

Интеллектуальный анализ работы хранилища данных на основании обработки ЛОГОВ

ТРЕК ОТ РОСТЕЛЕКОМ.

ЭКСПЕРТЫ - АНДРЕЙ ТЕЛЮКОВ, НИКИТА БОГДАНОВ.

КОМАНДА ВВУВУ - УЛЛУБИЙ КАГЕРМАНОВ, НИЗАМОВА ЛИЛИЯ,
ВЕРЕТНОВ ТИМОФЕЙ

Идея проекта:

- ▶ НАЙТИ ВЫСОКОНАГРУЖЕННЫЕ “ГОРЯЧИЕ” ОБЪЕКТЫ ХРАНИЛИЩА ДАННЫХ И СПРОГНОЗИРОВАТЬ ПОТЕНЦИАЛЬНЫЙ РОСТ НАГРУЗКИ.

Решение проблемы:

- ▶ Разработали алгоритм быстрой обработки данных.
- ▶ Выделили отдельные объекты, к которым чаще всего обращаются пользователи, как самых “тяжёлых” и самых “горячих”.
- ▶ Смоделировали прогнозирование нагрузки на систему в разные временные отрезки.

CJM ПОЛЬЗОВАТЕЛЯ ПРОДУКТА

1. Клиенты:

- ▶ пользователи, дата инженеры, администраторы баз данных

2. Для пользователей

- ▶ информация о приблизительном времени на запрос

3. Для инженеров и администраторов

- ▶ информация о "горячих" и "тяжелых" объектах и выделение под них дополнительных ресурсов,
- ▶ информация о периоде времени нагрузки на систему либо в процессе выполнения запросов либо в предсказательных целях



CJM ПОЛЬЗОВАТЕЛЯ ПРОДУКТА

Как использовать:

- ▶ телеграм бот или встроенный к код python скрипт
- ▶ ВВОД ДАННЫХ И ПОЛУЧЕНИЕ ОТВЕТА:
 1. топ(n) "горячих" таблиц"
 2. топ(n) "тяжелых таблиц"
 3. имя юзера(loguser) и запрос с получением среднего времени его исполнения.

