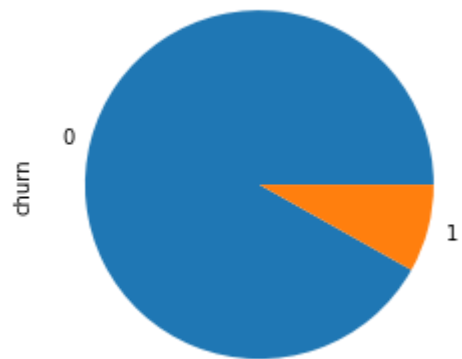


Telecom churn

Analyzing the churn Data

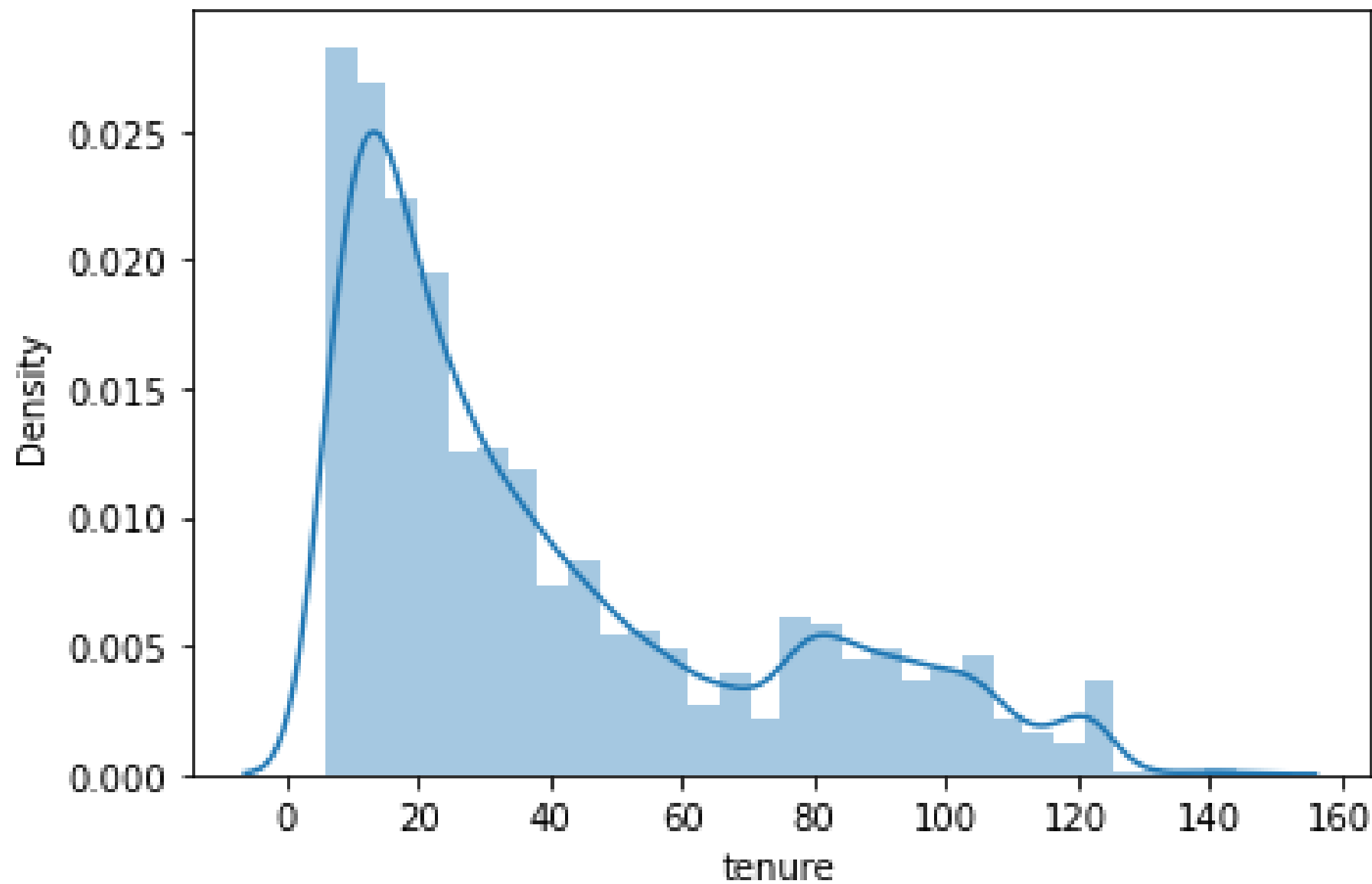
Content

1. Statistical Analysis of the Data
2. Handling missing values
3. Handling the other attributes with higher missing value percentage
4. Filtering the High Value Customer from Good Phase
5. *Defining Churn variable*
6. Deriving new variables to understand the data
7. Precision and recall trade-off
8. Metrics Evaluation



As we can see that 91% of the customers do not churn, there is a possibility of class imbalance

