

# **BISONS**

Поиск аномалий: Автоматизация процесса определения сбоя в поставках данных

### summary

данные

EDA

выбросы

ML модели

Для оценки валидности датасета была использована легко внедряемая модель Логистической регрессии. Для оценки модели использовались метрики ROC AUC и F1 score.

#### EDA

#### Найденные инсайты:

- Существует id, для которого не определен benchmark\_name
- Сезонность в количестве уникальных выгружаемых значений отсутствует
- медианное количество бенчмарков за одну поставку - 458

#### ML МОДЕЛИ

- Выделили топ-10 наиболее важных признаков
- точность выбранной модели:

- Провели сравнение 5 ML моделей
- Датасет от 15 марта валиден



#### Описание базы данных

Количество строк: 995 004

База данных содержит передачи данных с котировками бенчмарков от ДЗО с 09-04-2018 по 15-03-2024. То есть, каждый день подгружается отдельный пакет с котировками нескольких бенчмарков

\*ДЗО— дочернее зависимое общество, организация, в уставном капитале которой Сбер имеет контрольную долю



#### Задача кейса

Сделать ml-модель, которая будет своевременно выявлять ошибки в поставках данных, то есть находить в них выбросы



Количество столбов: 6

benchmark id - ID бенчмарка

benchmark\_name-Наименование бенчмарка

date\_of - Дата котировки

quote - Котировка

ctl\_loading - Идентификатор (id) загрузки

ctl\_loading\_date - Дата загрузки

В базе данных есть 4 названия бенчмарков, которым соответствует несколько id бенчмарков

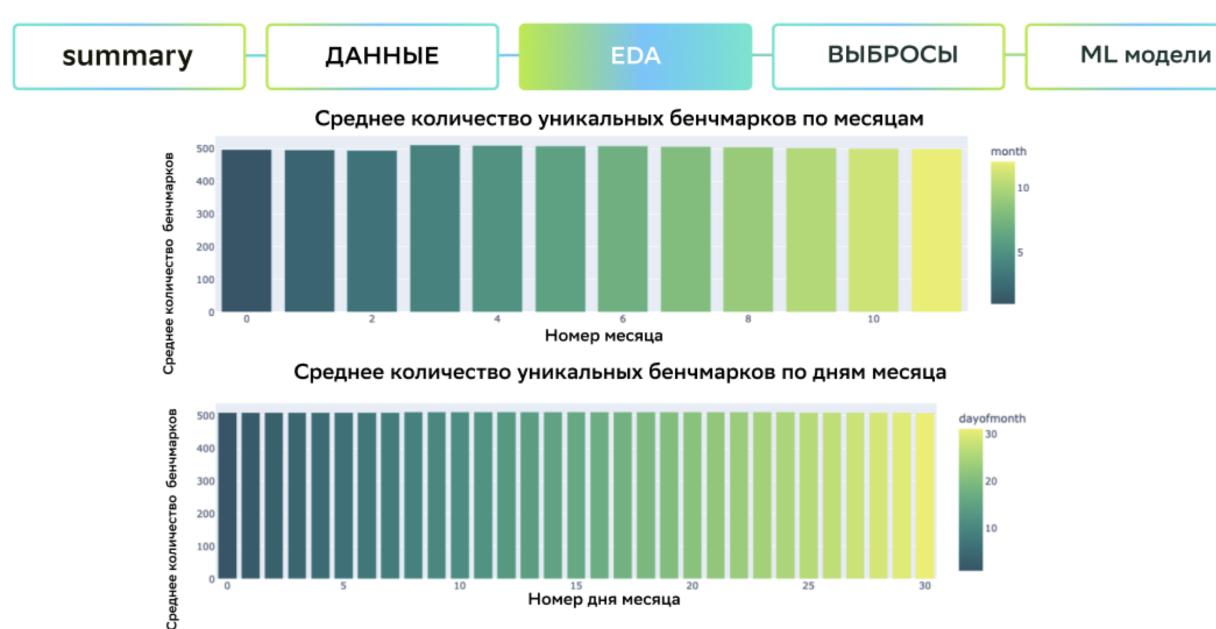
#### Пример:

бенчмарку "Technology Select Sector Index" соответствуют индексы

"8D8AADAB-A754-4C02-9C81-0E0B4890856В" "51D93A6F-358F-4AD9-9C70-C00BAE62B98C" Существует id, для которого не определен benchmark\_name

- Сделали единые id для бенчмарков с одинаковыми названиями
- Для каждой поставки создали массив с id бенчмарков из этой поставки

- Есть 176 дней, в которые было больше одной поставки за день из-за дополнительных поставок данных
  - → В эти дни были допоставки данных
  - При создании модели мы не классифицируем допоставки как ошибки





