

1. Вступление. Цель и значимость проекта для бизнеса.

Проект решает задачу прогнозирования оттока клиентов. Привлечение клиентов является непростой задачей. Компании, которые хотят удерживать устойчивые позиции на рынке среди конкурентов, готовы выделять значительный размер маркетингового бюджета на привлечение новых клиентов, повышение лояльности существующих клиентов и удержание клиентов, которые готовы прекратить сотрудничество с компанией. Уход клиентов к конкурентам значительно ослабляет позиции компании на рынке. Решение задачи прогнозирования оттока клиентов существенным образом помогает компании в рациональном распределении ресурсов и формировании оптимальной структуры бюджета.

2. Техническая часть.

- **Входные данные.**

Входные данные представлены 4 таблицами со столбцами:

personal.csv - данные о клиентах,

customerID - идентификатор клиента,
gender - пол клиента,
SeniorCitizen - наличие пенсионного статуса по возрасту,
Partner - наличие супруга(и),
Dependents - наличие иждивенцев,

contract.csv - данные о договорах,

BeginDate - дата начала пользования услугами,
EndDate - дата окончания пользования услугами,
Type - тип оплаты: ежемесячный/годовой и тп,
PaperlessBilling - безналичный расчет,
PaymentMethod - способ оплаты,
MonthlyCharges - ежемесячные траты,
TotalCharges - всего потрачено денег на услуги,

internet.csv - данные о пользовании интернетом,

InternetService - тип подключения: DSL или оптоволокно,
OnlineSecurity - наличие блокировки небезопасных сайтов,
DeviceProtection - наличие антивируса,
OnlineBackup - наличие облачного хранилища для резервного копирования данных,
TechSupport - наличие выделенной линии технической поддержки,
StreamingTV - наличие стримингового телевидения,
StreamingMovies - наличие каталога фильмов,

phone.csv - данные о пользовании телефонией,

MultipleLines - возможность подключения телефонного аппарата к нескольким линиям одновременно

