# KKBOX and the meaning of Churn

# The DATA

members

user logs

transactions

train

Data Preparation

# Data Cleaning and Transformation

Generating new attributes:
- total_churn based on transactions
- monthly aggregated usage activities

Removing abnormal values:
- age up to 2000?
- days below 0?

Treating missing and null values:
- categorical: replace with most common
- numeric: replace with mean

## Split and standardize train set

**Left branch:**

RobestScaler

↓

Undersample

↓

Apply ANN and DT

**Right branch:**

Handle outliers with isolation forest

↓

StandardScaler

↓

Undersample

↓

Apply ANN and DT

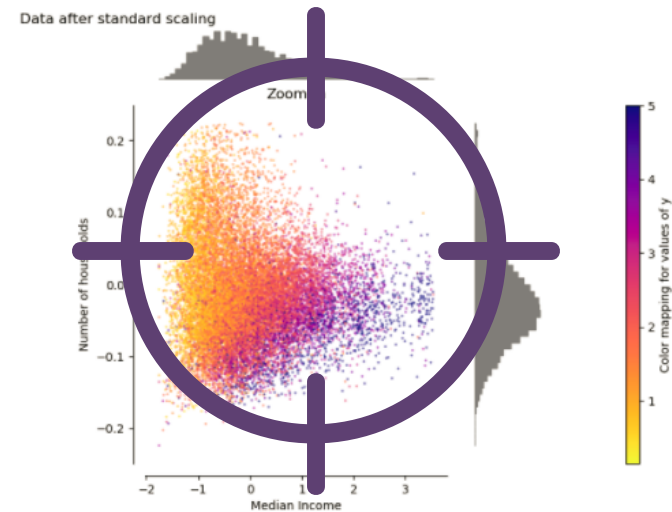# Balance train dataset and evaluate

## RobustScaler



**ANN: 90.82**
**DT: 88.59**

## StandardScaler



**ANN: 91.89**
**DT: 89.66**

# Reduce width: Feature selection

**154 features** →

PCA

RFE

MI

F-Classification

Random Forest

**59 features**

# Reduce width: Correlation Matrix



59

35

# Data Mining

**52 Models**

Logistic Regression
SVM
K-NN
Nearest Centroid
Naive Bayes
Decision Tree
ANN

# In Search for the best model...

Approach: start with a baseline model for each algorithm, fine-tune the parameters locally and then compare the chosen best models globally .

- K-NN :  optimal number for K?

- DT : combination of parameters, pure nodes removal?

- ANN : optimiser function, features, hidden layers/units?

- Logistic Regression : optimal features subset, C value?

- SVM : kernel, optimal gamma, C ?

- Naive Bayes/ Nearest Centroid : optimal features subset?
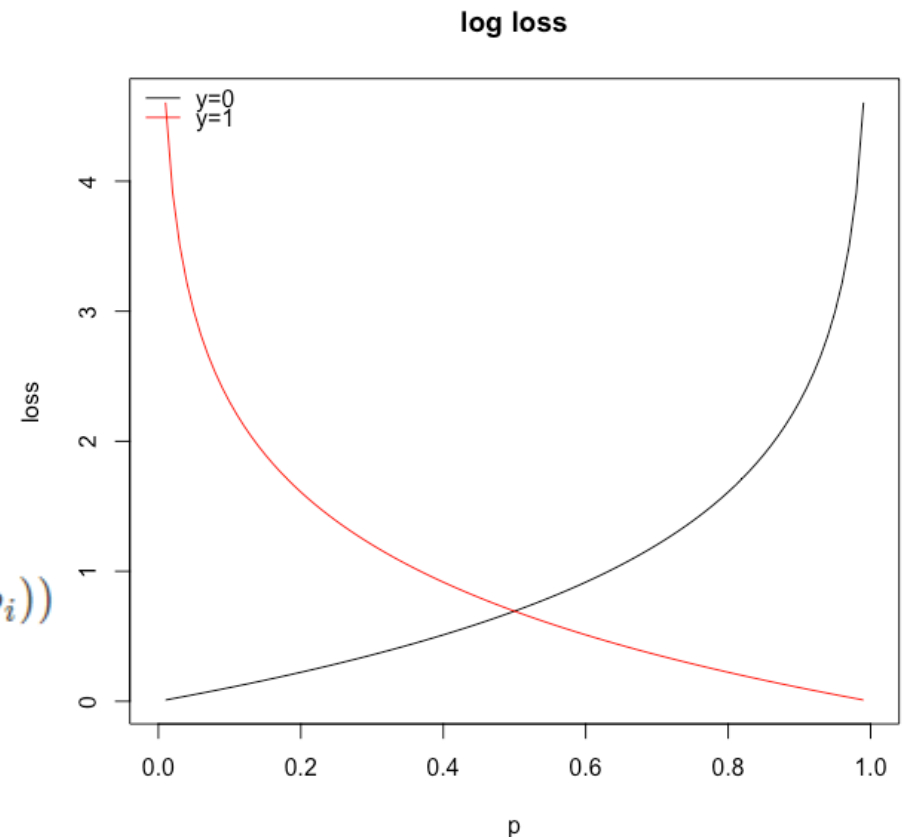
Grid Search