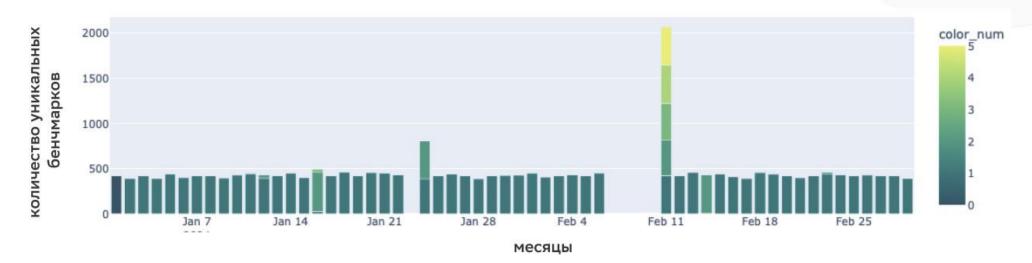


#### Текущий механизм отгрузки котировок бенчмарков

- происходит отгрузка бенчмарков
- если значение не соответствует норме (скорее всего это на текущий момент определяется вручную), отправляется запрос на доотправку пропущенных значений
- через несколько дней происходит повторная отправка пропущенных ранее данных, если выбросы заметили

ИНСАЙТ: иногда после выбросов доотправки не происходит, так как присутствует человеческий фактор)







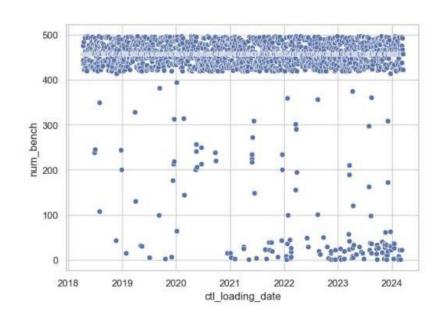


Мы измеряем с помощью межквантильного размаха [403, 503] Заметим четкую линию на уровне 458 Это медиана в нашей выборке!

Со временем количество выбросов увеличивается Особенно начиная с 2021 года.

Это может быть связано с ростом нагрузки на систему

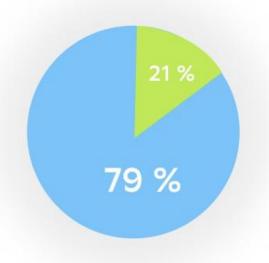






### **У** ЧТО МЫ СЧИТАЕМ ВЫБРОСОМ?

Выброс - это отклонение , которое которое не входит в размах



21% выбросов компенсировалиотдельной поставкой

79% выбросов которые не компенсировали отдельной поставкой

- добавили новый признак, чтобы модель отличала выброс от допоставки
- мы заметили, что медиана количества бенчмарков в рамках одной поставки на протяжении всех 6 лет - это 458.
- можно предположить, что существует всего 458 бенчмарков, по каждому из которых каждый день как минимум один раз должна выгружаться котировка

- если значение бенчмарков в рамках одной поставки совсем маленькое через несколько дней происходит допоставка --> в сумме получается 458 (+- 10)
- при этом наша команда нашла случаи когда есть некритичные недопоставки, и в таком случае допоставок не происходит, что может негативно отразиться на бизнес-процессах.



ДАННЫЕ

**EDA** 

ВЫБРОСЫ

ML модели



Постановка задачи: бинарная классификация пакетов

### Пайплайн модели



4. Выделение признаков для модели
5. Кодируем категориальные признаки (OneHotEncoder)
6. Масштабирование численных признаков (StandartScaler)
8. Проверка модели на тестовой выборке

9. Оценка

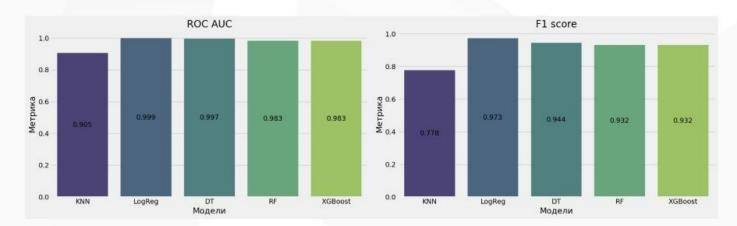


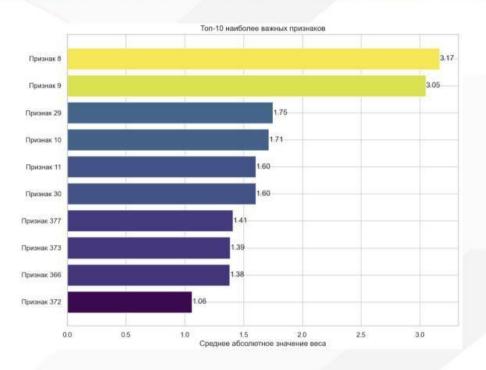
• Логистическая Регрессия

Оптимальный гиперпараметр С = 2,78

#### Плюсы модели:

- Скорость работы
- Простота реализации
- Устойчивость к переобучению





$$p_{+} = \frac{1}{1 + exp(-x^{T}w)} = \sigma(x^{T}w)$$

Вероятность выброса

Вектор Весов

(значения см. на графике выше)



🕑 Валидность датасета от 15 марта

Да





# **Google Colaboratory**

https://colab.research.google.com/drive/1pRhzQ1Hy3Pc3 QOi\_k\_zM9FDiH7m5T\_zF?usp=sharing

# Команда





Алиев Руслан
Аналитик
ниу вшэ
Бизнес-Информатика



Ванин Виктор
Аналитик
ниу вшэ
Бизнес-Информатика



Ковалева Дарья
Аналитик
НИУ ВШЭ
Бизнес-Информатика



Панкова Злата
Аналитик
НИУ ВШЭ
Бизнес-Информатика



Кислова Алена

Дизайнер

В&D

Цифровой дизайн