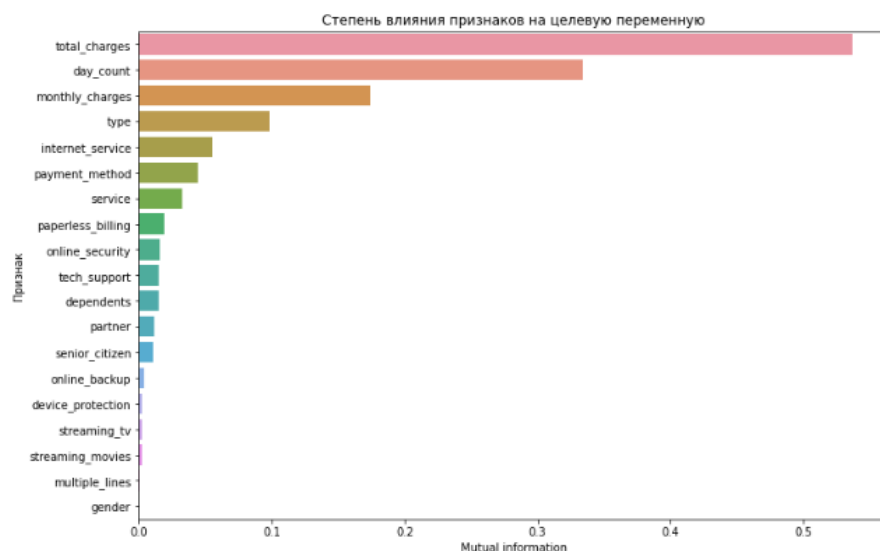
- **Предобработка данных и анализ.**

- ◆ Преобразование типа данных столбца TotalCharges из строкового в float64, пропуски по только что пришедшим клиентам заполнены значением из MonthlyCharges
- ◆ Данные в столбцах BeginDate и EndDate приведены к типу данных datetime64. Значения столбца EndDate по лояльным клиентам вместо 'no' заполнены датой выгрузки данных из базы (задается константой в разделе загрузки необходимых библиотек и модулей)
- ◆ В столбцах строкового типа данных устранены неявные дубликаты путем приведения данных к нижнему регистру с удалением лишних пробелов
- ◆ Данные в столбце SeniorCitizen заполнены 'yes/no' для унификации отображения бинарных признаков.
- ◆ Выделены новые признаки:
 - целевая переменная churn – признак ухода клиента (0/1) из данных столбца EndDate,
 - признак day_count – признак количества дней сотрудничества клиента (разница EndDate и BeginDate),
 - признак service – тип предоставляемых услуг (internet/phone/both),
 - service_count – количество предоставляемых опций на услугах интернета и телефонии
- ◆ Таблицы объединены по столбцу customerID. Если клиент не пользуется услугой телефонии или интернета по соответствующим признакам из таблиц internet и phone проставляется значение 'no'.
- ◆ Исследовательский анализ включает
 - Проверку данных на наличие аномалий: boxplot (количественных переменных) и value_counts() (категориальных переменных)
 - Проверку границ BeginDate и выгрузки данных по уходам EndDate
 - Анализ распределений количественных переменных TotalCharges, MonthlyCharges и day_count в разрезе данных churn
 - Проведена оценка влияния параметров категориальных признаков (по большей части предоставляемых опций на услугах интернета и телефонии) на увеличение/снижение вероятности ухода клиента
 - Составлен “портрет” склонного к оттоку клиента
- **Построение модели и анализ.**

Чтобы избежать утечки данных удалим признаки customerID, BeginDate, EndDate, сгенерированные для анализа признаки begin_month и end_month.

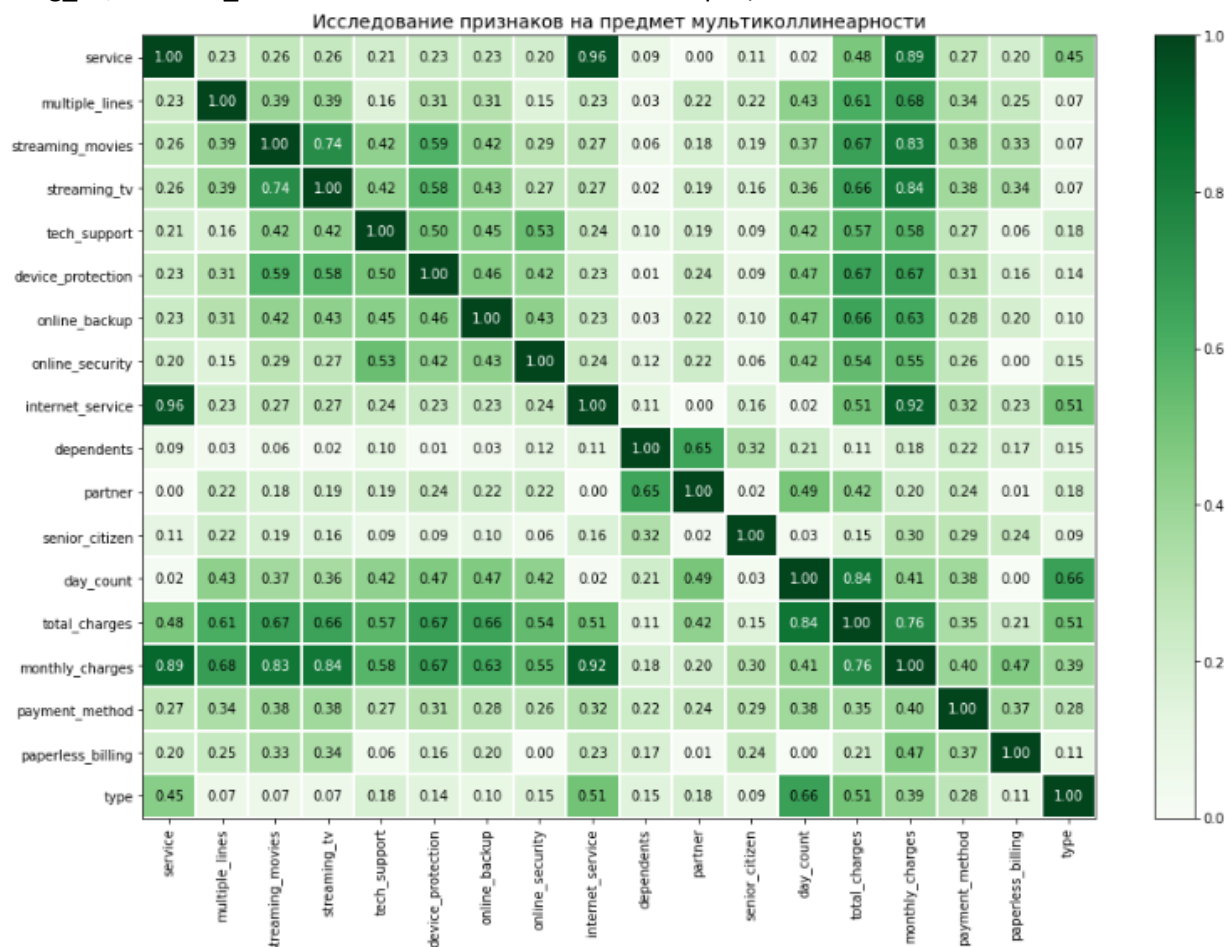
Проверка на наличие полных дубликатов выявила наличие таковых. Удалили.

