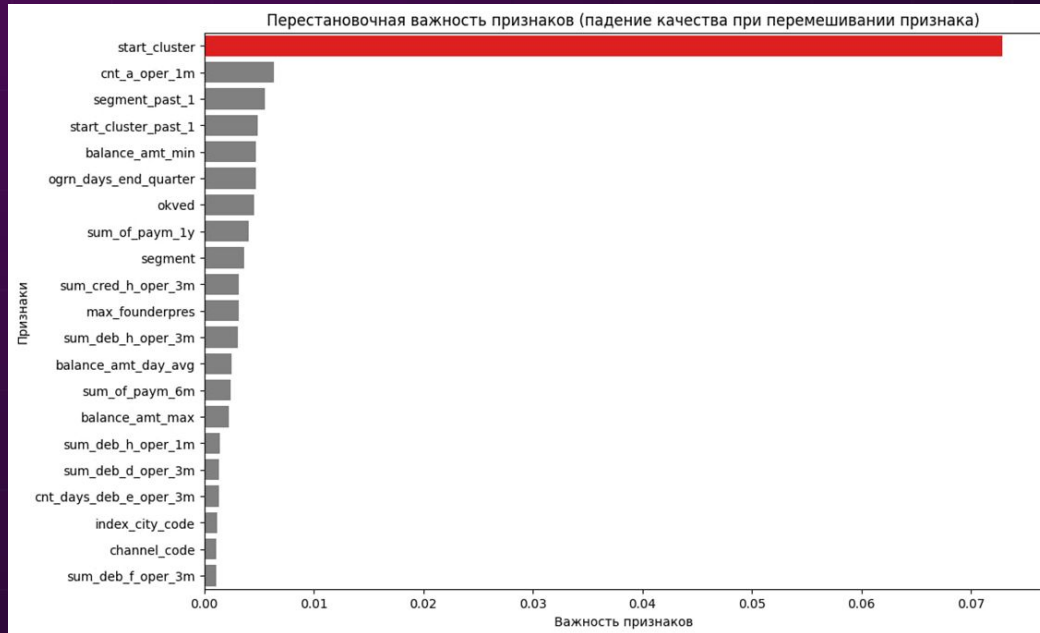


Решение бизнес-задач, связанных с CLTV

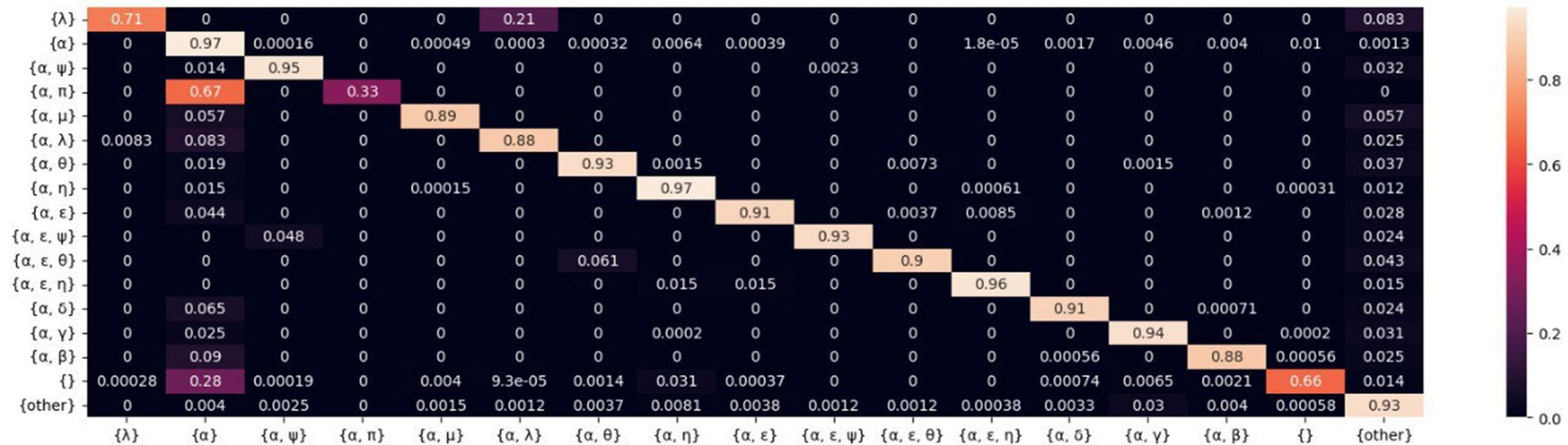
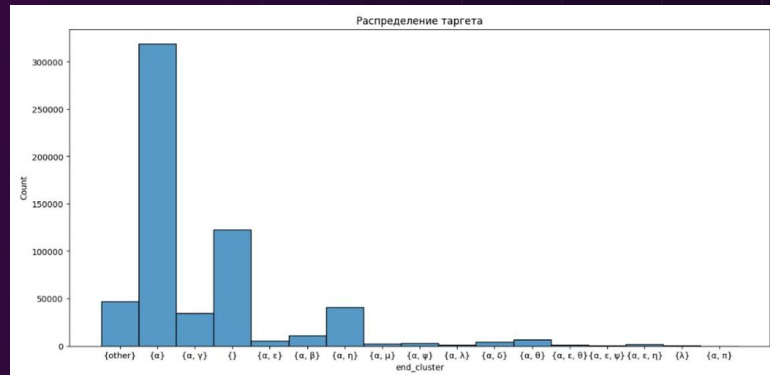
No Loss March