0 2-5 Yrs 7 2-5 Yrs 8 6-12 Months 21 1-2 Yrs 23 1-2 Yrs



0

2-5 Yrs

7

2-5 Yrs

8

6-12 Months

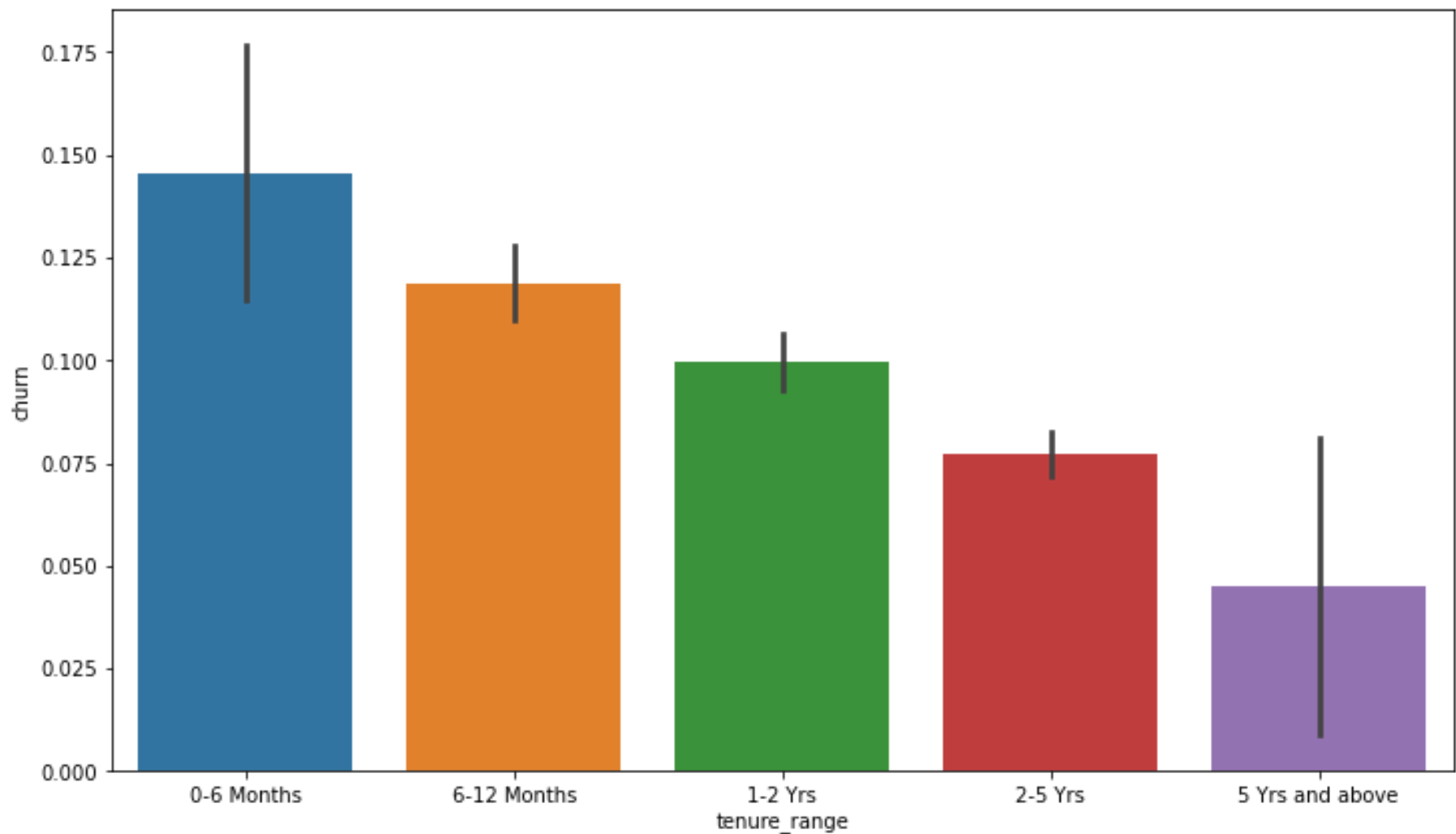
21

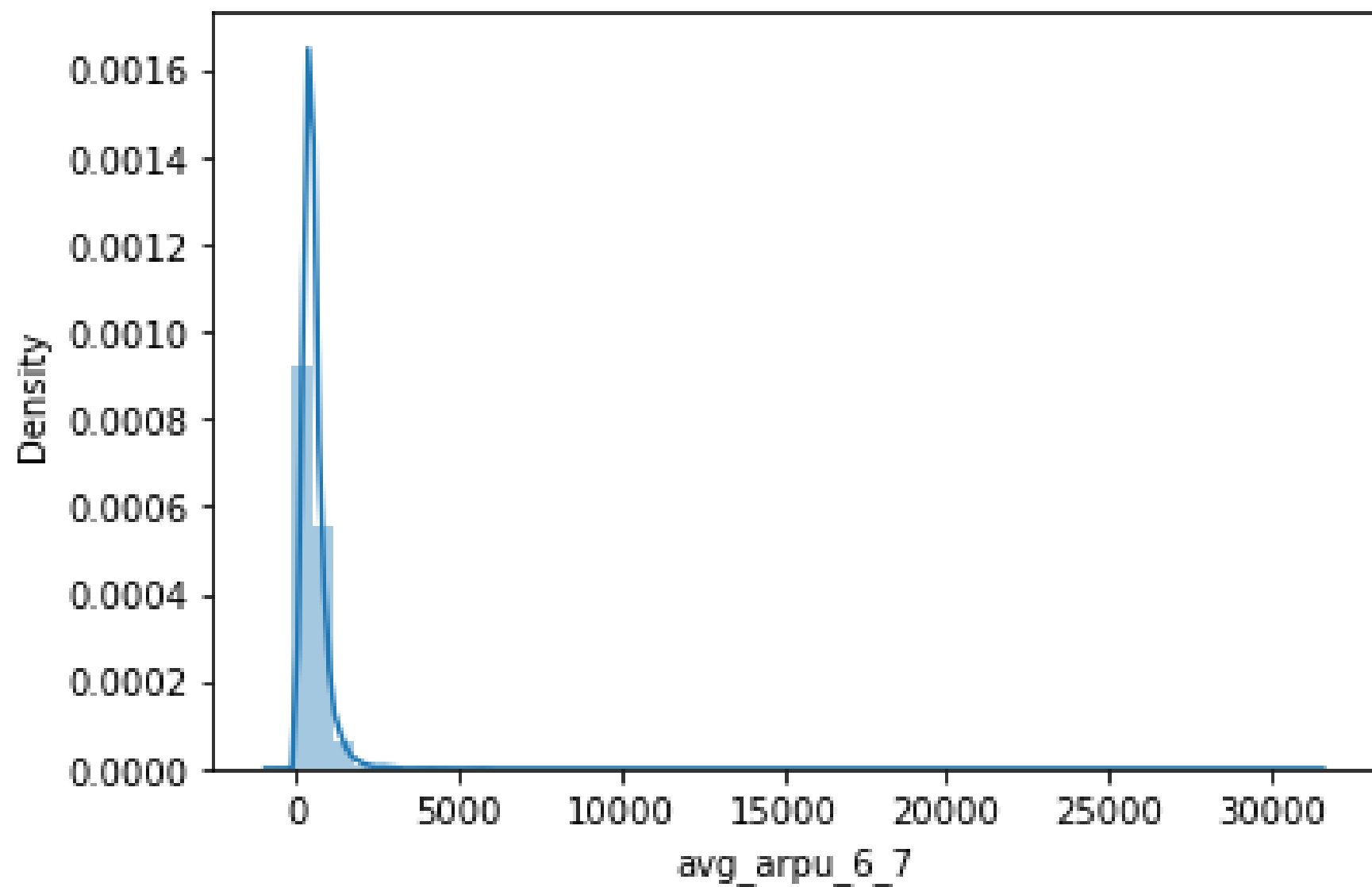
1-2 Yrs