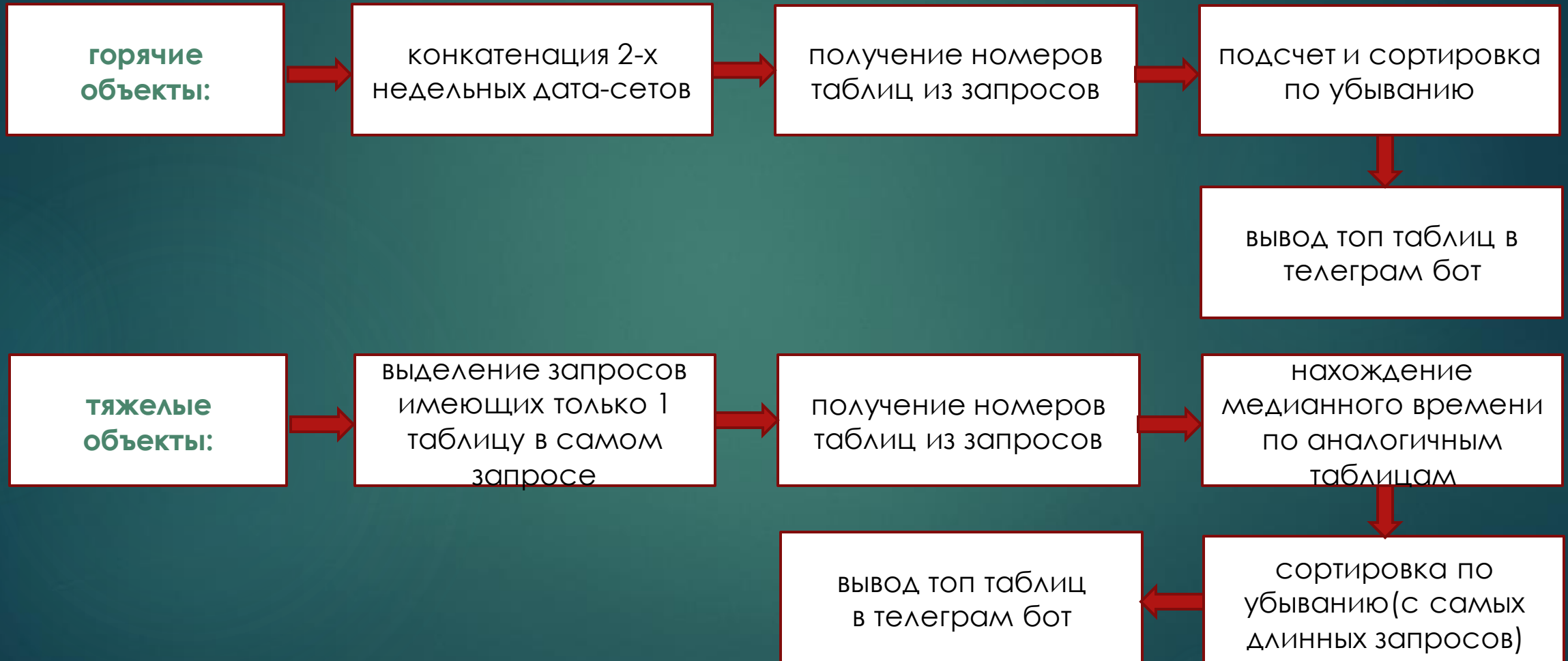


Перечень ИСПОЛЬЗУЕМЫХ ТЕХНОЛОГИЙ

► Используемый стек: Python

pandas, numpy, collections.Counter,
catboos (CatBoostRegressor, cv, Pool),
sklearn.model_selection (train_test_split), re,
sklearn.feature_extraction.text (CountVectorizer),
scipy.sparse, telebot, config.tokens

