Итоговая аттестация Проект №2 Анализ рынка акций

Kypc 1T: Data Engineer Тарасов Сергей Евгеньевич

Задание

Общая задача: создать ETL-процесс формирования витрин данных для анализа изменений курса валют.

- Разработать скрипты загрузки данных в 2-х режимах:
 - ▶ Инициализирующий
 - Инкрементальный
- Организовать правильную структуру хранения данных:
 - Сырой слой данных
 - ▶ Промежуточный слой
 - ▶ Слой витрин
- ▶ Источник данных: API https://www.alphavantage.co

План реализации проекта

- 1. Анализ источника данных
- 2. Проектирование ER-диаграммы
- 3. Определение стека технологий
- 4. Реализация инициирующей загрузки данных
- 5. Реализация инкрементальной загрузки данных

Анализ источника данных

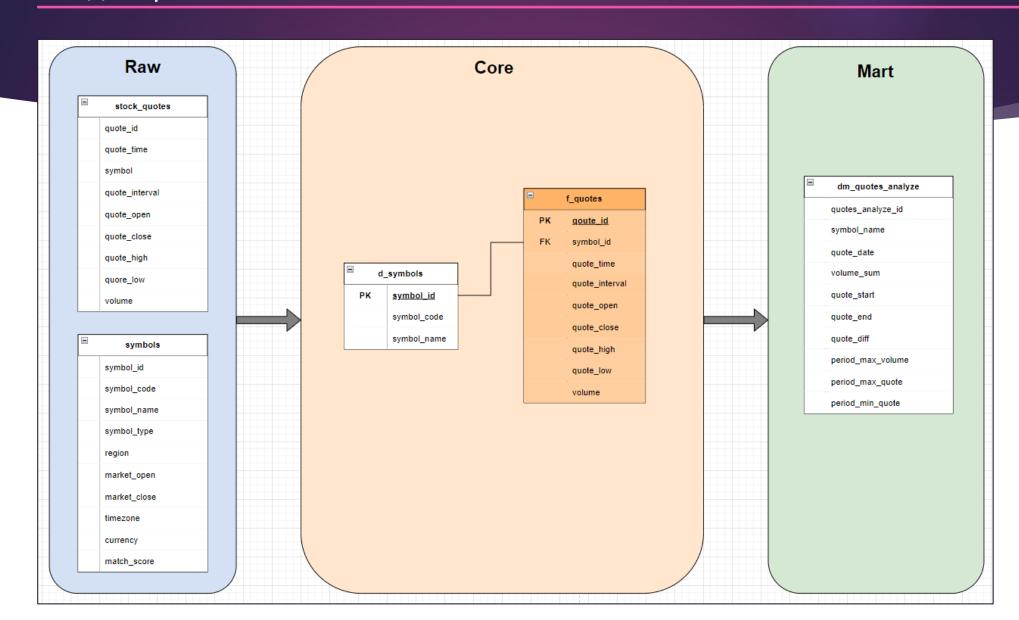
В качестве источника данных используется API сервиса https://www.alphavantage.co.

В процессе изучения API Alpha vantage, были определены параметры, которые позволяют получить необходимую для построения витрины информацию.

Параметры:

- function: TIME_SERIES_INTRADAY
- symbol: аббревиатура акции (напр. IBM)
- ▶ Interval: 5min
- ▶ month: месяц, за который необходимо получить данные (напр. 2023-12)
- apikey: индивидуальный ключ для использования API

ER-диаграмма



Стек технологий

Python: Используется для написание скриптов получения данных по API, трансформации данных а также записи данных в $Б\Delta$.

Postgres: Используется для хранения данных. Выбрал postgres, т.к. предполагается не значительный объем базы данных. А также из-за простоты и удобства работы для разработчика.

Docker: Используется для быстрой и удобной сборки и развертывания проекта на различных машинах.

Crontab: Очень хотелось использовать airflow для создания отдельных дагов для инициализации и инкрементальной загрузки данных, но к сожаления на моем ПК не удалось настроить работу даже версии 1.х. Пришлось воспользоваться crontab для автоматизации этих процессов.

