

KKBox Churn Prediction

Springboard Data Science Career Track
Capstone Project 3

Marco Riva
August 10, 2021

Client

- KKBox: Asia's leading music streaming service
- Business model: Subscription model
- <https://www.kkbox.com/>



Motivation & Business Impact

- Churn is of great concern for subscription model businesses
- Acquisition cost generally higher than retention cost
- Understanding customer behavior is key to drive value and growth
- Predictive tools for churn prediction are fundamental to maximizing customer retention



Problem Formulation

- Supervised Learning
- Binary Classification: Churn VS Retention
- Imbalanced dataset ~5% Churn - ~95% Retention

Objectives

- Utilize historical data provided by KKBBox regarding user demographics, transaction and logs to predict customer churn in March 2017 (test set)
- Identify features correlated with churn and provide actionable insights

Evaluation Metric

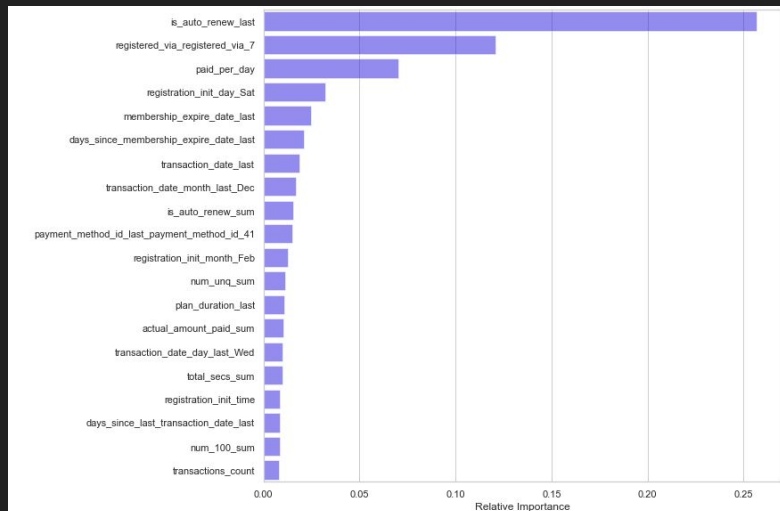
- Log Loss

Data Source

- <https://www.kaggle.com/c/kkbox-churn-prediction-challenge/data>

High Level summary and Key Findings

- User demographic, transactional and logs data were processed to generate predictors (features) for each user
- CatBoost, XGBoost and LightGBM models trained for churn prediction
 - Performance ~50% better than chance, Log Loss = 0.11 on test set (Kaggle)
- Key factors affecting churn:
 - Auto Renew Enabled
 - Price (averaged over day)
 - Registration Method
 - Temporal Features / Seasonality
 - Canceled Transactions
 - Usage: number of songs / time played



This presentation includes details of my work

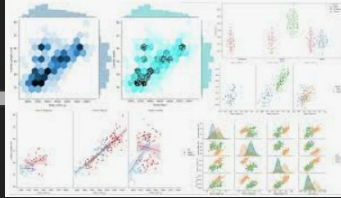
Data
Acquisition



Data
Wrangling



Exploratory
Data Analysis



Feature
Engineering



Modeling



Insights



Data

- Member info
~5M records
- Transactions
~22M records
- Logs
~400M records

Outliers

- Errors
- High Cardinality

Clean Data
Exploration

Processing
Aggregation
Encodings

Binary
Classification

What we've
learned

Available Data

- members.csv: User information
- transactions.csv: User transaction until end of February 2017
- transactions_v2.csv: User transaction in March 2017
- user_logs.csv: User logs until end of February 2017
- user_logs_v2.csv: User logs in March 2017
- train.csv: User ids and target variable for January 2017
- train_v2.csv: User IDs and target variable for February 2017
- submission.csv: User IDs and target variable for March 2017 (test set)

Data Wrangling

What Features?