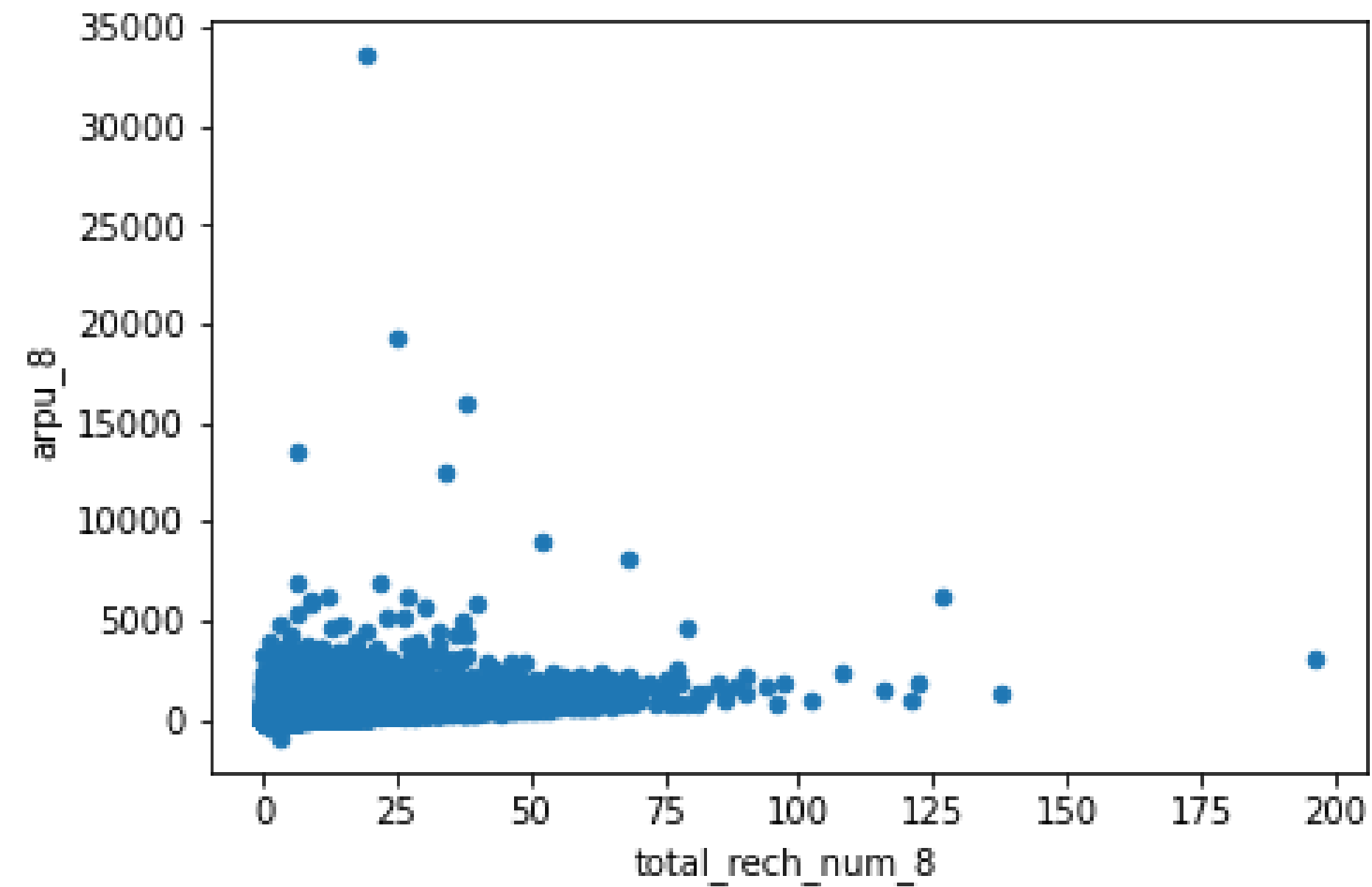
23

1-2 Yrs

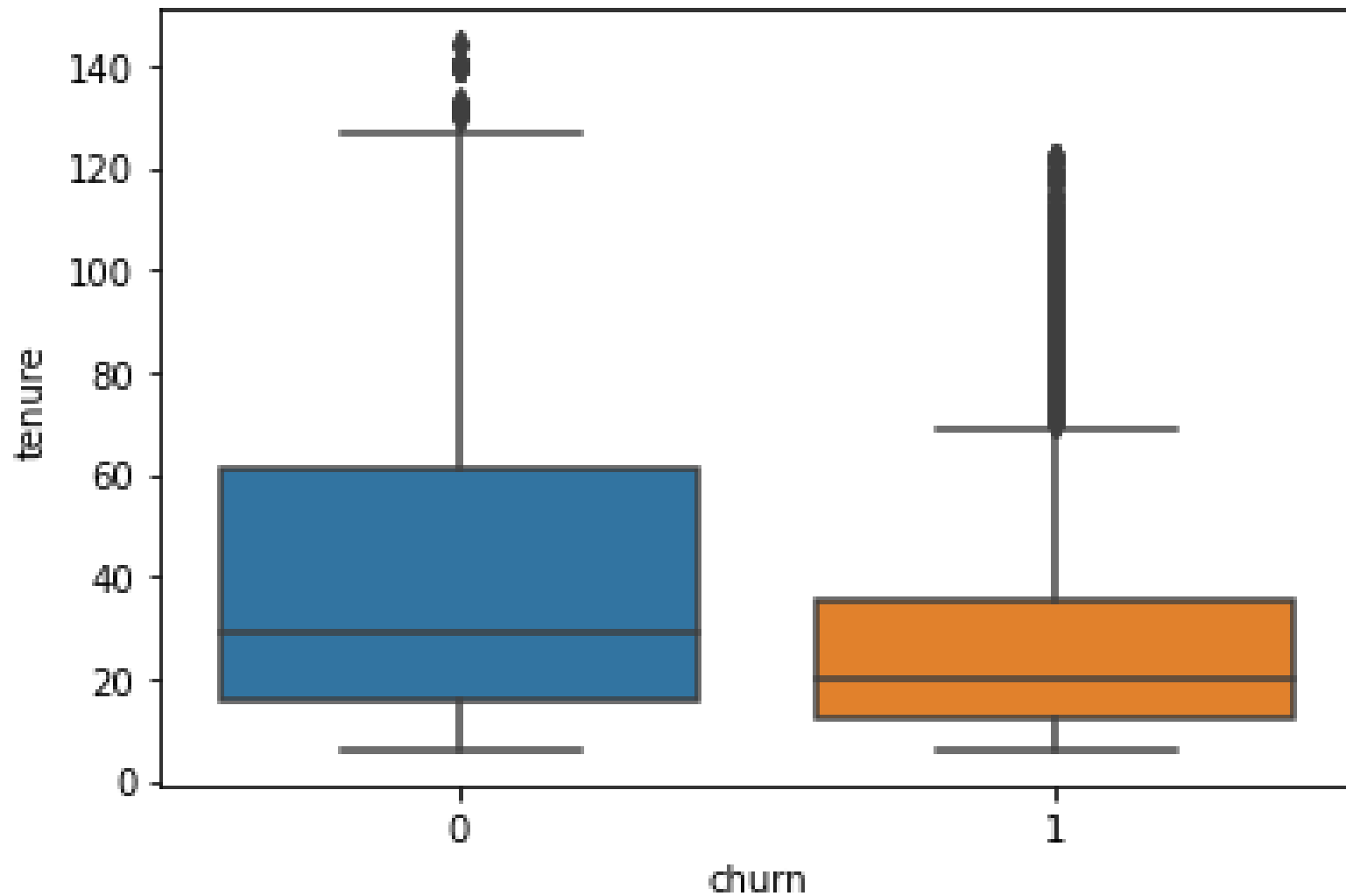
It can be seen that the maximum churn rate happens within 0-6 month, but it gradually decreases as the customer retains in the network.