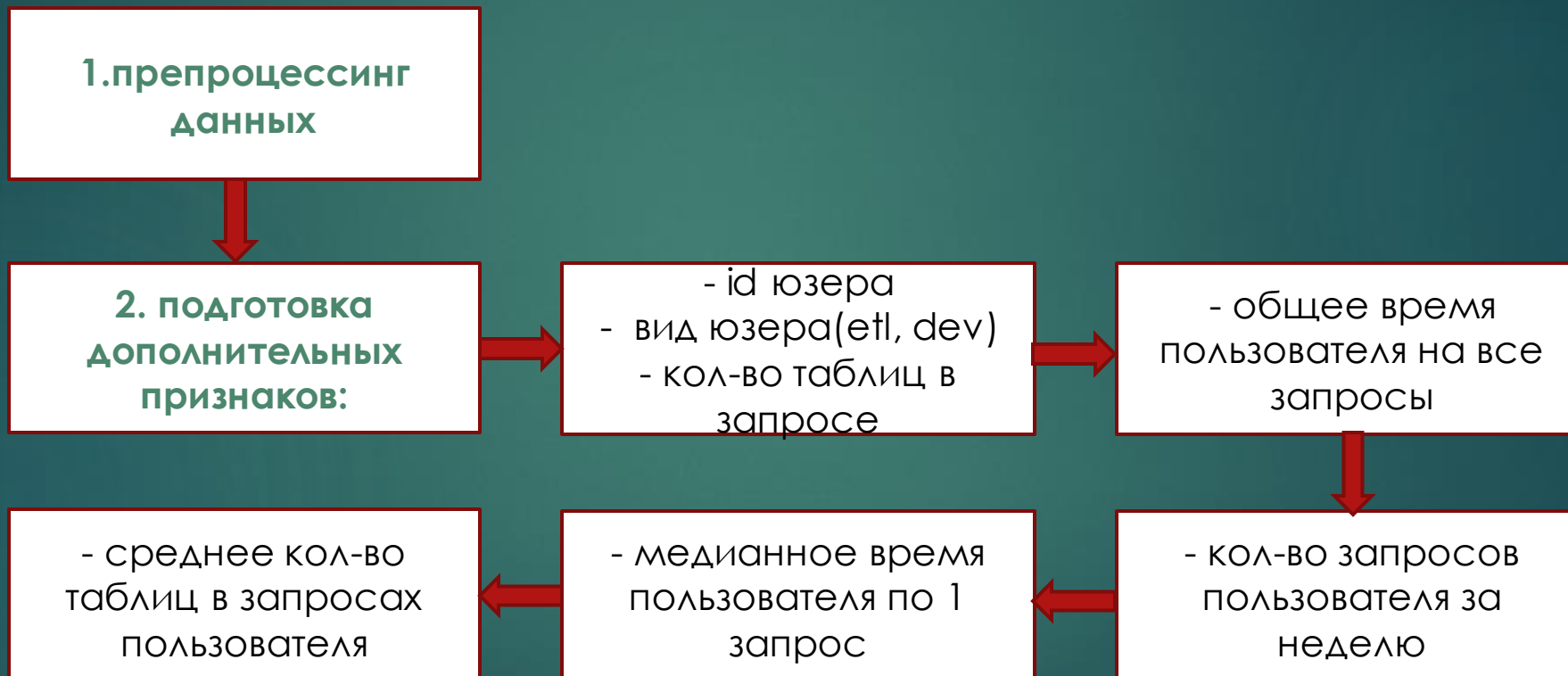
Технические характеристики решения



Технические характеристики решения



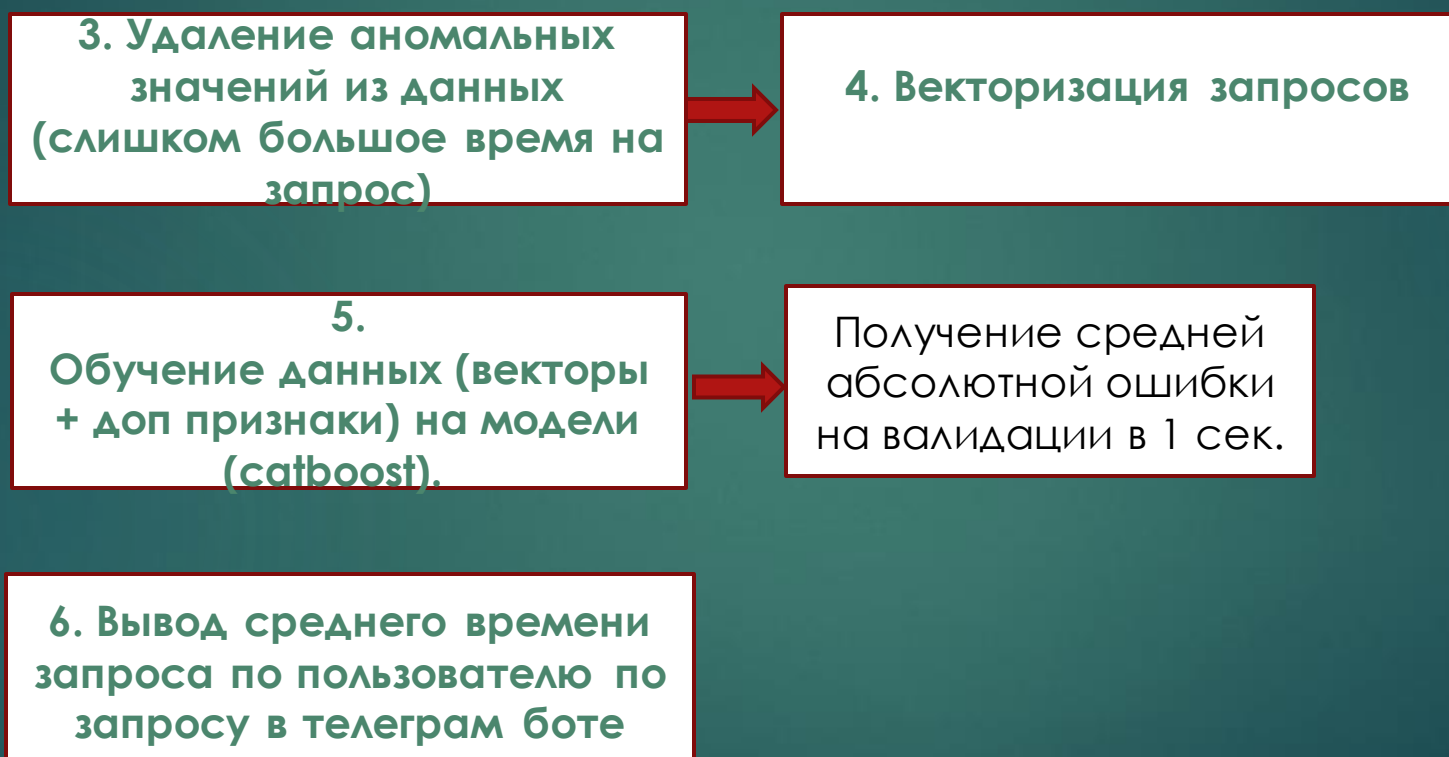
Прогнозирования времени исполнения запроса по
предоставленным данным



Технические характеристики решения



Прогнозирования времени исполнения запроса по предоставленным данным





Дополните- льная информация

Для отображения полученных нами данных мы используем телеграм бота

Функционал:

- ▶ оценка времени
- ▶ по шаблонным запросам
- ▶ по введёному запросу
- ▶ получения Топ Горячих и Тяжелый Таблиц
- ▶ расчет среднего времени выполнения запроса в зависимости от таблицы и пользователя

Что не успели сделать

- ▶ оформить весь ресёрч из jupyter notebook в продакт версию с py файлами,
- ▶ связать телеграм бота с нашими алгоритмами и моделью



Как можно
было бы
улучшить, имея
дополнительные
данные:

- ▶ улучшить точность прогнозирования времени запроса, имея полный запрос с агрегациями и сортировками.
- ▶ прогнозировать загруженность таблиц по дням, неделям и времени суток, имея timestamp
- ▶ прогнозировать загруженность нагрузки на систему(ЦМ, RAM, SDD), имея данные о нагрузках