# Evaluation Metrics

- Not accuracy: count of correct predictions

- Recall , Precision and F1-score

- ROC curve and AUC

- **Log-loss** : account for the uncertainty of your prediction
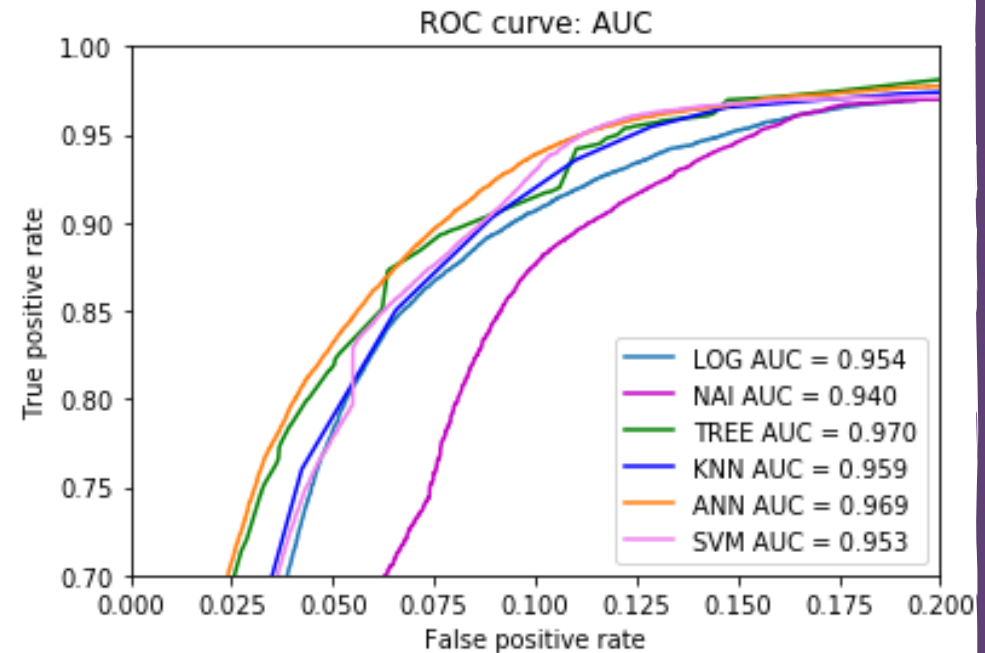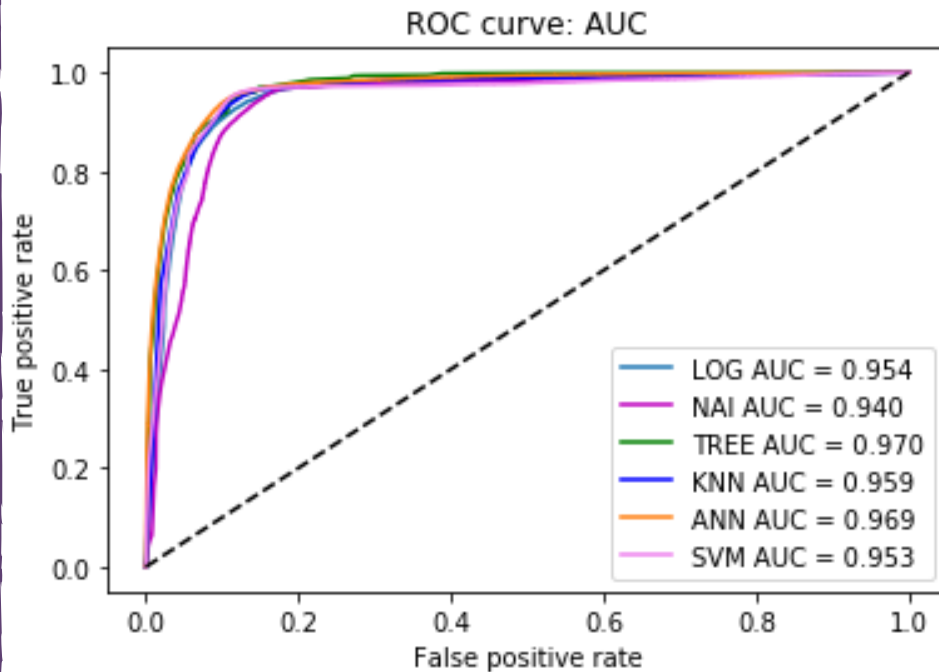
$$logloss = -\frac{1}{N} \sum_{i=1}^{N} (y_i \log(p_i) + (1 - y_i) \log(1 - p_i))$$