- Member Info
- Transactions for each member
- Logs for each member

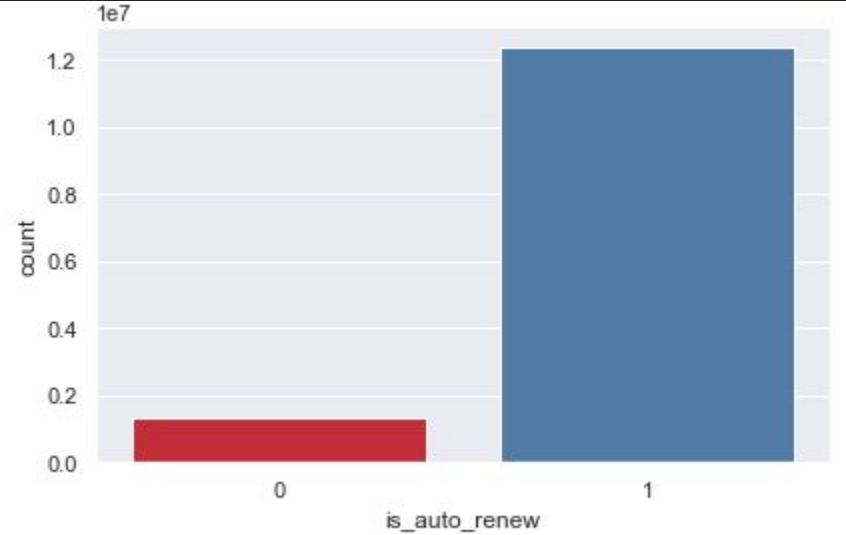
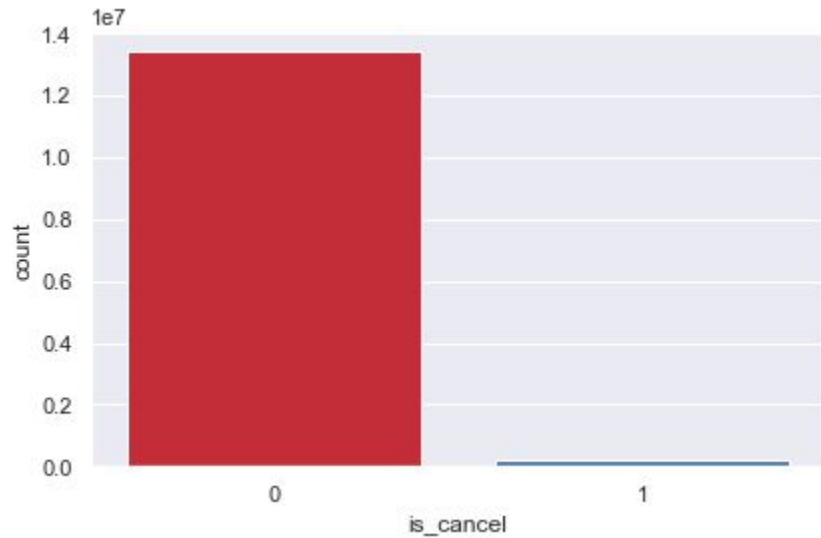
What Target?