Проведена оценка степени влияния каждого признака на целевую переменную по Mutual Information



Для обучения оставлены признаки 'type', 'paperless_billing', 'gender', 'payment_method', 'monthly_charges', 'total_charges', 'day_count', 'senior_citizen', 'partner', 'dependents', 'internet_service', 'online_security', 'online_backup', 'device_protection', 'tech_support', 'streaming_tv', 'streaming_movies', 'multiple_lines', 'service'.

Анализ мультиколлинеарности показал значительную корреляцию между признаками 'monthly_charges', 'streaming_tv', 'internet_service': затемненные области на матрице из библиотеки Phik.



Произведено удаление коррелирующих признаков. Корреляция по оставшимся признакам не достигает 70%.

Произведен анализ на предмет сбалансированности целевого показателя. Наблюдается дисбаланс классов 0 к класс - 73%, 1 класс – 27%. Произведен расчет классов методом compute_class_weight с целью дальнейшей передачи в модели при обучении.

Тк планируется составление пайплайна со встроенным кроссвалидатором, то разбиение итоговых данных произведено на обучающую и тестовую выборку с долей тестовых данных 25%, со стратификацией по целевому показателю, перемешиванием shuffle=True и random_state=270323 для воспроизводимости результата.

В рамках проекта был построен ряд Pipeline с перебором параметров с помощью RandomizedSearchCV для моделей LogisticRegression, RandomForestClassifier, LGBMClassifier, CatBoostClassifier.

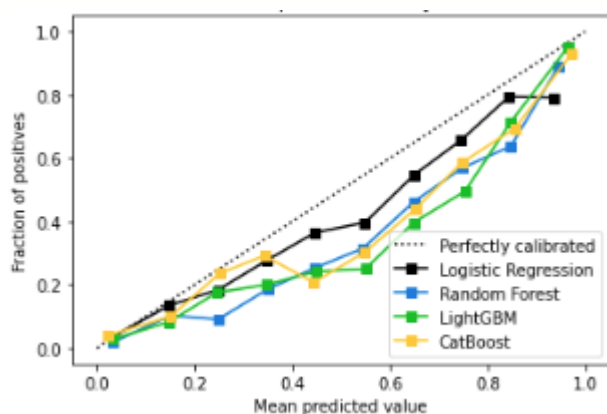
Для линейных моделей было проведено масштабирование численных признаков методом StandardScaler, кодирование категориальных признаков методом OneHotEncoder(drop='first'). Борьба с дисбалансом производилась встраиванием метода SMOTE в пайплан, передачей параметра class_weights='balanced' в модель или непосредственно заранее рассчитанных долей классов методом compute_class_weight. Для бустинговых моделей производился подбор оптимального шага обучения. Для деревянной модели случайного леса был произведен подбор количества деревьев в лесу, глубины дерева, критерия gini или entropy. В RandomizedSearchCV передавались пайплайны с перечнем перебираемых параметров с разбиением обучающего датасета на 5 частей для кроссвалидации, ограничением итераций до 100. Максимизируемая метрика ROC-AUC.

В результате обучения получены следующие величины метрик по разным моделям (см. табл)

	Method	Roc_auc	Время обучения, сек
0	Pipeline+RSCV+LogReg	0.84	39
1	Pipeline+RSCV+RandForest	0.85	1.57
2	Pipeline+RSCV+LightGBM	0.89	11
3	Pipeline+RSCV+CatBoost	0.91	1.22

Проведено исследование надежности модели (устойчивости на тестовых данных) с помощью plot_calibration_curve().