log loss

# Performance Evaluation

| Evaluation | K-NN | NC | LogR | DT | ANN | NB | SVM |
|---|---|---|---|---|---|---|---|
| AUC | 0,96 | - | 0,95 | 0,97 | 0,97 | 0,94 | 0,95 |
| Log Loss | 0,82 | - | 0,27 | 0,20 | 0,21 | 0,63 | 0,27 |
| Precision-avg/total | 0,95 | 0,95 | 0,95 | 0,95 | 0,96 | 0,95 | 0,95 |
| Recall-avg/total | 0,89 | 0,89 | 0,89 | 0,90 | 0,91 | 0,84 | 0,89 |
| F1 Score -avg/total | 0,91 | 0,91 | 0,91 | 0,92 | 0,92 | 0,87 | 0,91 |
| Precision-class 1 | 0,37 | 0,36 | 0,37 | 0,39 | 0,41 | 0,27 | 0,36 |
| Recall-class 1 | 0,89 | 0,86 | 0,92 | 0,92 | 0,93 | 0,97 | 0,96 |
| F1 Score-class 1 | 0,53 | 0,50 | 0,53 | 0,54 | 0,57 | 0,43 | 0,53 |

# ROC curve AUC



ROC curve: AUC

LOG AUC = 0.954
NAI AUC = 0.940
TREE AUC = 0.970
KNN AUC = 0.959
ANN AUC = 0.969
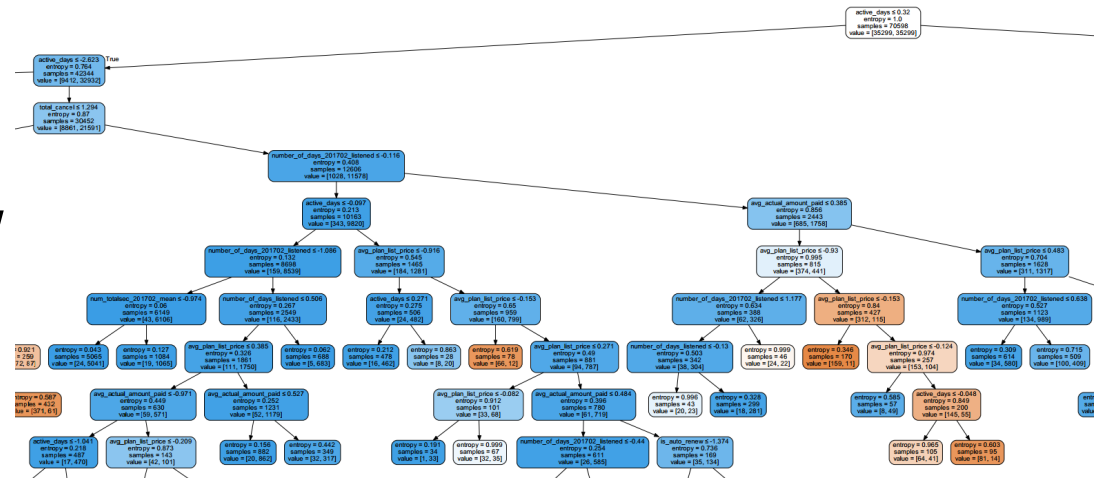SVM AUC = 0.953

Conclusion and Outlook

# Decision Tree

**#1**

Optimal parameters:
- criterion='entropy',
- min impurity decrease=0.000015,
- max depth=15,
- max leaf nodes=None,
- min samples leaf=10,
- min samples split=20



## Predicted

|  |  | 0 | 1 |
|---|---|---|---|
| **Actual** | 0 | 219,970 | 25,580 |
|  | 1 | 1436 | 15,857 |

# Discussion

Drawbacks:
- other balancing approaches: oversampling
- scale out to the whole timeframe (3 years)

Outlook:
- hybrid model*

*Lee, Jae, and Jin Lee. "Customer churn prediction by hybrid model." Advanced Data Mining and Applications (2006): 959-966.

# Thank You !

Presented by
Team 9:

Simona Doneva
Na Gong
Bengi Koseouglu
Siying Liu
Liam Pang

# Questions?

# Decision Tree on train data before standardization and balancing

Optimal parameters:
- criterion='entropy',
- min impurity decrease=0.000015,
- max depth=15,
- max leaf nodes=None,
- min samples leaf=10,
- min samples split=20

|  |  | Predicted | |
|---|---|---|---|
|  |  | 0 | 1 |
| Actual | 0 | 243,240 | 2310 |
|  | 1 | 8212 | 9081 |

# Top Features

- active days
- total churn
- total cancel
- actual payment
- registration method
- total_sec_listened
- total_days_listened
....