- Imbalanced binary variable
- Positive class implies churn

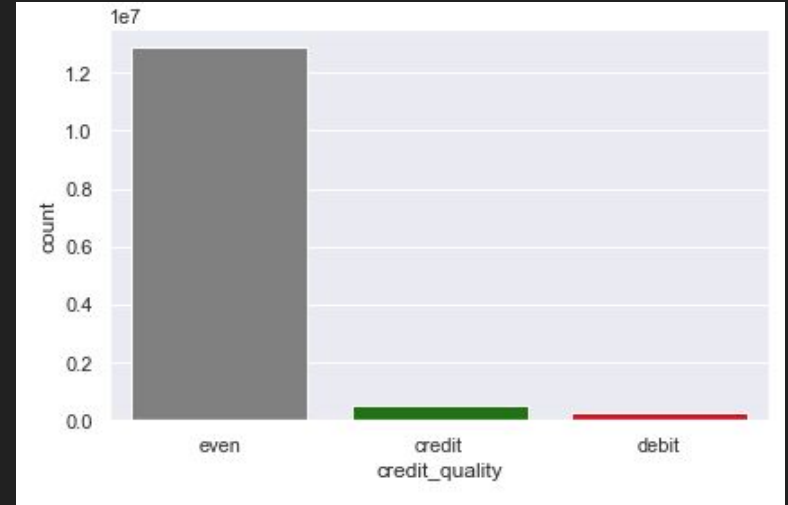
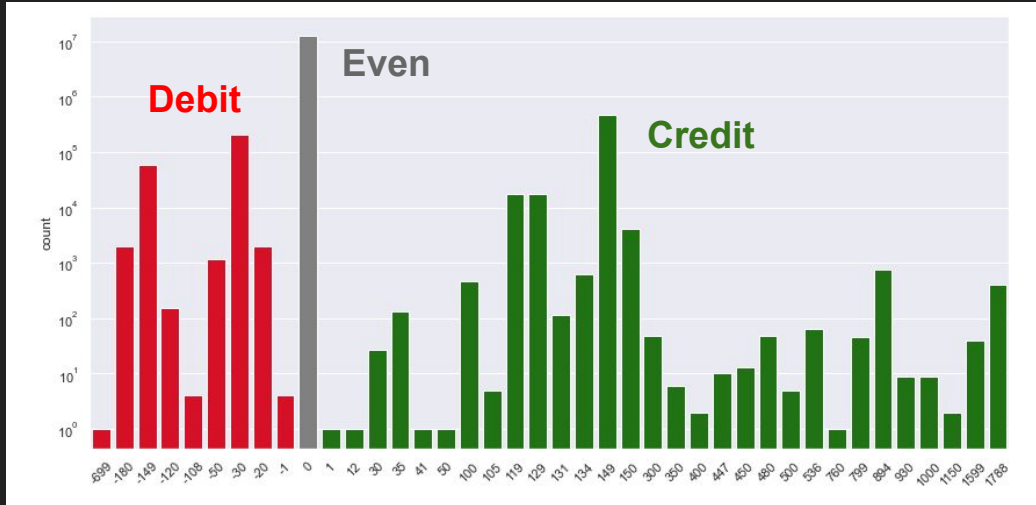
Data Cleaning

- Removed obvious outliers
- Dropped features with more than 40% missing values
- Reduced levels of high cardinality features
- No duplicates founds

