

Metrics:

The main metrics is log space rmse because the revenue is right-skewed and the revenue consists of a lot of zero values. I have got ~25% off in revenue predictions for positive spenders.

MAE = 0.69 means most of the time the prediction is near zero correct.

For R2 I have got 0.4753. It shows how that 47 per cent of the variability the model is able to explain.

RMSE = 20.22. it is high because of big spenders where the error can be hundred of dollars.

Pinball Loss (focuses on different parts of the distribution):

Median loss = 0.3448 shows that average excess is ~0.34

Top K Spenders

We aim to identify the big spenders. If we target the top 1% of users by the predicted spend, we will successfully include about 68% of the actual top-spending customers.

Limitations and Future improvements

Limitations:

- we use only Day 1 actions;
- lack of platform context, information about refunds;
- zero-inflation handled simplistically.
- Single monthly hold-out
- no hyperparameter tuning;

data improvements:

- enrich user actions with 7 days activity. we may identify users who start slow but is willing to pay after few days of usage.
- find out what the platform is and what it sells.
- add refunds/cancellations;
- add info if there was a campaign to reach out users;
- try different encoding techniques to handle string columns;
- check for more time related features

Modelling:

I think there are more approaches worth trying:

two stage model- predict for probability of payment + predicted spend;

try other boosting models;

handle the target distribution with zero values in another way;

Hyperparameter tuning;

ensemble models;

try k fold validation splitted by months (to check if the distribution amongst months is different);

Ensemble methods.