

Расширенное Задание: Анализ Рыночных Цен Автомобилей с Построчной Обработкой и Анализом в Hive

Цель: Создать программу на Python для анализа средних рыночных цен на автомобили различных марок с последующей обработкой данных в Apache Hive для определения самой дорогой и самой дешевой марки автомобиля.

Датасет: Используйте доступный на Kaggle датасет, включающий марку, модель, цену и другие характеристики автомобилей

(<https://www.kaggle.com/datasets/CooperUnion/cardataset>).

Основные Задачи:

Загрузка Данных: Автоматическая загрузка датасета, если он отсутствует локально.

Построчная Обработка Данных:

- Используйте map и reduce для обработки данных из CSV-файла построчно.
- mapper: извлекает марку и цену из каждой строки.
- reducer: накапливает данные по маркам и ценам.

Анализ Данных:

- Вычисление средней цены для каждой марки.

Сохранение Результатов:

- Сохраните результаты в новый CSV-файл.

Обработка Данных в Hive:

- Загрузите полученный CSV-файл в таблицу Hive.
- Используйте HiveQL для анализа данных и нахождения самой дорогой и самой дешевой марки автомобиля.

Дополнительные Условия:

- Построчная Обработка: Эффективная обработка больших объемов данных.
- Интеграция с Hive: Умение работать с Big Data и проводить аналитические запросы.

Результат: Скрипт Python для определения средних цен автомобилей, сохранение результатов в формате CSV, последующая обработка данных в Hive и вывод самой дорогой и самой дешевой марки автомобиля.