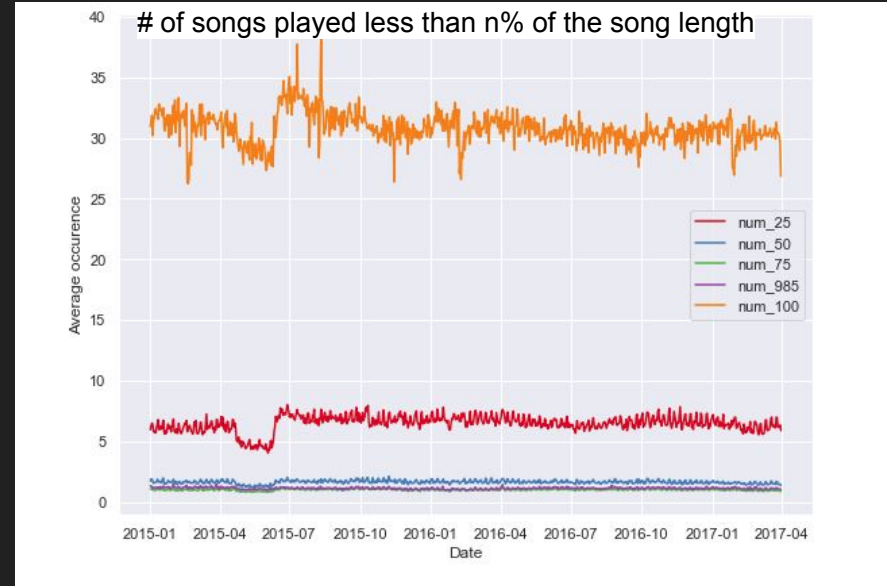
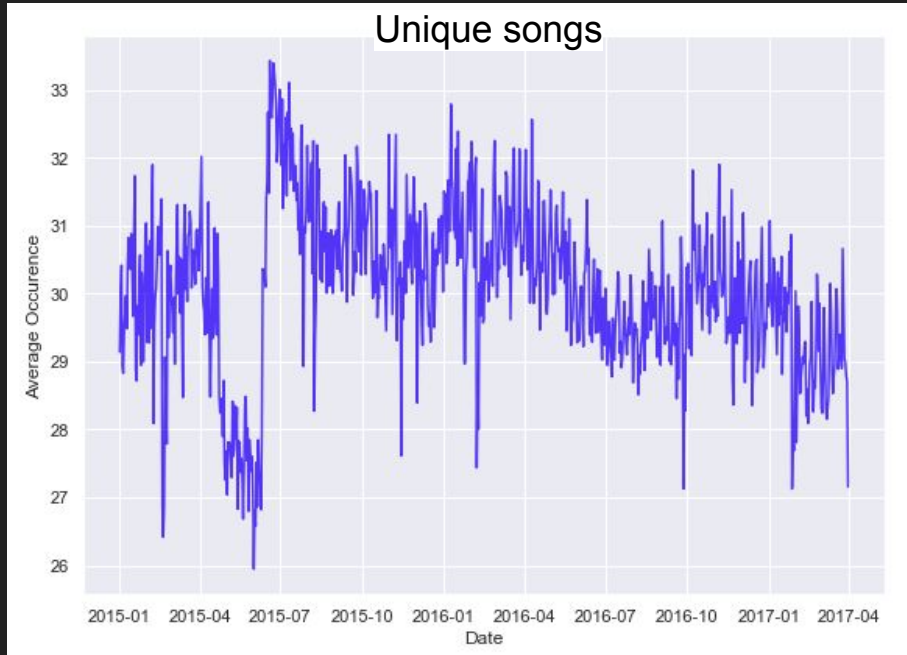
Exploratory Data Analysis