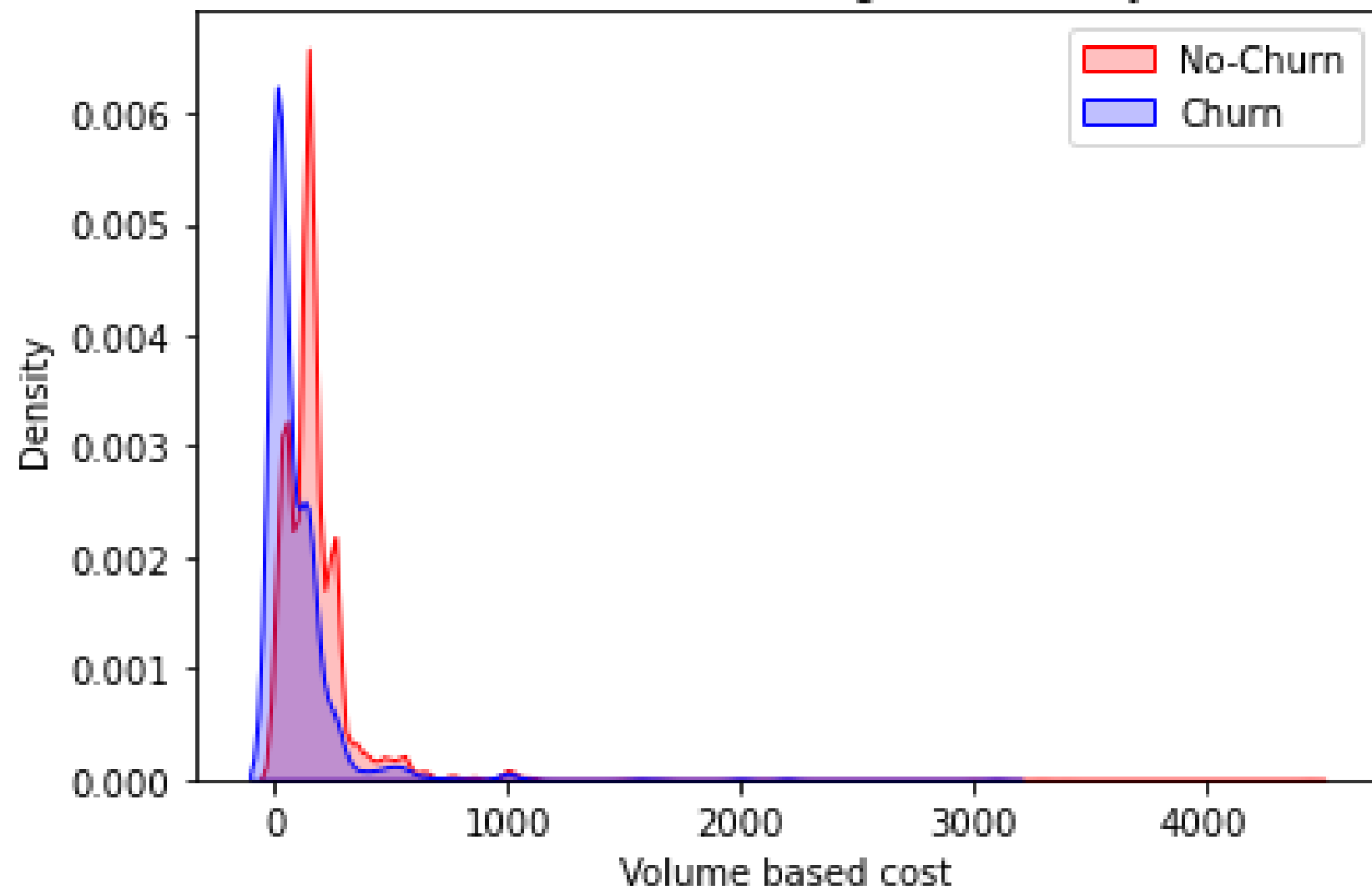




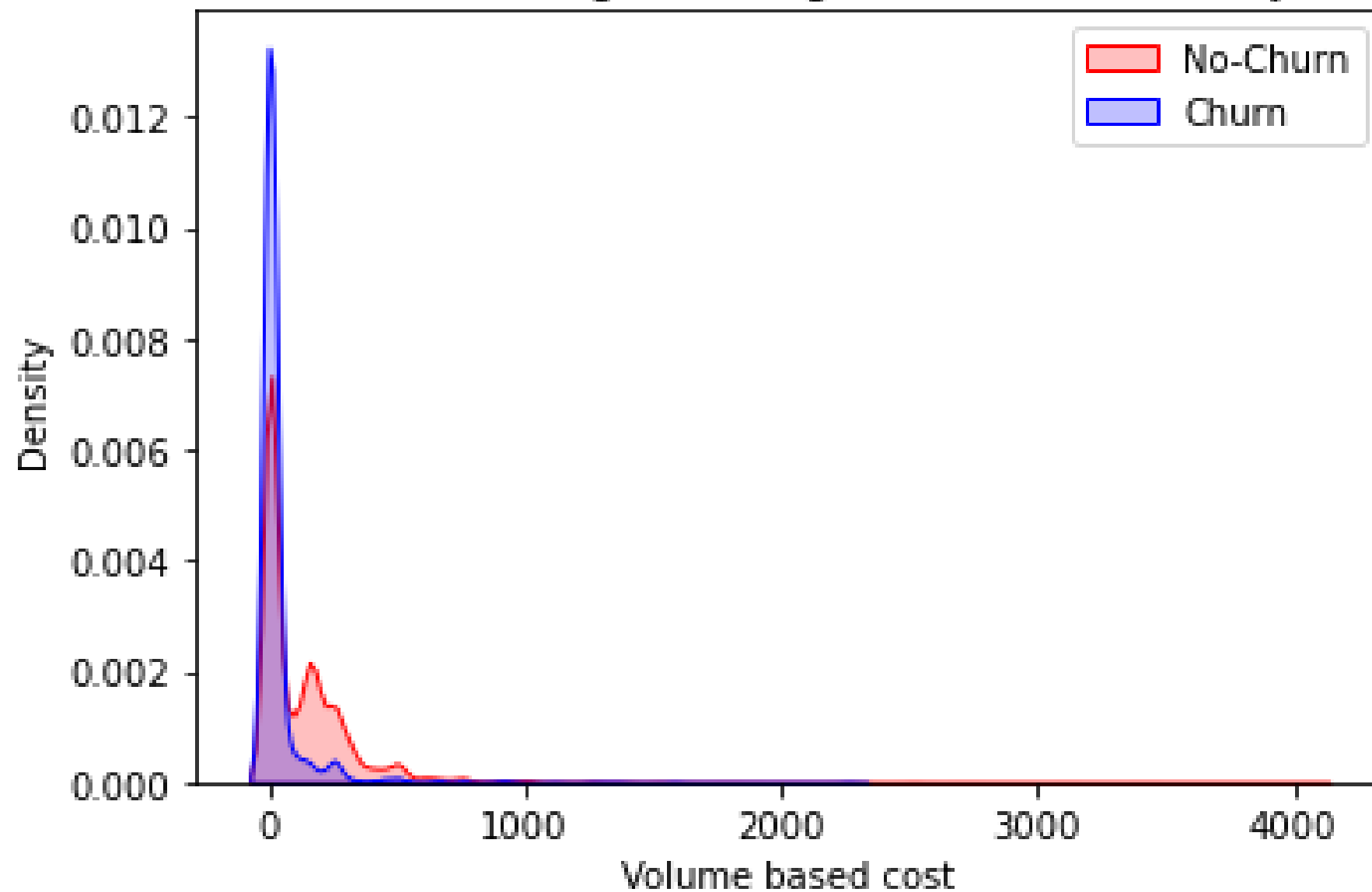
From the above plot , its clear tenured customers do no churn and they keep availing telecom services

