

# **Проект "Подключения услуги – «Мегафон»"**

## Обзор данных

В качестве исходных данных представлена информация об отклике абонентов на предложение подключения одной из услуг. Каждому пользователю может быть сделано несколько предложений в разное время, каждое из которых он может или принять, или отклонить.

Отдельным набором данных является нормализованный анонимизированный набор признаков, характеризующий профиль потребления абонента. Эти данные привязаны к определенному времени, поскольку профиль абонента может меняться с течением времени.

Данные train и test разбиты по периодам – на train доступно 4 месяцев, а на test отложен последующий месяц.

Итого, в качестве входных данных представлены:

- **data\_train.csv:**
  - id,
  - vas\_id
  - buy\_time
  - target
- **features.csv.zip:**
  - id
  - feature\_list
- **data\_test.csv:**
  - id
  - vas\_id
  - buy\_time

### Описание датасета

- **id** - идентификатор абонента
- **vas\_id** - подключаемая услуга
- **buy\_time** - время покупки, представлено в формате timestamp, для работы с этим столбцом понадобится функция `datetime.fromtimestamp` из модуля `datetime`.
- **target** - целевая переменная, где 1 означает подключение услуги, 0 - абонент не подключил услугу соответственно.

# Информация о модели, ее параметрах, особенностях и основных результатах

## Задача

Требуется на основании имеющихся данных построить алгоритм, который для каждой пары пользователь-услуга определит вероятность подключения услуги.

## Модель CatBoost

Для решения задачи применена модель CatBoost со следующими параметрами:

Константные параметры:

- |  |  |
|--|--|
| • <code>loss_function='Logloss'</code>       | - показатель, используемый для обучения                                      |
| • <code>eval_metric='F1'</code>              | - метрика, используемая для обнаружения переобучения                         |
| • <code>auto_class_weights='Balanced'</code> | - автоматический подбор весов для балансировки классов                       |
| • <code>random_state=42</code>               | - случайное зерно, используемое для обучения                                 |
| • <code>logging_level='Verbose'</code>       | - вывод оптимизированных метрик, затраченного и оставшегося времени обучения |
| • <code>task_type='GPU'</code>               | - используется CPU или GPU. По умолчанию стоит CPU                           |
| • <code>cat_features=f_categorical</code>    | - массив с категориальными признаками  |
| • <code>one_hot_max_size=20</code>           | - максимальное количество уникальных значений среди категориальных признаков |
| • <code>early_stopping_rounds=50</code>      | - отслеживание переобучения  |

Лучшие подбираемые параметры (с использованием сетки гиперпараметров):

- |  |  |
|--|--|
| • <code>depth=10</code>                | - глубина дерева   |
| • <code>learning_rate=0.03</code>      | - скорость обучения  |
| • <code>iterations=100</code>          | - максимальное количество построенных деревьев                   |
| • <code>l2_leaf_reg=20.0</code>        | - коэффициент при члене регуляризации L2 функции потерь          |
| • <code>bagging_temperature=2.0</code> | - настройка интенсивности байесовского бутстрапа, по умолчанию=1 |

Подбор гиперпараметров модели CatBoost выполняется при помощи рандомизированного поиска по сетке с использованием кросс-валидации, проверяется 30 наборов гиперпараметров.

## Результаты модели CatBoost

F1 = 0.47 по качеству прогноза для класса 1 – подключение услуги абонентом.

AUC\_ROC = 0.859

F1 = 0.92 по качеству прогноза для класса 0 – не подключение услуги абонентом.

AUC\_PR = 0.353

# Обоснование выбора модели и ее сравнение с альтернативами

Модель CatBoost имеет чуть лучший показатель F1 по сравнению с альтернативной моделью логистической регрессии.

## Модель логистической регрессии

Построение модели логистической регрессии выполняется с автоматической балансировкой классов `class_weight='balanced'` с применением пайплайнов:

- Предобработка: `StandardScaler()`, `OneHotEncoder()`
- Селекция: `SelectPercentile()`
- Модель: `LogisticRegression()`

Подбор гиперпараметров модели `LogisticRegression` выполняется при помощи поиска по сетке с использованием кросс-валидации.

Подобранные параметры модели логистической регрессии:

- `model__C=5` - обратная сила регуляризации
- `selector__percentile=5` - процент лучших признаков

## Результаты модели логистической регрессии

F1 = 0.46 по качеству прогноза для класса 1 – подключение услуги абонентом.

AUC\_ROC = 0.845

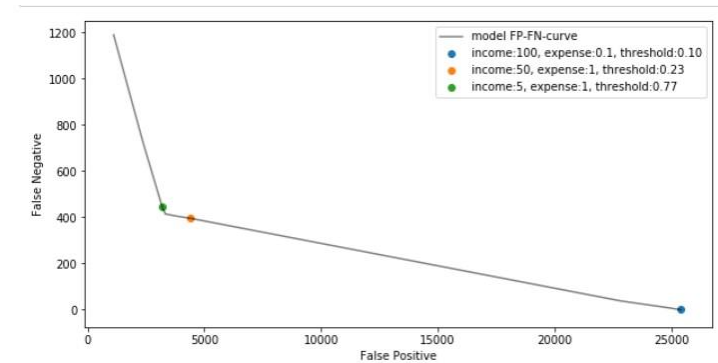
F1 = 0.92 по качеству прогноза для класса 0 – не подключение услуги абонентом.

AUC\_PR = 0.350

# Принцип составления индивидуальных предложений для выбранных абонентов

Предлагаемый принцип - максимизация получаемой прибыли при заданных значениях дохода от подключенной услуги и затраты на рассылку предложения.

**График количества  $NFN$  и  $NFP$  при разных порогах на предсказаниях модели Catboost**



*Прибыль оператора от положительного отклика клиента на услугу = Доход от клиента - Затраты на рассылку предложения этому клиенту*

- ошибка первого рода отражает доход, который оператор потерял, не отправив предложение.
- ошибка второго рода отражает затраты оператора на рассылку, которые оказались напрасными.

С учетом того что доход от подключаемой услуги выше затрат, а так же как правило абонент продолжает пользоваться услугой в последующем (продлевает услугу), то ошибка первого рода сильнее влияет на упущенную выгоду. При этом затраты на предложение с дальнейшим отказом от услуги несут только затраты на само предложение, но при высоких затратах могут привесить доход от лояльных клиентов.