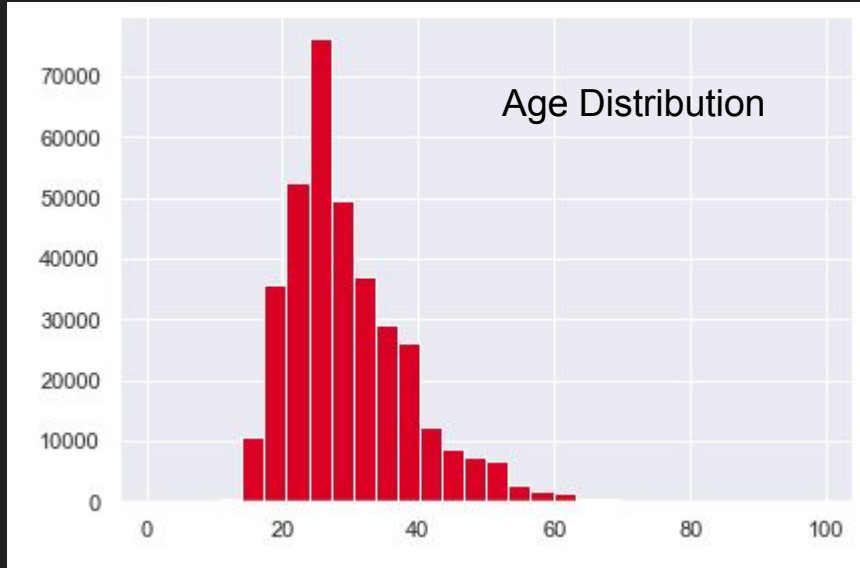
Exploratory Data Analysis



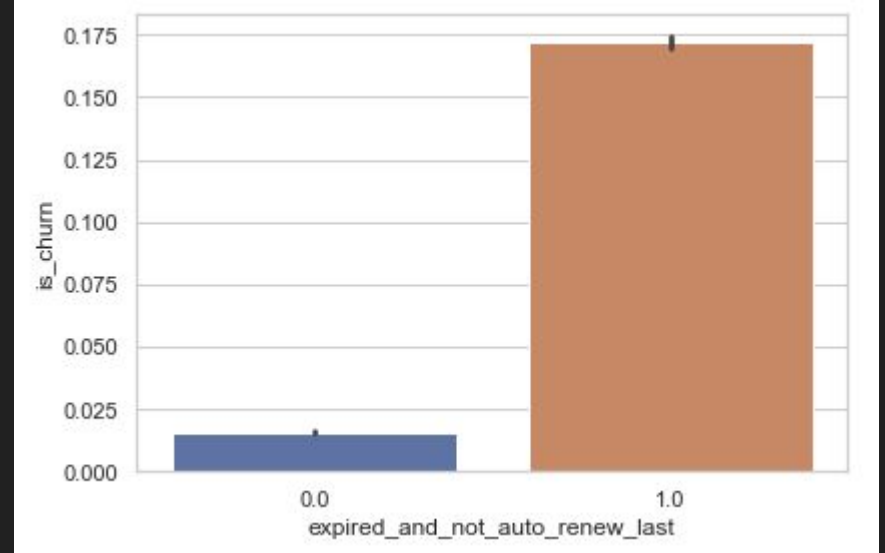
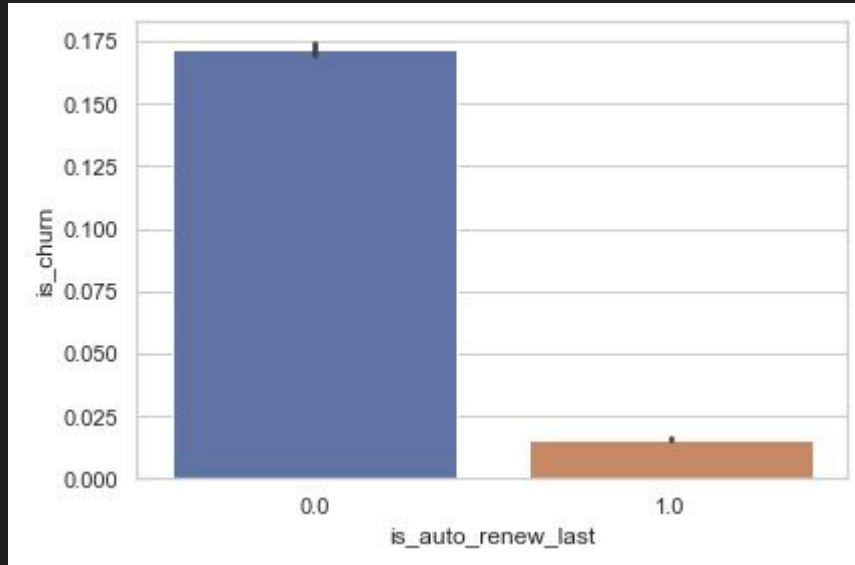
Exploratory Data Analysis