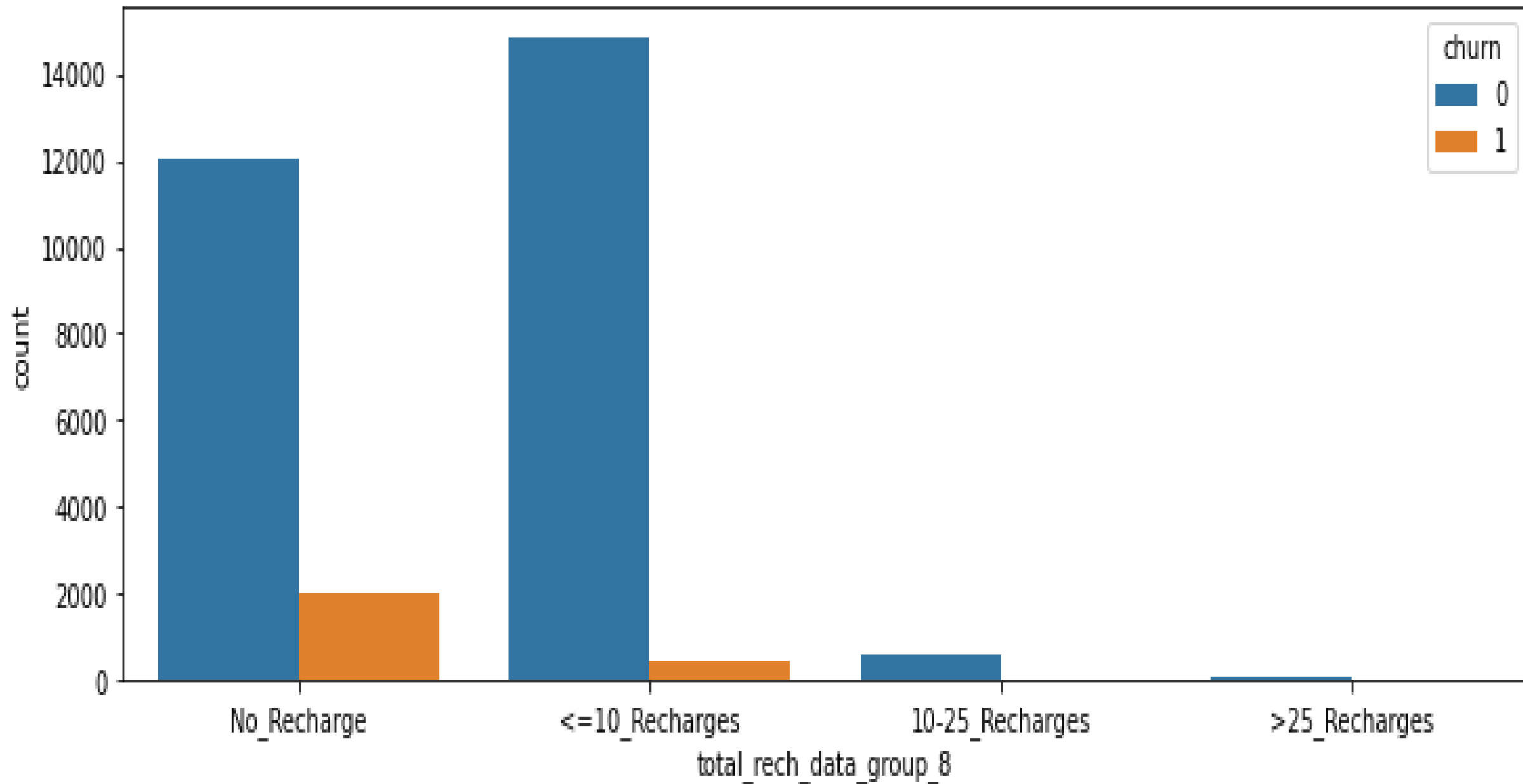


Distribution of Max Recharge Amount by churn

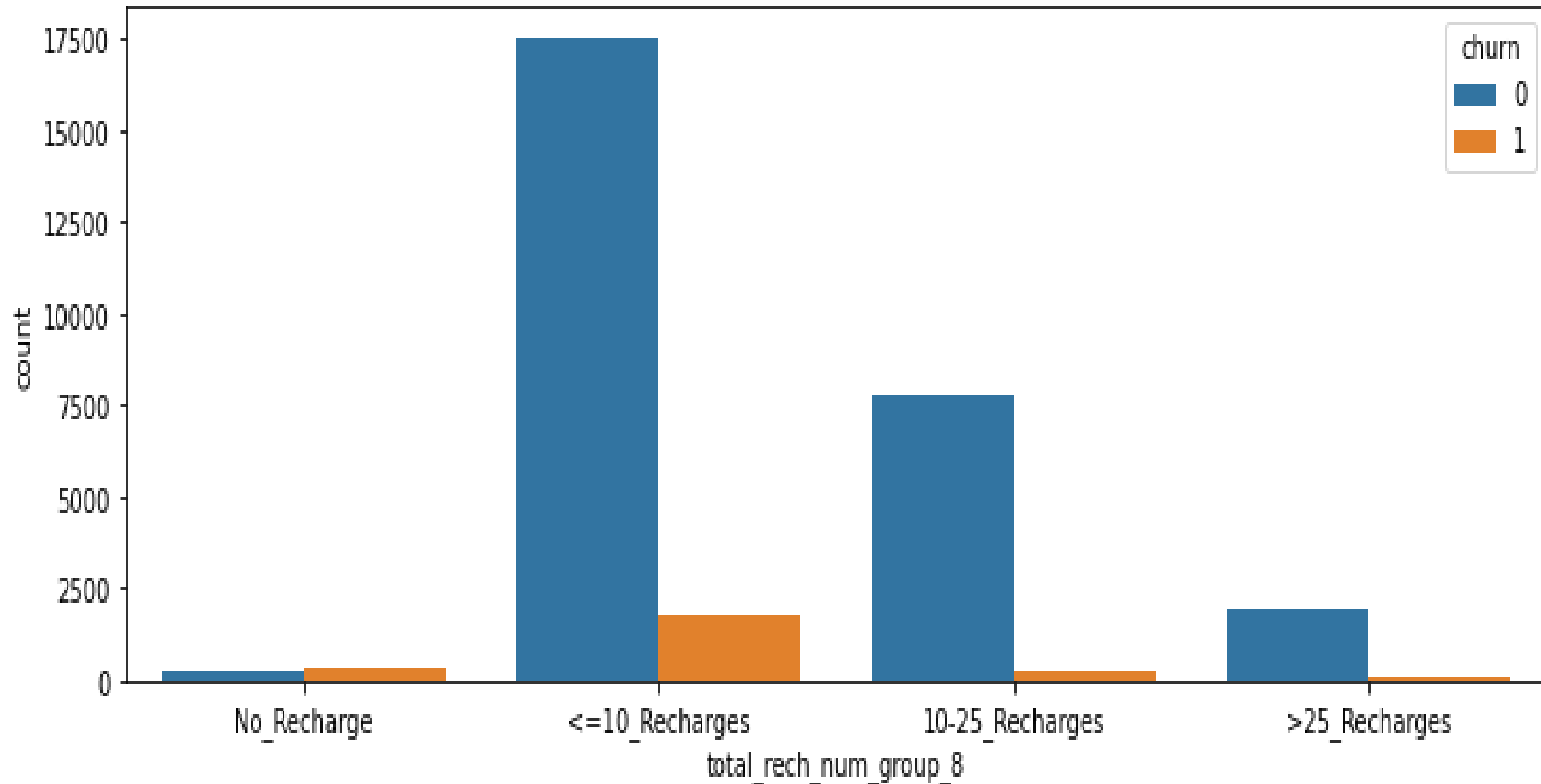


Distribution of Average Recharge Amount for Data by churn

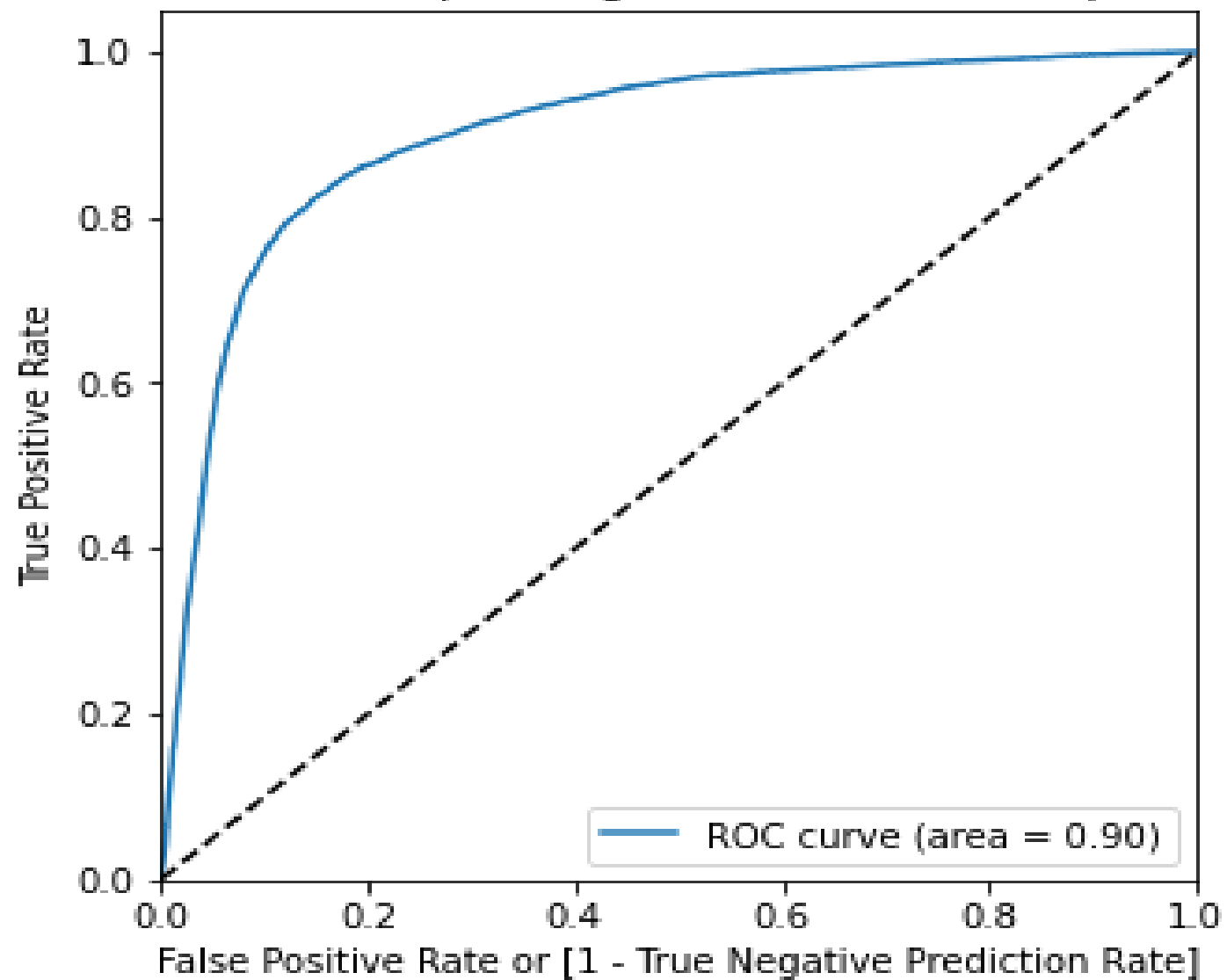