Start cluster – самый важный признак.

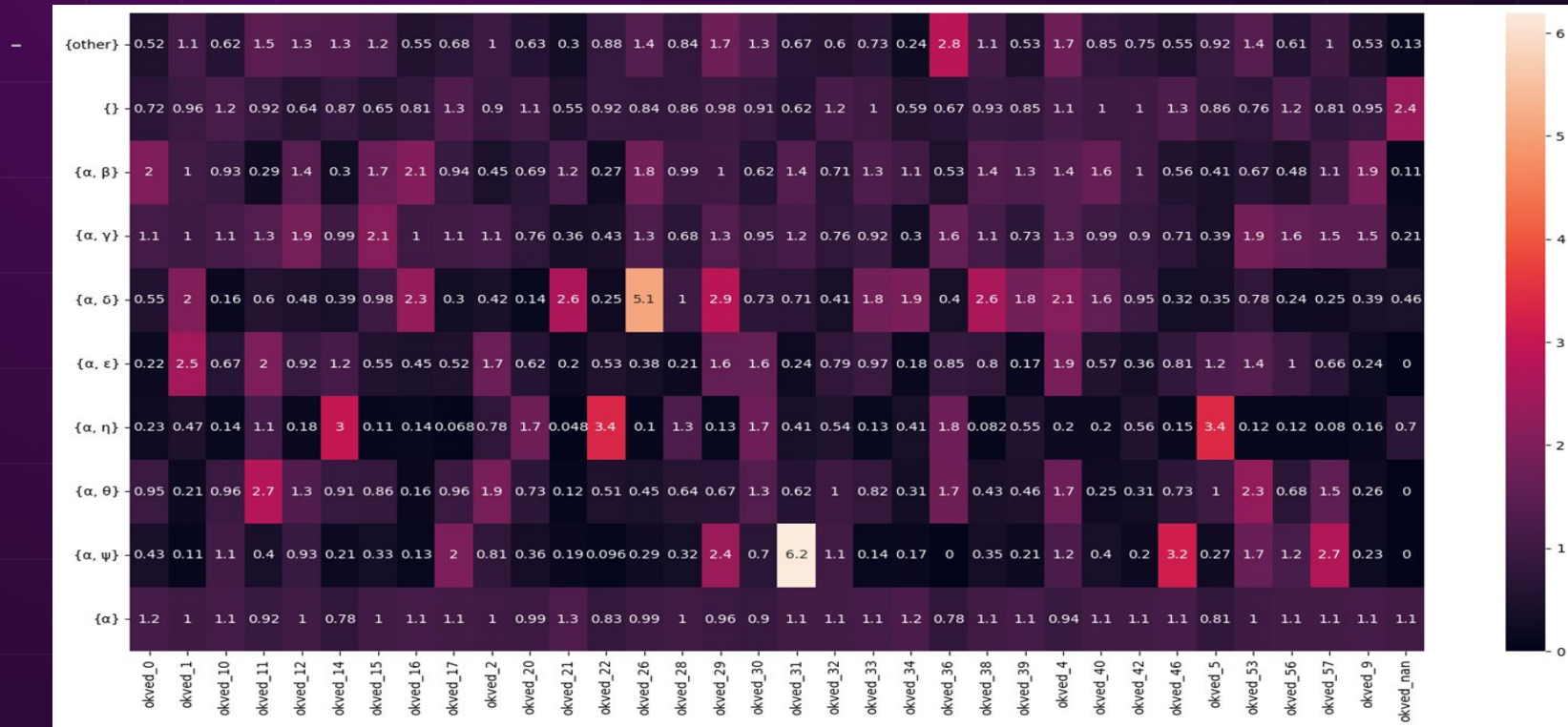
82% - точность константного baseline, основанного на Start_cluster