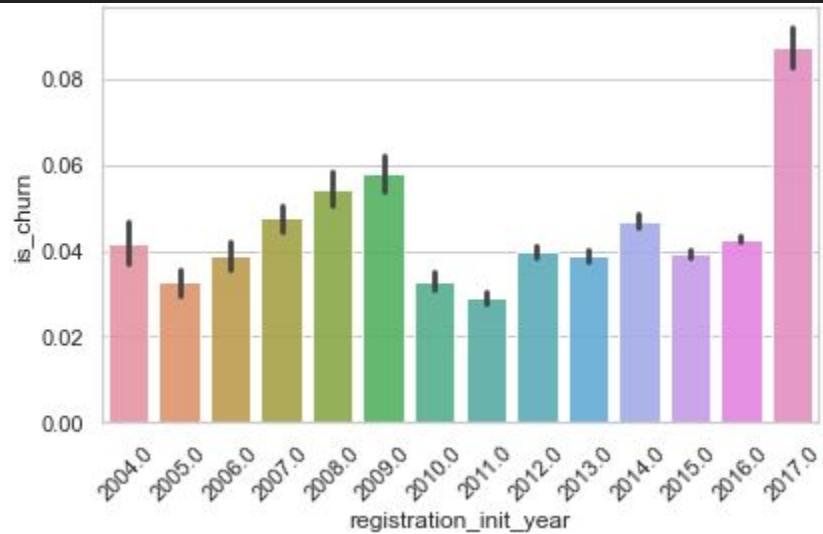
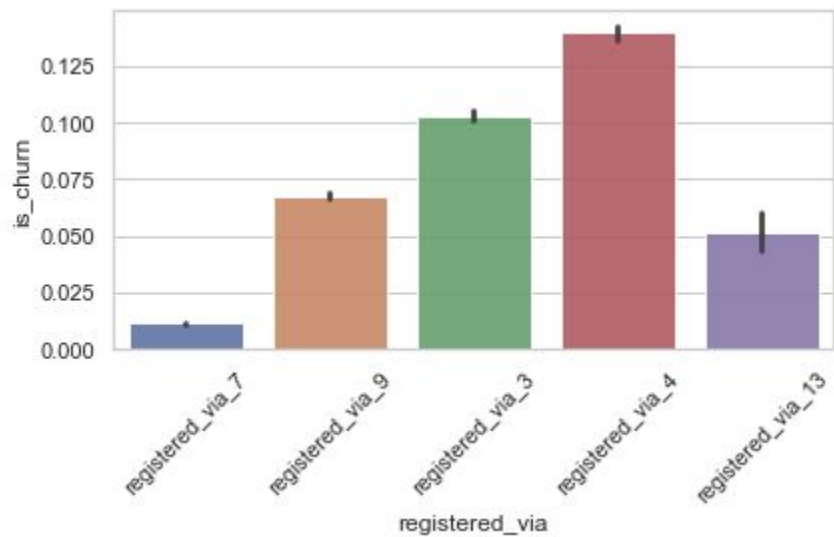


Exploratory Data Analysis

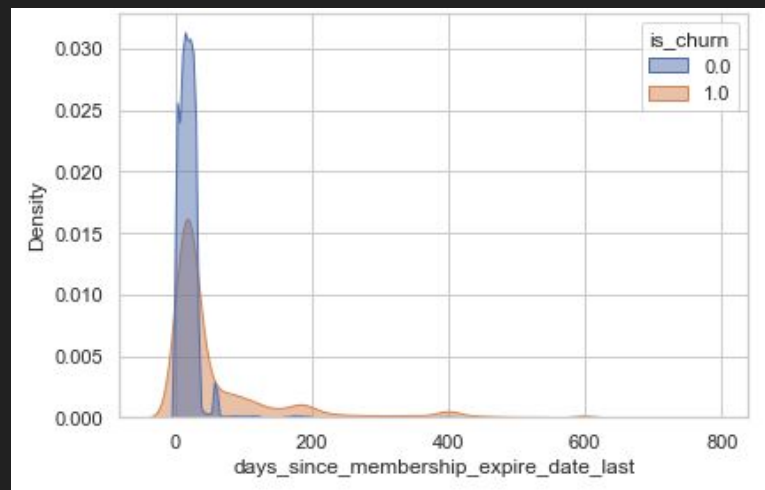
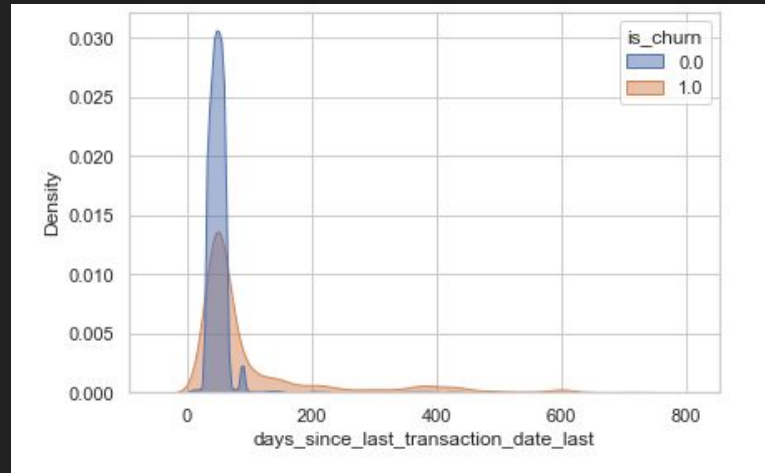
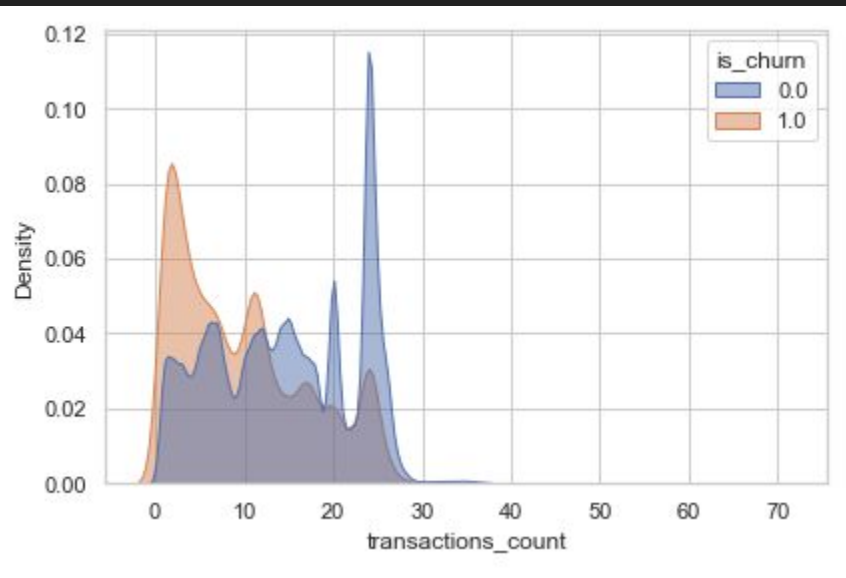