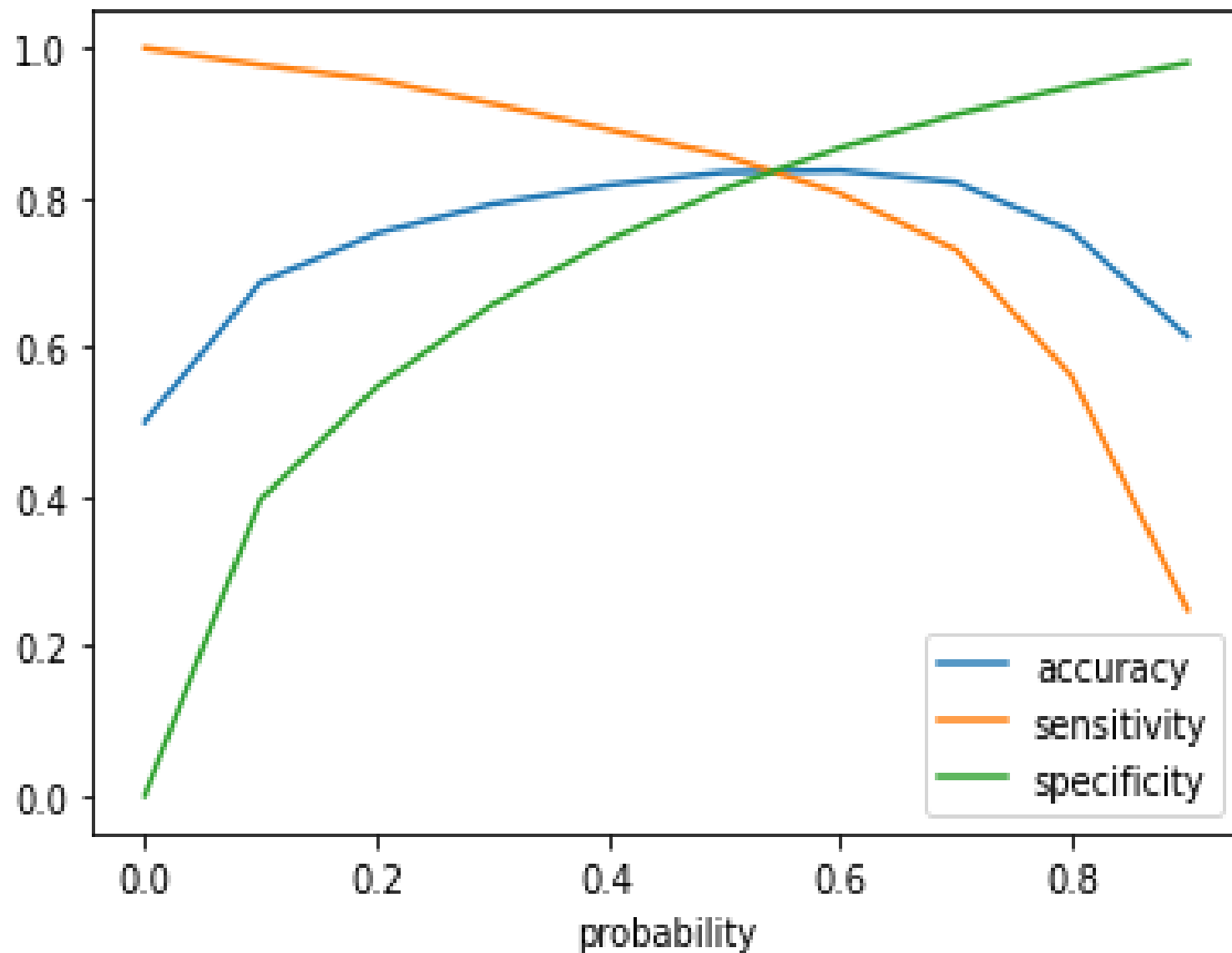


As the number of recharge rate increases, the churn rate decreases clearly.



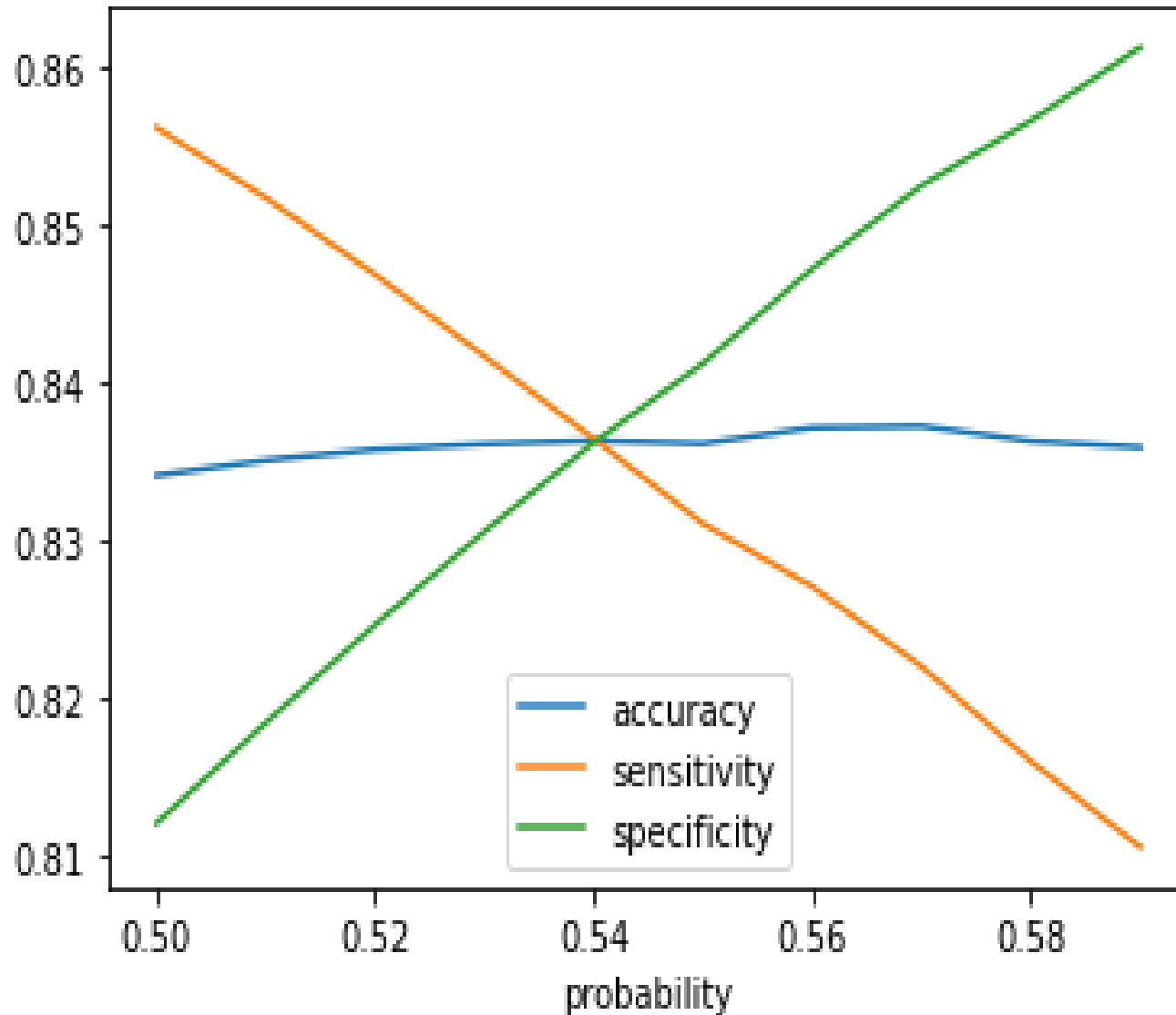
Receiver operating characteristic example