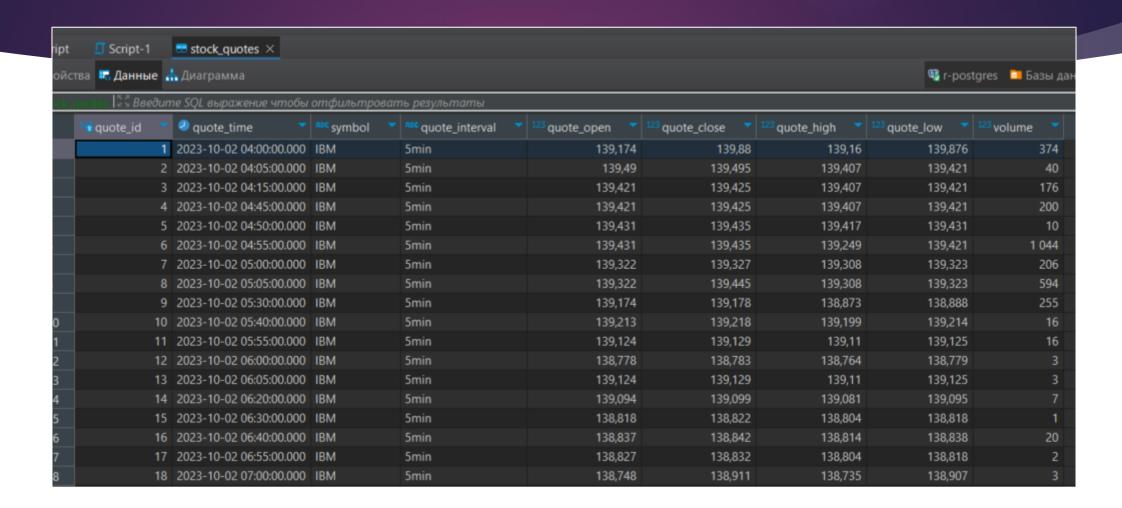
Этап инициализации данных

- Создаются три слоя данных, и таблицы в них.
- Таблицы наполняются предзагруженными данными из csv файлов.
 Эти файлы были созданы, для обхода ограничения API на 25 запросов в сутки. В них собраны данные из предыдущих периодов
- Происходит актуализация данных, догружаются данные, которых не хватает в csv-файлах.
- Происходит наполнение всех слоев данными, и расчет итоговой витрины

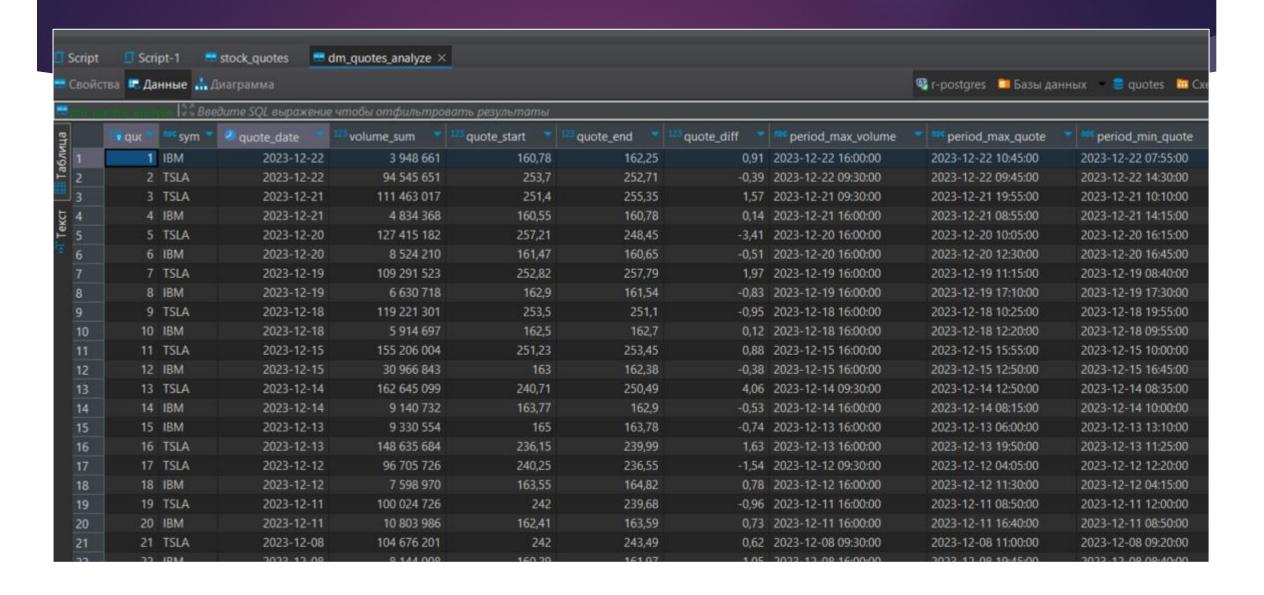
Этап инкрементальной загрузки данных

- Инкрементальная загрузка догружает данные за последний день
- ▶ Происходит пересчет итоговой витрины
- ▶ Запускается с помощью cron, настроена на запуск в 6:00

Пример входных данных (raw слой)



Пример данных в итоговой витрине (mart слой)



Выводы

- ▶ В процессе разработки был создан ETL-процесс получения, преобразования и анализа курса акций
- Полученные результаты можно использовать для анализа курса акций и принятия решений о торговых операциях на бирже акций.
- ▶ В процессе выполнения работы закреплены навыки проектирования DWH и разработки ETL-процессов.