Candidate Features



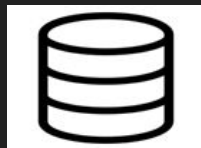
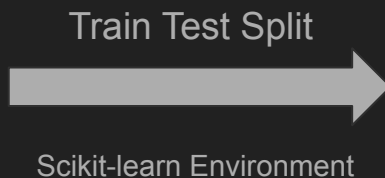


Candidate Features

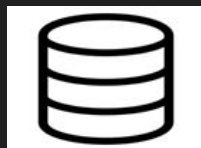


Feature Engineering

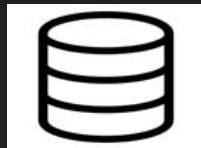
- Aggregated transaction and logs data to generate summary features (sum, average, count, etc.)



Train Set - Churn in February



Validation Set - Churn in March



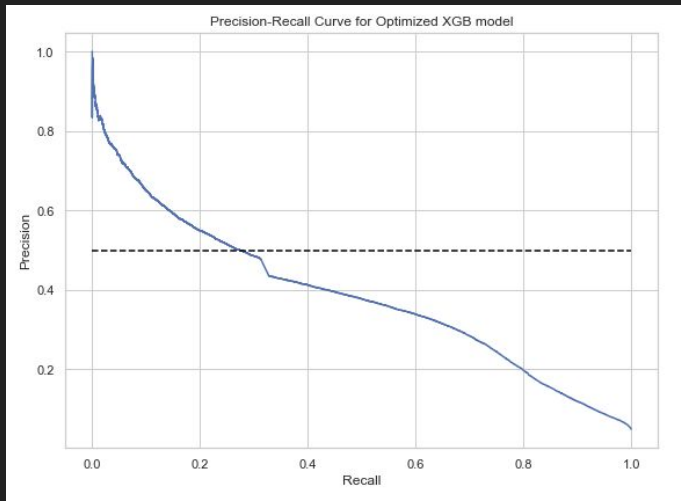
Test Set - Churn in April

Preprocessing Pipeline

- Missing values imputation
- One Hot Encoding
- 133 Features Total

Modeling

- Churn Rate 5%: dummy classifier Log Loss 0.19
- Boosting Models
 - CatBoost
 - XGBoost
 - LightGBM
- Tuned XGBoost provided best Log Loss of 0.11 on Test Set