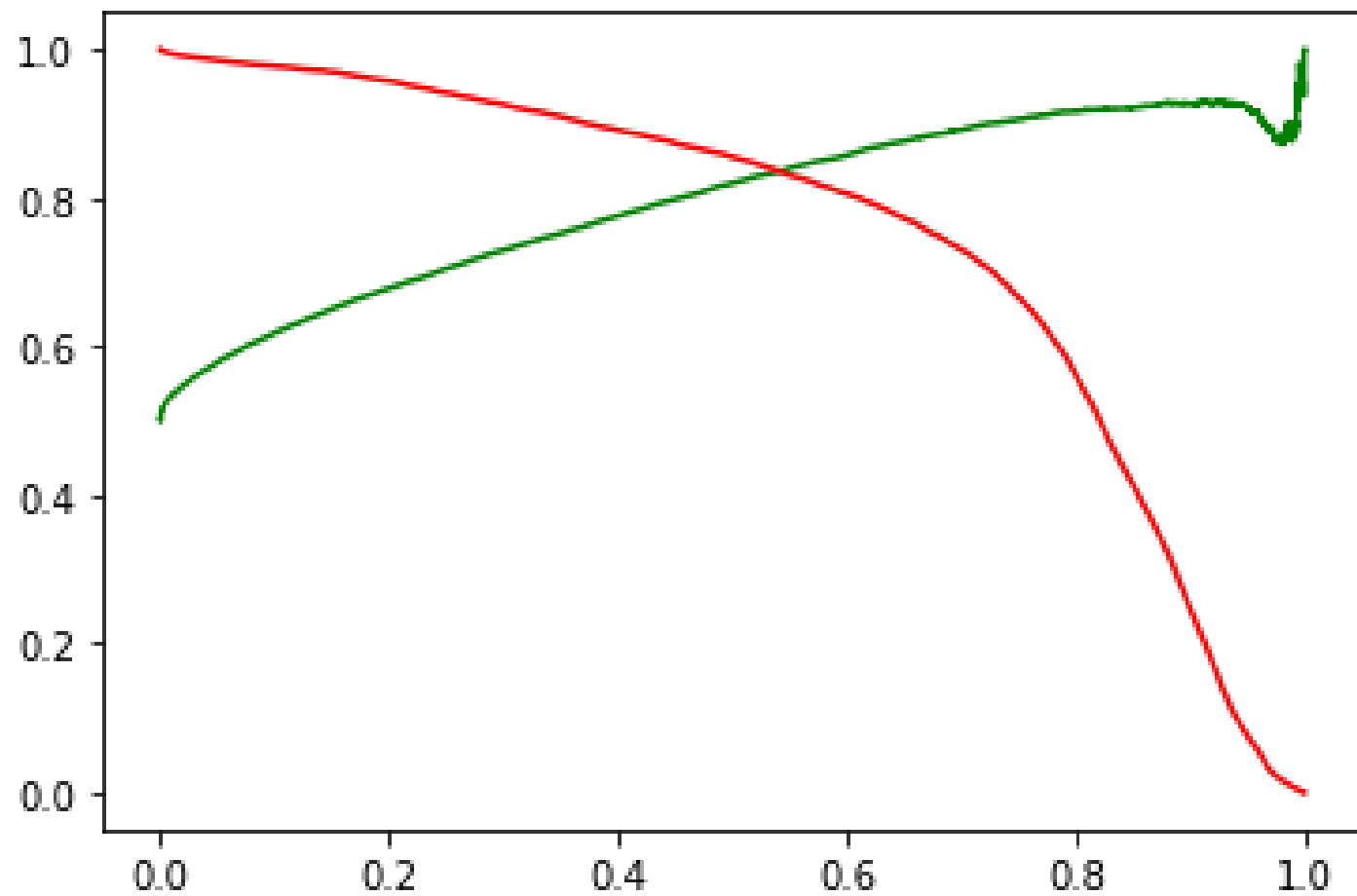
Initially we selected the optimum point of classification as 0.5.

From the above graph, we can see the optimum cutoff is slightly higher than 0.5 but lies lower than 0.6. So lets tweak a little more within this range.



From the above graph we can conclude, the optimal cutoff point in the probability to define the predicted churn variable converges at 0.54

Precision and recall trade-off



Confusion Matrix

```
[[6860 1412]
```

```
[ 145 584]]
```

Sensitivity=0.8010973936899863

Specificity = 0.8293036750483559

False Positive Rate = 0.1706963249516441

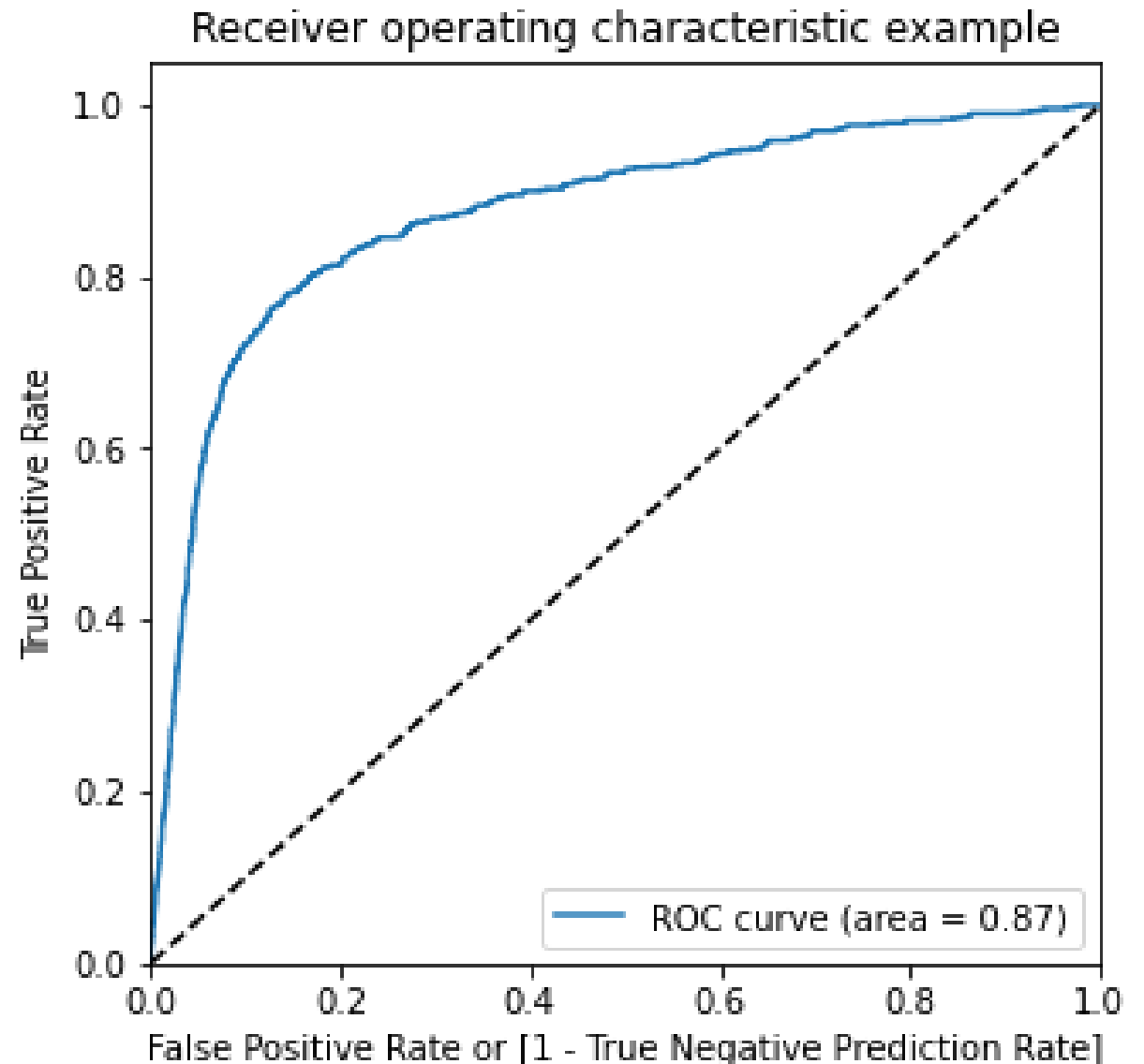
Precision = 0.2925851703406814 True Negative