Частотный анализ Start cluster



Предпочтение различных ОКВЕД к кластеру

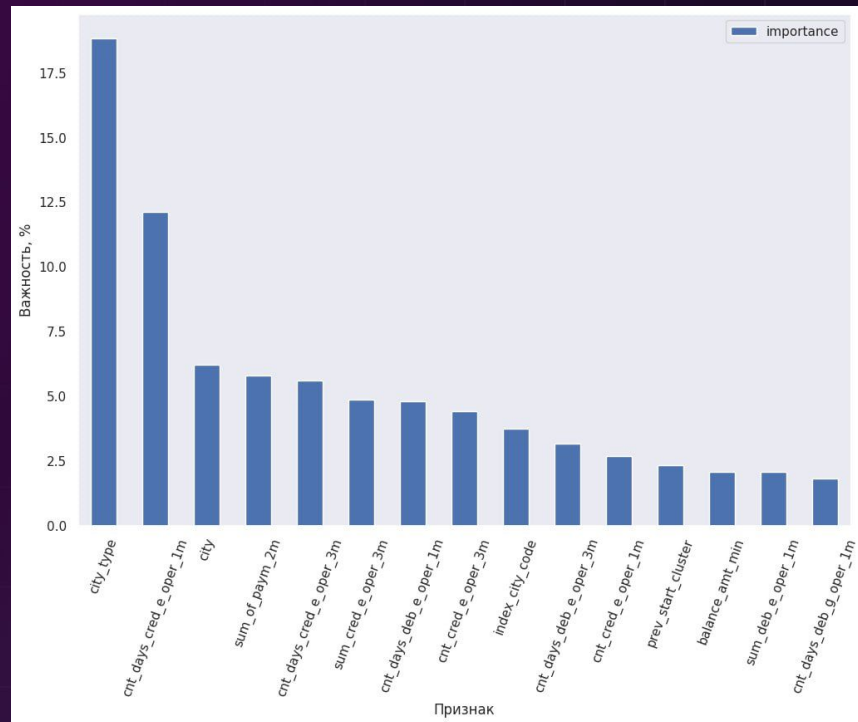


Другие интересные замечания из EDA

- У **новых** пользователей свой **ограниченный набор кластеров**
- Среднее время до получения ОГРН в трэйн ~ 0 , в тест $> 0 \Rightarrow$ **train идет после test** по времени
- Город, индекс города и тип города, **перекрывают пропуски** в данных друг у друга. Выбрали **оставить** все **несмотря на** кажущуюся **дублируемость** информации.
- В крупных городах России распределение категорий доходности отличается от малых.