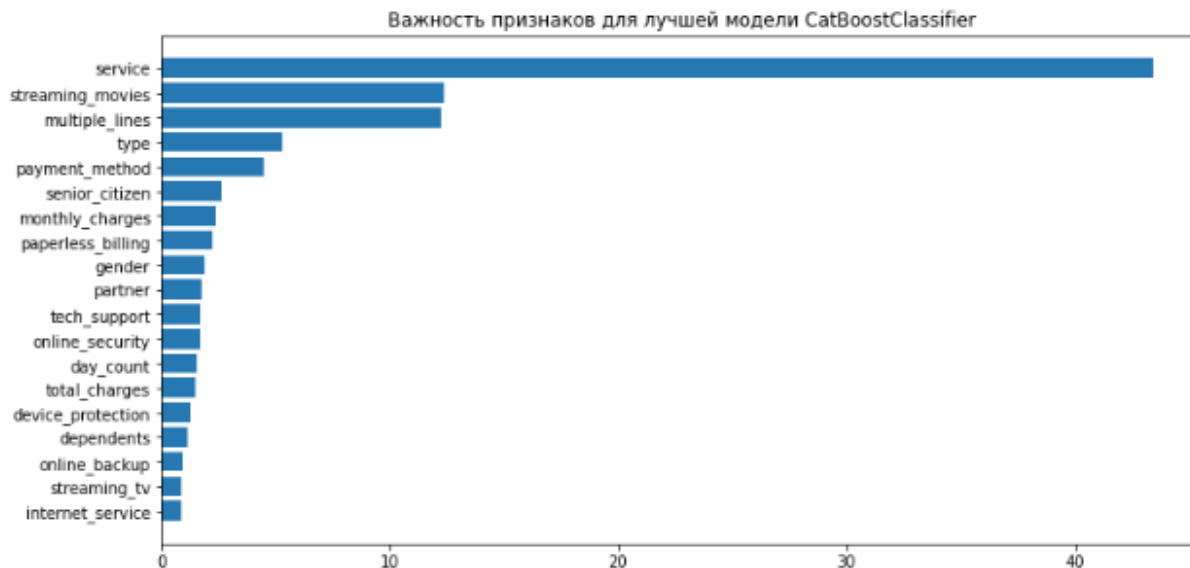
Степень калибровки построенных моделей



По степени калибровки лучшей является модель Logistic Regression, затем идет CatBoost и Random Forest. Учитывая величину метрики, можно утверждать, что CatBoost является лучшей моделью.

Лучшей моделью оказалась бустинговая модель CatBoostClassifier с шагом обучения 0,1 и весами для классов {0: 0.68, 1: 1.88}, встроенная в pipeline с кодированием категориальных переменных методом OrdinalEncoder. Метрика ROC_AUC составила как на обучающей, так и на тестовой выборке 91%.

Рассмотрим, на какие признаки в основном опиралась модель при расчете прогноза вероятности ухода клиента.

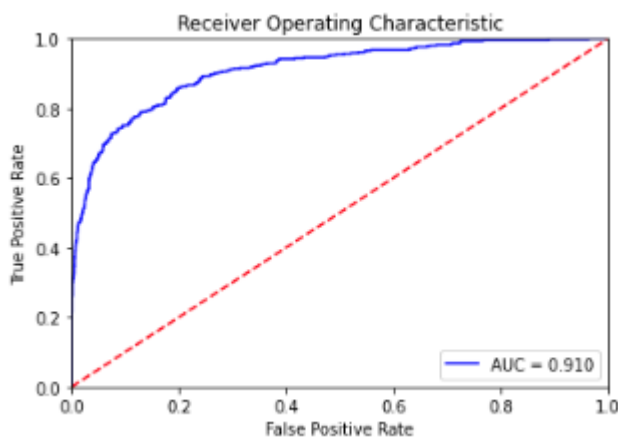


Наиболее важный признак для прогнозирования модели оказался service, затем в равной степени streaming_movies и multiple_lines (который ранее был практически незначимый по матрице Phik), в еще меньшей - type и payment_method и далее все остальные.

- Тестирование.

3. Анализ результатов работы модели. Бизнес-анализ альтернативных издержек. Рекомендации руководству компании.

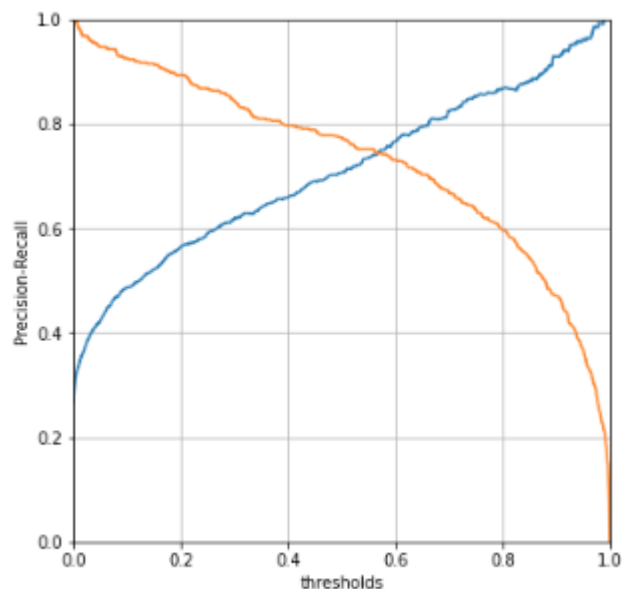
Для анализа соотношения True Positive Rate и False Positive Rate была построена ROC-кривая



Произведен расчет precision, recall, F-меры и ROC-AUC на разных порогах классификации.