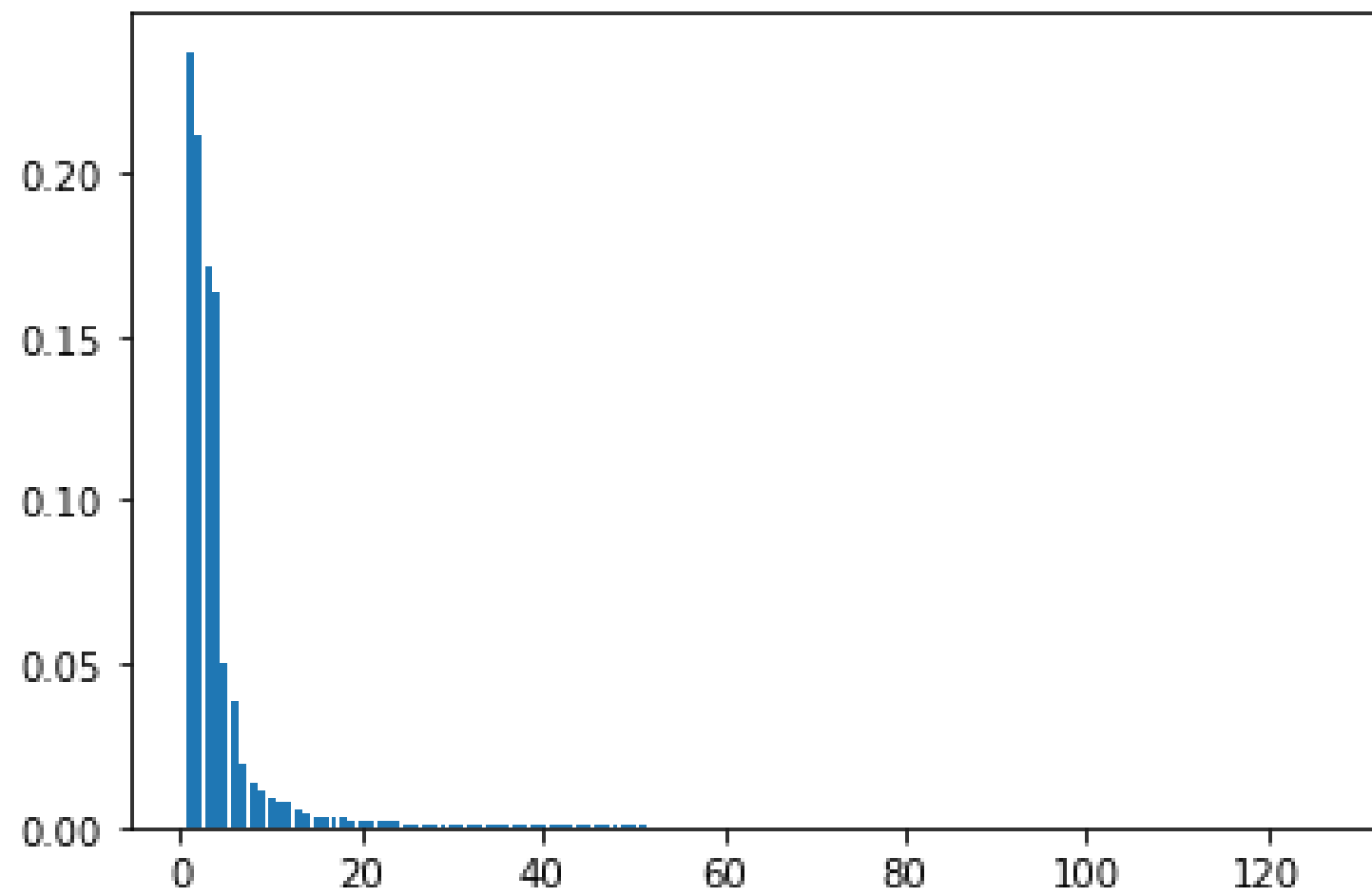
Prediction Rate = 0.979300499643112

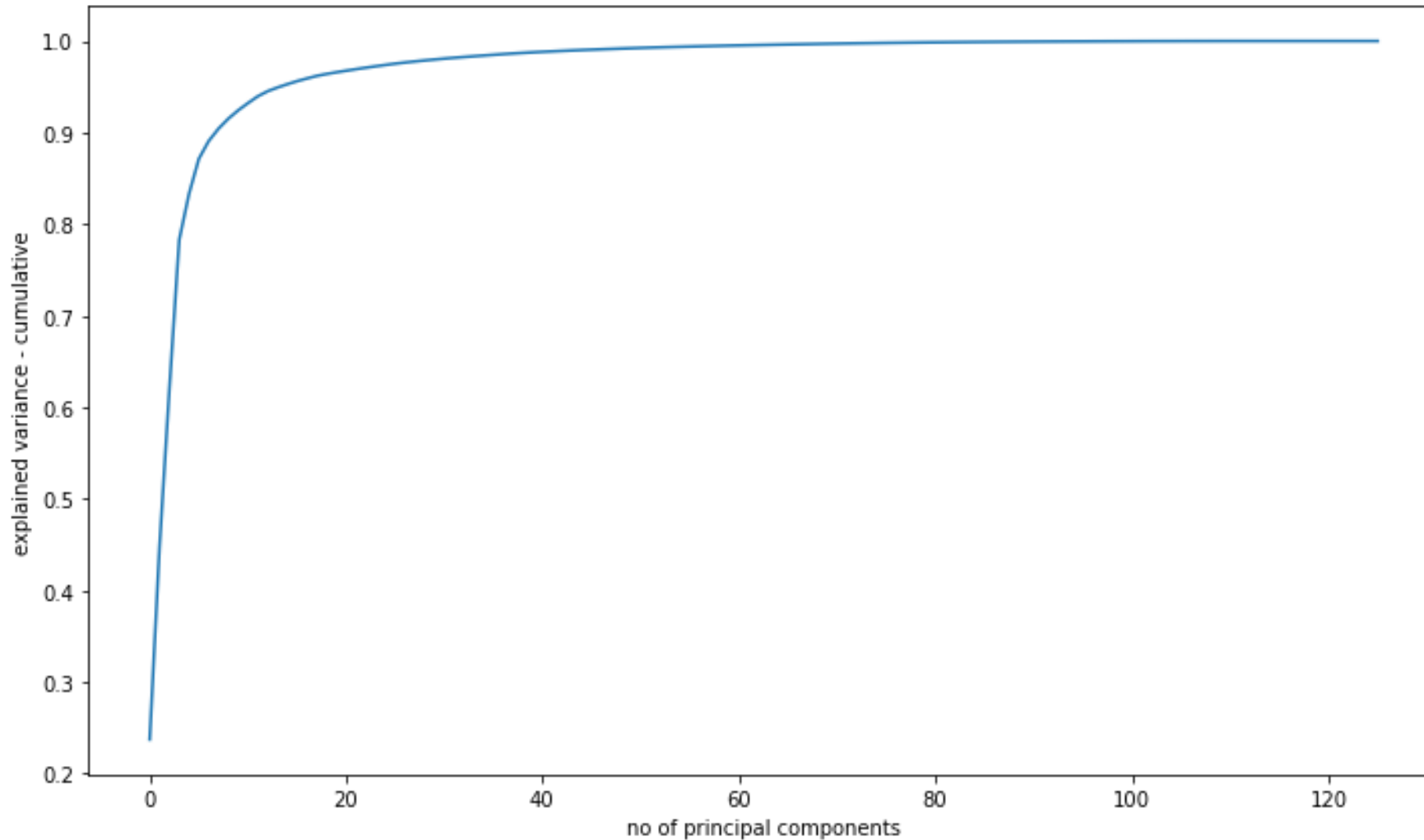
Results

The accuracy of the predicted model is: 83.0 % The sensitivity of the predicted model is: 80.0 % As the model created is based on a sensitivity model, i.e. the True positive rate is given more importance as the actual and prediction of churn by a customer

- The AUC score for train dataset is 0.90 and the test dataset is 0.87.
- This model can be considered as a good model.







Confusion Matrix for y-test & y - pred

```
[[6250 2022]  
 [ 185  544]]
```

Accuracy of the logistic
regression model with
PCA:
0.7548050216642596