

Project Report

Bohdan Vey, Maksym Kuzushun

Big Data

Precomputed reports and streaming data

Для читання стімінгу ми вирішили використовувати Kafka, оскільки це основне завдання Kafka, для якого вона створювалась.

Дані ми вирішили зберігати в двох форматах:

- 1) Pandas Dataframe - дану таблицку ми вирішили використовувати, для погодинних запитів, оскільки ми можемо зберігати необхідну інформацію про запити за останню годину локально на комп'ютері, через що наш сервер не буде витрачати час на те щоб зв'язатись з іншим сервером і дістати необхідні дані

Ми також зберігаємо дані в одному форматі для всіх трьох запитів і використовуємо threads(тут також можна використовувати інші сервери), для того щоб обробити ці дані і дати відповідь на запитання, які нас цікавлять.

Самі відповіді на запити ми зберігаємо в json форматі, так як цього потребує завдання, в папці data в файлах question_{1..3}.json

Ad-hoc queries

Кожен раз, коли ми отримуємо новий запис в кафці, ми відправляємо його в функцію, яка буде записувати новий запис у таблиці Cassandra. Ми маємо 5 різних таблиць, і 5 різних функцій для запису в ці таблиці, скрипт можна знайти [тут](#). Схему таблиць можете бачити на діаграмі внизу. [Посилання на cql файл](#) для ініціалізації keyspace.

Для відповідей на запити у нас є Python скрипт який вміє давати відповіді на будь-який з п'ятих запитів. [Посилання на папку з скриптом](#)

