## Сравнительный анализ трех задач

Кузнецов Г.И.

### Введение

Проведем сравнение нескольких выполненных работ по нескольким ключевым позициям.

Для начала кратко вспомним основные концепции и задачи, поставленные в предыдущих работах. После чего сравним их согласно плану.

#### План сравнения:

- 1.Понимание бизнес-целей (Business Understanding)
- 2. Начальное изучение данных (Data Understanding)
- 3.Подготовка данных (Data Preparation)
- 4. Моделирование (Modeling)
- 5.Оценка (Evaluation)
- 6.Внедрение (Deployment)

# Работа 1: Создание предиктивной модели рейтинга мобильных приложений

- 1. Провести анализ данных
- 2. Ответить на вопросы:
  - как информация о приложении влияет на рейтинг пользователей?
  - чем отличается статистика приложений для разных групп?

|   | Unnamed:<br>0 | id        | track_name   | size_bytes | currency | price | rating_count_tot | rating_count_ver | user_rating | user_rating_ver |
|---|---------------|-----------|--|------------|----------|-------|------------------|------------------|-------------|-----------------|
| 0 | 1             | 281656475 | PAC-MAN<br>Premium   | 100788224  | USD      | 3.99  | 21292            | 26               | 4.0         | 4.5             |
| 1 | 2             | 281796108 | Evernote -<br>stay<br>organized                            | 158578688  | USD      | 0.00  | 161065           | 26               | 4.0         | 3.5             |
| 2 | 3             | 281940292 | WeatherBug - Local Weather, Radar, Maps, Alerts            | 100524032  | USD      | 0.00  | 188583           | 2822             | 3.5         | 4.5             |
| 3 | 4             | 282614216 | eBay: Best<br>App to Buy,<br>Sell, Save!<br>Online<br>Shop | 128512000  | USD      | 0.00  | 262241           | 649              | 4.0         | 4.5             |
| 4 | 5             | 282935706 | Bible  | 92774400   | USD      | 0.00  | 985920           | 5320             | 4.5         | 5.0             |

часть датасета

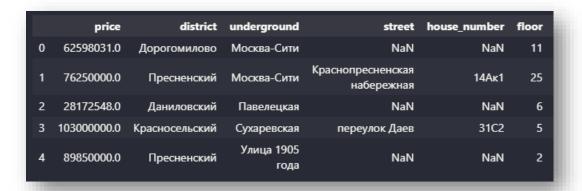
## Работа 2: Построение предиктивной модели оценки надежности заемщика

- 1. Провести разведочный анализ данных по данным скоринга.
- 2. Определить достаточность данных для определения кредитного скоринга, социального скоринга.

```
Data columns (total 62 columns):
    Column
                                   Dtype
                                   int64
                                   int64
    rn
    pre since opened
                                   int64
    pre since confirmed
                                   int64
    pre pterm
                                   int64
    pre fterm
                                   int64
    pre till pclose
                                   int64
    pre_till_fclose
                                   int64
    pre loans credit limit
                                   int64
    pre loans next pay summ
                                   int64
10 pre loans outstanding
                                   int64
11 pre loans total overdue
                                   int64
12 pre loans max overdue sum
                                   int64
13 pre loans credit cost rate
                                   int64
14 pre loans5
                                   int64
15 pre loans530
                                   int64
16 pre loans 3060
                                   int64
17 pre loans6090
                                   int64
18 pre loans90
                                   int64
19 is zero loans5
                                   int64
60 fclose flag
                                   int64
61 flag
                                   int64
```

## **Работа 3:** Разработка многофакторной модели для оценки стоимости недвижимости

- 1. Собрать данные по стоимости недвижимости в выбранном регионе.
- 2. Выбрать признаки, влияющие на стоимость разных типов недвижимости в регионе.
- 3. Вдвинуть гипотезы о том, какие признаки наиболее и наименее влияют на стоимость разных типов недвижимости, проверить гипотезы.
- 4. Смоделировать модель прогнозирования стоимости разных видов недвижимости в выбранном регионе на основе разных методов обучения, сравнить эффективность моделей.
- 5. Дать рекомендации покупателям и продавцам недвижимости.



#### часть датасета

| floors_count | total_meters | rooms_count | residential_complex | uri                                      |
|--------------|--------------|-------------|---------------------|--|
| 18           | 45.8         | 2           | Бадаевский ЖК       | https://www.cian.ru/sale/flat/313877981/ |
| 61           | 61.0         | 2           | Capital Towers      | https://www.cian.ru/sale/flat/315795006/ |
| 27           | 40.5         | 2           | Эра ЖК              | https://www.cian.ru/sale/flat/301260611/ |
| 5            | 152.0        | 3           | NaN                 | https://www.cian.ru/sale/flat/317646469/ |
| 22           | 99.7         | 2           | Лайф Тайм ЖК        | https://www.cian.ru/sale/flat/316363605/ |

## Понимание бизнес-целей (Business Understanding)

Перед началом любого анализа данных важно четко определить цели проекта.

В этих трех задачах мы работали с разными бизнес-потребностями: от предсказания рейтинга мобильных приложений до оценки кредитоспособности заемщиков и прогнозирования стоимости недвижимости.

Каждый проект требовал уникального подхода к постановке задачи и интерпретации результатов.

| Задача                 | Цель   | Особенности  |
|------------------------|--|--|
| Рейтинг<br>приложений  | Анализ влияния данных на рейтинг, сравнение групп приложений | Простой датасет,<br>мультиклассовая/рег<br>рессия                |
| Скоринг заемщиков      | Оценка надежности клиента (кредитный/социальный скоринг)     | Сильный дисбаланс<br>(97% vs 3%), сложный<br>формат данных (.pq) |
| Оценка<br>недвижимости | Прогнозирование<br>цены на основе<br>гибридных подходов      | Самостоятельный<br>сбор данных,<br>проверка гипотез              |

## Начальное изучение данных (Data Understanding)

На этом этапе мы провели первичный анализ данных, оценили их структуру, качество и потенциальные проблемы.