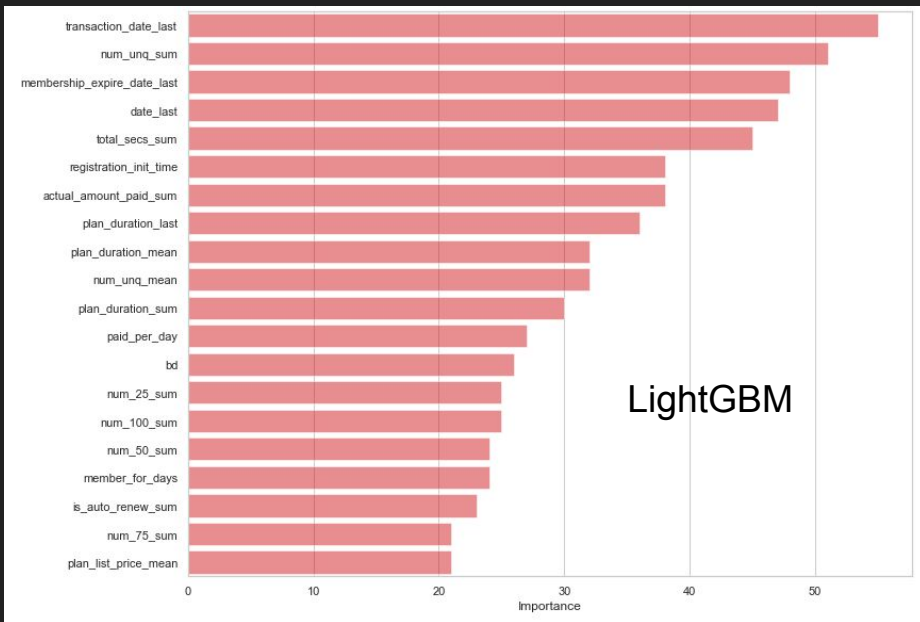
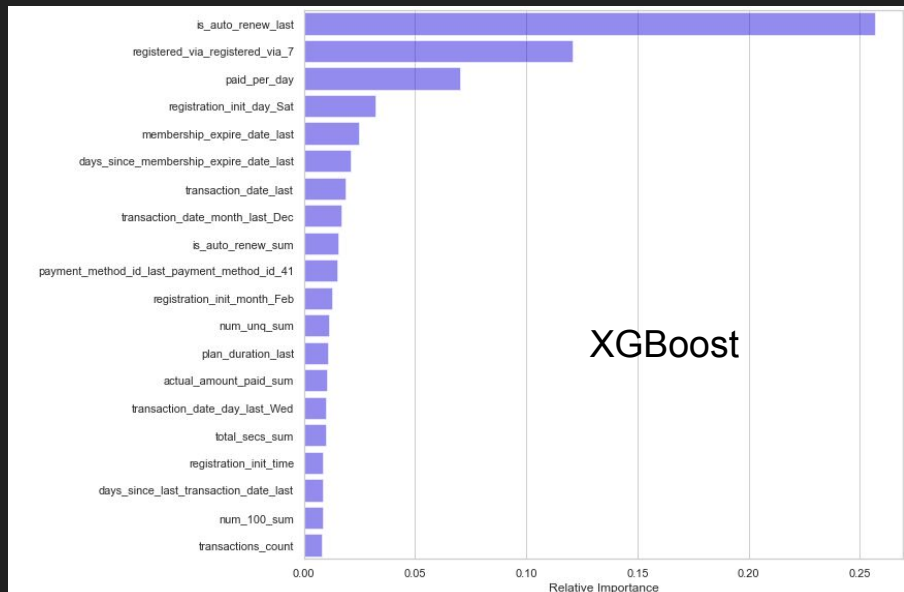
CatBoost

XGBoost



LightGBM

Feature Importance



Future Work

- Adjust sample weights to account for potential differences in feature / target distribution between train/valid/test sets
- Apply Oversampling techniques and/or SMOTE to improve predictions on minority class
- Adopt parallel computing solutions to investigate a wider range of hyperparameters

Thank you!