Train vs Test

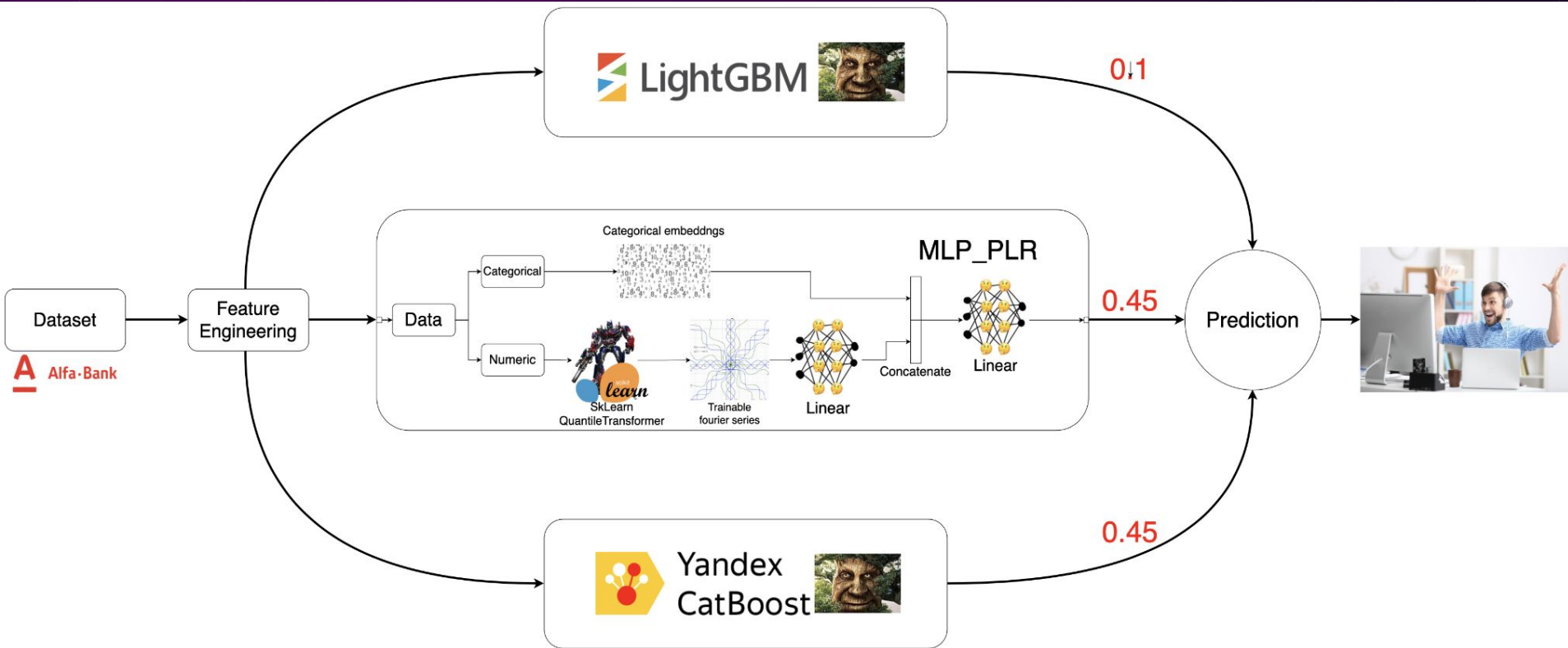
- Сильные различия в признаках указывающих на прошедшее время с момента.
- В тесте есть новые города
- Числовые признаки испытывают небольшой сдвиг по времени(Операция “Е” изменяется сильнее всего, возможно в купе с типом города)



Неудачные эксперименты

- Добавление истории. RNN.
- Tied Embeddings.
- Комбинирование признаков.
- Доменная адаптация(Предобучение)
- Добавление весов примерам

Модель



Конкурентные Преимущества.

Высокая точность



Легко встроить и поддерживать



Конкурентные Преимущества.

Легкая масштабируемость



- Легко применить к большому количеству данных
- Борьба с ковариационным сдвигом
- Извлечение пользы из неразмеченного набора данных

Над проектом работали



Дмитрий Харчев

ML-разработчик

@KharDim08



Тохчуков Данил

ML-разработчик

@makriot



Черемискин Егор

ML-разработчик

@he_is_already_here



Лапиков Владислав

ML-разработчик

@What_is_Love_iss



Кадченко Иван

Аналитик

@KadchenkoIE