В зависимости от задачи, данные были представлены в разных форматах — от простых CSV-файлов до сложных паркеттаблиц.

Особое внимание пришлось уделить проблеме дисбаланса классов в задаче кредитного скоринга и самостоятельному сбору актуальных данных по недвижимости.

| Задача                | Источник данных                        | Проблемы  |
|-----------------------|--|---|
| Рейтинг<br>приложений | 2 CSV-файла                            | Простые, чистые<br>данные                               |
| Скоринг<br>заемщиков  | Паркет-файлы<br>(.pq)                  | Дисбаланс,<br>сложность<br>загрузки                     |
| Недвижимость          | Парсинг<br>актуальных<br>данных (2025) | Отсутствие<br>готового датасета,<br>feature engineering |

### Подготовка данных (Data Preparation)

Качество данных напрямую влияет на результат моделирования.

В каждом проекте мы выполнили очистку данных, обработку пропусков и преобразование признаков.

Для задачи с дисбалансом классов (3% положительных примеров) был применен метод undersampling. В случае с недвижимостью — проведен тщательный feature engineering на основе собранных данных.

| Задача                | Методы<br>обработки   | Фичи                                  |
|-----------------------|---|---------------------------------------|
| Рейтинг<br>приложений | Работа с<br>категориальными<br>признаками, в<br>остальном чисто | Категориальные +<br>числовые          |
| Скоринг<br>заемщиков  | Undersampling<br>(только 1 +<br>случайные 0),<br>обработка .pq  | Социальные +<br>кредитные<br>признаки |
| Недвижимость          | Парсинг, очистка,<br>отбор признаков                            | Геоданные,<br>параметры<br>объектов   |

## Моделирование (Modeling)

Для решения поставленных задач мы тестировали различные алгоритмы машинного обучения.

В проектах с рейтингом приложений и стоимостью недвижимости сравнивались четыре модели, тогда как в задаче кредитного скоринга из-за особенностей данных использовался только Random Forest. Выбор оптимальной модели осуществлялся на основе анализа метрик.

| Задача                | Модели  | Лучшая модель     |
|-----------------------|---|-------------------|
| Рейтинг<br>приложений | Линейная, Decision<br>Tree, Random Forest,<br>Gradient Boosting | Gradient Boosting |
| Скоринг<br>заемщиков  | Только Random Forest<br>(из-за сложности<br>данных)             | Random Forest     |
| Недвижимость          | Линейная, Decision<br>Tree, Random Forest,<br>Gradient Boosting | Gradient Boosting |

## Оценка (Evaluation)

Оценка качества моделей проводилась с учетом специфики каждой задачи. Для регрессии использовались метрики MSE и R<sup>2</sup>, для классификации — precision, recall и F1-мера.

Анализ важности признаков позволил выявить ключевые факторы, влияющие на целевую переменную в каждом случае.

| Задача                | Метрики   | Выводы  |
|-----------------------|---|---|
| Рейтинг<br>приложений | RMSE, R <sup>2</sup> , MAE  | Важность описания,<br>скриншотов  |
| Скоринг<br>заемщиков  | F1, ROC-AUC,<br>Classification report<br>(precision, recall,<br>accuracy) | Важность отдельных платежей, типа кредита, социальных признаков               |
| Недвижимость          | RMSE, MAE, R2   | Сильное влияние<br>площади,<br>инфраструктуры<br>(район, близость к<br>метро) |

### Внедрение (Deployment)

Результаты каждого проекта имеют практическую ценность для бизнеса.

Модель оценки рейтинга может помочь разработчикам улучшать свои приложения, кредитный скоринг — снизить риски банка, а прогнозирование цен на недвижимость — поддержать принятие решений покупателями и продавцами.

Успешное внедрение требует не только технической реализации, но и понятной интерпретации результатов для конечных пользователей.

| Задача                | Где применимо                 | Рекомендации                         |
|-----------------------|-------------------------------|--------------------------------------|
| Рейтинг<br>приложений | Рекомендации<br>разработчикам | Улучшение<br>описаний и<br>категорий |
| Скоринг<br>заемщиков  | Банковский скоринг            | Учет социальных<br>данных            |
| Недвижимость          | Риелторские<br>сервисы        | Динамическое<br>ценообразование      |

### Выводы (по сравнительному анализу)

#### 1. Разнообразие задач — разные подходы

Несмотря на общую структуру CRISP-DM, каждый проект потребовал уникальных решений: от борьбы с дисбалансом данных в скоринге до самостоятельного сбора актуальной информации по недвижимости.

#### 2. Качество данных — основа успеха

Наиболее сложным этапом оказалась подготовка данных, особенно при работе с «сырыми» источниками (парсинг, паркет-файлы). Глубокий EDA и feature engineering критически важны для построения качественных моделей.

#### 3. Выбор модели зависит от контекста

Gradient Boosting показал лучшие результаты в задачах регрессии. Простые модели (линейная регрессия) могут быть полезны для интерпретируемости.

#### 4. Практическая ценность

Все модели решают конкретные бизнес-задачи: от снижения рисков банка до рекомендаций по ценообразованию. Важно не только построить модель, но и обеспечить её внедрение в рабочие процессы.

#### Спасибо за внимание