

Apache Airflow

1. Написание кода DAG'а.

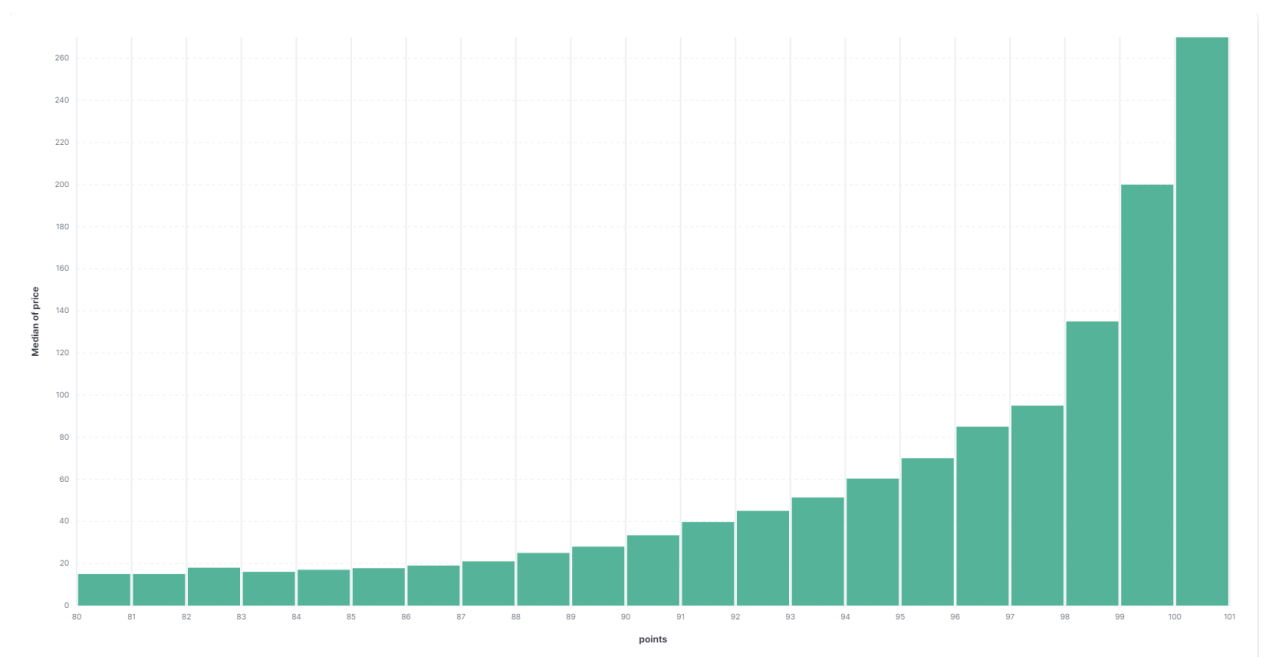
Были написаны python функции, которые были потом привязаны к задачам.

В целом проблем не было за исключением, что было изначально непонятно как передавать данные между задачами и проблемы с kibana. В первой реализации я попробовал использовать xcom, но иногда вылетали ошибку из-за недостаточного количества памяти, поэтому я переделал, чтобы промежуточные данные сохранялись .csv файлы в папке data, и следующие задачи брали эти данные из этих .csv файлов. Проблема с kibana заключалась в том, что данные не сохранялись, решилась она предварительным переформатированием строки в json.

2. Apache nifi.

Мягко говоря, apache nifi мне совсем не понравился в связи с огромным количеством проблем. Одна из них – это крайне скудная документация к процессорам, да и в целом информации о данном инструменте довольно мало. По ходу выполнения были выявлены неочевидные особенности, такие как: невозможность работы с большим количеством файлов в рамках одного процессора, приходилось заранее объединять их в подгруппы, невозможность объединить несколько файлов с разными названиями, необходимо предварительно переименовать все файлы на одно имя.

Используя данные, обработанные с помощью пайплайна, в была построена данная гистограмма.



График