая. Подход демонстрирует, как помочь LLM различать похожие категориальные переменные, что важно для понимания ограничений моделей. Потенциал для адаптации: Высокий. Метод можно адаптировать для различных типов данных и задач, требующих четкого различения категорий.

Итеративная генерация данных:

Прямая применимость: Высокая. Пользователи могут напрямую использовать многократную генерацию для создания более репрезентативных наборов данных. Концептуальная ценность: Средняя. Подход демонстрирует ограничения LLM в охвате полного распределения данных в одном запросе. Потенциал для адаптации: Высокий. Метод можно применять для различных генеративных задач, где требуется разнообразие выходных данных.

Применение для несбалансированных классов:

Прямая применимость: Очень высокая. Метод напрямую решает распространенную проблему в ML - несбалансированные данные. Концептуальная ценность: Высокая. Исследование показывает, как правильное форматирование промптов может помочь LLM генерировать качественные данные для миноритарных классов. Потенциал для адаптации: Высокий. Принципы можно применить к другим задачам с несбалансированными распределениями.

Эффективность без дообучения:

Прямая применимость: Очень высокая. Подход не требует специальных технических навыков или вычислительных ресурсов для дообучения. Концептуальная ценность: Высокая. Демонстрирует возможности in-context learning в LLM. Потенциал для адаптации: Высокий. Методологию можно применить к широкому спектру задач, где нежелательно или невозможно дообучение.

Prompt:

Использование исследования EPIC в промптах для GPT ## Ключевые принципы из исследования

CSV-стиль форматирования - более эффективен, чем текстовое описание
Сбалансированное представление классов - помогает модели понять все классы равномерно
Группировка по классам - улучшает распознавание паттернов в данных
Уникальное отображение переменных - помогает LLM различать похожие категориальные переменные
Триггеры завершения вместо прямых инструкций ##
Пример промпта на основе EPIC

[=====] Сгенерируй синтетические данные для классификации клиентов банка по кредитоспособности.

Вот примеры данных, сгруппированные по классам:

#	Класс:	Высокая	кредитоспособность
возраст,доход,стаж_работы,кредитная_история,текущие_кредиты,			результат
42,120000,15,A7X,0,высокая	38,95000,12,B3Y,1,высокая	51,150000,20,A2Z,0,высокая	

#	Класс:	Средняя	кредитоспособность
возраст,доход,стаж_работы,кредитная_история,текущие_кредиты,			результат
35,65000,8,C4X,2,средняя	45,72000,10,B9Y,1,средняя	29,58000,4,C1Z,2,средняя	

#	Класс:	Низкая	кредитоспособность
возраст,доход,стаж_работы,кредитная_история,текущие_кредиты,			результат
27,35000,2,D8X,3,низкая	52,42000,25,F2Y,5,низкая	33,38000,3,E5Z,4,низкая	

Сгенерируй 15 новых записей, сбалансированных по всем классам (по 5 для каждого класса).

возраст,доход,стаж_работы,кредитная_история,текущие_кредиты,	результат
[=====]	

Почему этот промпт работает эффективно

Структура CSV позволяет GPT четко понимать формат данных и экономит токены
Группировка примеров по классам помогает модели лучше понять характеристики каждого класса
Сбалансированное представление - каждый класс представлен равным количеством примеров
Уникальное кодирование для категориальной переменной "кредитная_история" (A7X, B3Y и т.д.)
Триггер завершения - последняя строка содержит только заголовки, что естественно подсказывает модели продолжить генерацию в том же формате
Такой подход позволяет получить более качественные синтетические данные, сохраняющие характеристики исходного распределения, особенно для задач с несбалансированными классами, где традиционные методы часто показывают низкую эффективность.