Порог	Точность	Полнота	F1	ROC-Auc
0.28	0.608	0.859	0.712	0.829
0.30	0.618	0.842	0.713	0.827
0.32	0.629	0.829	0.715	0.826
0.34	0.634	0.812	0.712	0.821
0.36	0.647	0.807	0.718	0.824
0.38	0.656	0.805	0.723	0.826
0.40	0.660	0.797	0.722	0.824
0.42	0.669	0.792	0.725	0.825
0.44	0.685	0.788	0.733	0.829
0.46	0.692	0.779	0.733	0.827
0.48	0.702	0.777	0.738	0.829
0.50	0.706	0.773	0.738	0.829
0.52	0.717	0.760	0.738	0.826
0.54	0.733	0.752	0.742	0.826
0.56	0.740	0.749	0.745	0.827
0.58	0.751	0.741	0.746	0.826
0.60	0.766	0.730	0.748	0.825
0.62	0.779	0.724	0.750	0.825
0.64	0.786	0.709	0.745	0.820
0.66	0.796	0.700	0.745	0.818

Построена кривая Precision-Recall с расчетом наилучшего порога классификации.



Best threshold = 0.568

При пороге 0,568 наблюдаем оптимальное соотношение параметров: высокая метрика precision (0.74) при сохранении на должном уровне recall (0.749) и остальных метрик: на лчень хорошем уровне F1 = 74.5%, ROC-Auc = 82.7%. На этом пороге модель максимизирует долю истинных положительных и долю истинных отрицательных результатов, минимизируя вероятность ложной тревоги - ошибки первого рода. На этом пороге показатели точности и полноты 0.610 и 0.627 соответственно. На дальнейших порогах точность продолжает расти, остальные метрики начинают постепенно падать.

С другой стороны в данной задаче определения вероятности ухода клиента более важно максимизировать метрику precision и ложно отправить скидку или письмо клиенту, который на самом деле не хотел уходить из компании, чтобы оставить его. Чем не предоставить скидку и потерять клиента. Здесь нужно смотреть выручку, которую приносит клиент (определить категорию клиента по типу тарифа), возможный размер предоставляемой скидки и пр. Эти вопросы нужно рассматривать исходя из требований заказчика.

Далее произведен расчет вероятности ухода для каждого клиента из тестового набора данных и выведен список клиентов, отнесенных моделью к 1 классу.

```
probabilities_test
```

```
array([[0.93349766, 0.06650234],  
       [0.84105922, 0.15894078],  
       [0.99756652, 0.00243348],  
       ...,  
       [0.06210931, 0.93789069],  
       [0.41805488, 0.58194512],  
       [0.41627428, 0.58372572]])
```

На основании величины вероятности ухода клиента, можно сгруппировать «отточных» клиентов по величине вероятности и каждой группе предложить различные системы скидок. На основе этих данных можно провести оптимизацию маркетинговой политики, предлагаемых клиенту продуктов и условий сотрудничества.

Имея выгрузку данных по клиентам, которые склонны к оттоку, можно уже рассчитать альтернативные издержки на предоставление скидки и упущенную выгоду от ухода клиента с проведением соответствующего анализа и составлением прогнозов для компании в обоих случаях. Таким образом, можно добиться сокращения неоптимальных расходов компании и целенаправленной точечной работой с отдельными группами клиентов.

4. Дальнейшие пути развития проекта.

- Хотела, но не реализовала пока: выбор лучшей модели осуществлять не только на основании величины метрики и калибровки модели, но также исходя из принципов минимизации разброса возможных вариантов метрики вокруг среднего генеральной совокупности (мы же не знаем, какие тестовые данные будут предоставлены), для этого нужно определить доверительные интервалы для метрик всех построенных моделей.
- Рассчитать альтернативные издержки бизнеса с учетом предлагаемых скидок для удержания клиента. Составить прогноз выручки и чистой прибыли компании в обоих случаях.
- Провести развертывание построенной модели машинного обучения в рабочей среде в качестве API с помощью Flask.