# Personal Loan Status Prediction using XGBoost

# PROBLEM STATEMENT
# &
# UNDERSTANDING DATASET

o The collection of the dataset

- LendingClub Notes platform

- Duration: 2007–2020 | USD | Kaggle

o Content of the data

- (1) Demographic data such as employment status, employment length and house ownership,

- (2) Loan Characteristics such as loan amount, term and purpose and

- (3) Behavioral data related to historical payments

o Dimension: 2925493 rows and 141 columns

o Data choosen for this project: Sub–sample of 700000/2925493

o Univariate EDA

- Numerical data

- Categorical data

o Bivariate EDA
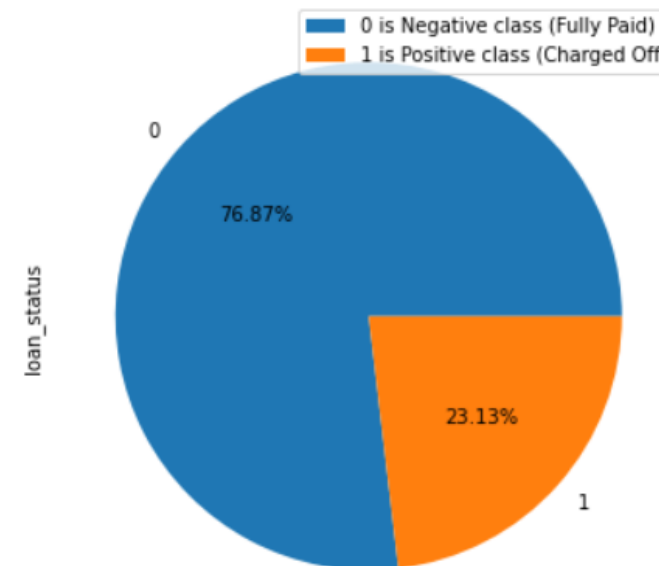
Loan status

Loan amount

Purpose

Interest rates

o Conclusion of EDA

- Most of accepted borrowers: a **good** profile (good average FICO scores and no public records of bankruptcies), mainly with **financial reasons** (Debt consolidation and credit card).

- Most borrowers: a shorter loan term (**36 months**, compared to 60 months) and do **not** enroll into the **hardship** plan, but the percentage of people who are Charged Off and default accounts for nearly **one third** of the data sample.

- Checking bivariance relationship **cannot** show that the features such as annual income, verified status, house ownership, credit line, etc. have a specific effect on Charged Off, Default or Late status.

- People with **lower interest rate** and **lower loan amount** likely to pay off loans successfully.

# Data Preparing

o Cleaning columns

o Cleaning rows

o Convert data types (one-hot coding for setting dummies)

o Set up predictors (43 variables) and responding variable (loan_status)

o Data split: Stratified style with proportion of train-test set (70%–30%)

o Imbalanced data:

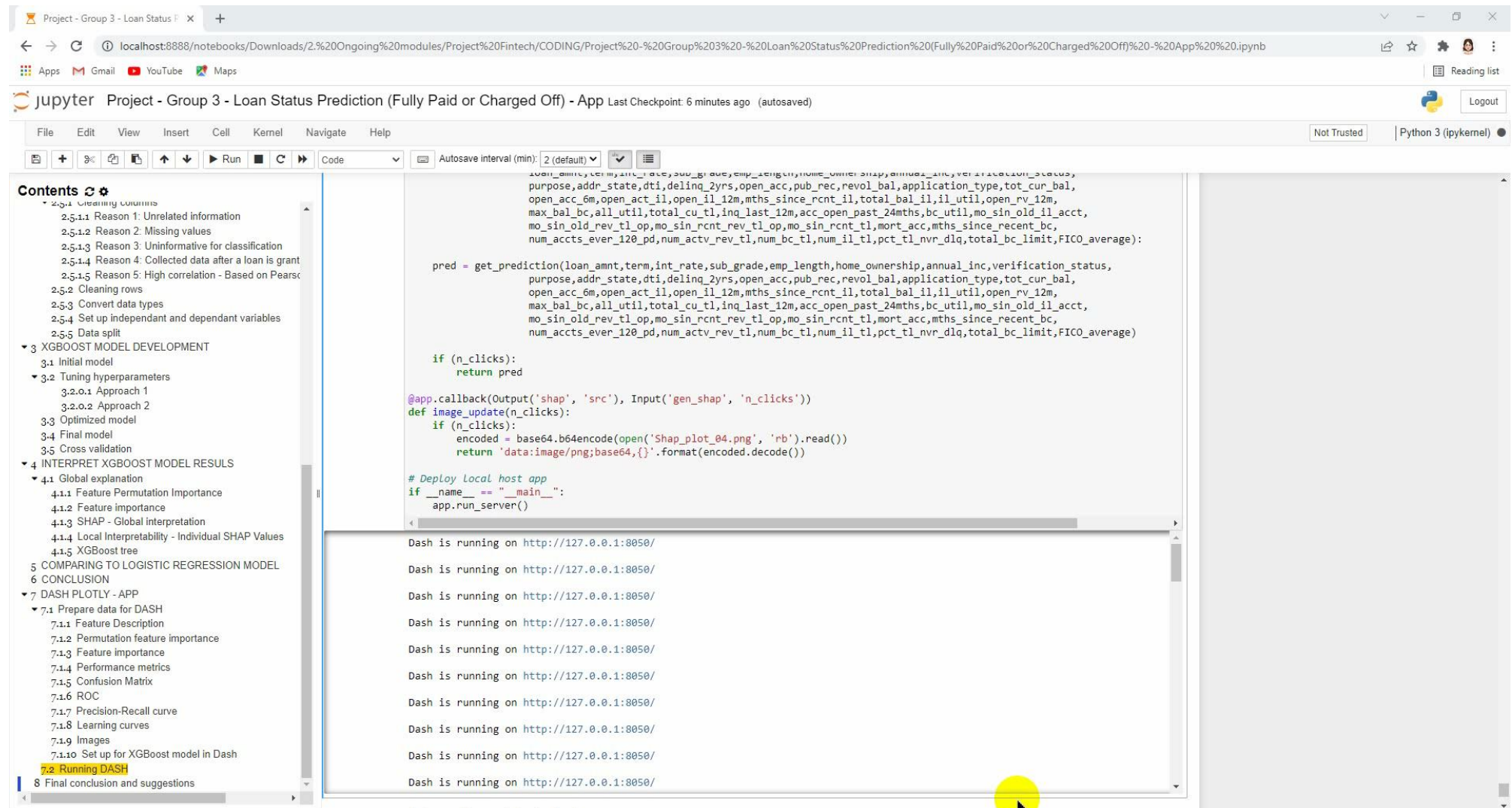Classes of Charged off and Fully Paid in Loan status



0 is Negative class (Fully Paid)
1 is Positive class (Charged Off)

0

76.87%

loan_status

23.13%

1

# DASH Plotly (APP)

- It is a localhost web

- Conducted using DASH plot library

- Has 3 different tabs

# DASH Plotly (APP) – Link

# HOW XGBOOST WORKS

# Overview



- Ensemble algorithm - creating a model by combining a number of baser learners
- Boosting: Sequential ensemble



*Image Source: edureka.co*

# XGBoost math in brief

- Extreme gradient boosting (XGBoost) is a decision-tree-based ensemble algorithm
- Log loss function: the negative log-likelihood for **Binary Classification** (Daniel & Martin(2021). Speech and Language Processing. 5.pdf (stanford.edu). P.152)

$$L(y_i, p_i) = -[y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]$$

- The general objective function *(Chen & Guestrin, 2016)* is:

$$L(\phi) = \sum_i l(\hat{y}_i, y_i) + \sum_\kappa \Omega(f_k)$$

$fk$ is an independent tree,
T is the number of leaves in a tree
w is leaf weight,
γ and λ are hyperparameters.

$$\Omega(f) = \gamma T + \frac{1}{2}\lambda\|w\|^2$$

a loss function measuring
how well model fit on the
training data

regularization to measure the
complexity of trees

OTTO VON GUERICKE
UNIVERSITÄT
MAGDEBURG

IFM

LEHRSTUHL FÜR
INNOVATIONS- UND
FINANZMANAGEMENT

Financial
Technology
ego.-INKUBATOR

# How XGBoost works

*Fomulas source (in this page):*
Tianqi Chen and Carlos Guestrin.
2016.

$$L(\phi) = \sum_i l(\hat{y}_i, y_i) + \sum_{\kappa} \Omega(f_k)$$

Boosting

$$\hat{y}_i^{(t)} = \hat{y}_i^{(t-1)} + f_t(x_i)$$

Plug in

$$\mathcal{L}^{(t)} = \sum_{j=1}^{n} l(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \Omega(f_t)$$

Taylor
approximation $\quad f(x + \Delta x) \cong f(x) + f'(x)\Delta x + \frac{1}{2}f''(x)\Delta x^2$

$$\mathcal{L}^{(t)} \cong \sum_{i=1}^{n} [l(y_i, \hat{y}_i^{(t-1)}) + g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i)] + \Omega(f_t)$$

$$\text{where } g_i = \partial_{\hat{y}^{(t-1)}} l(y_i, \hat{y}^{(t-1)}) \text{ and } h_i = \partial_{\hat{y}^{(t-1)}}^2 l(y_i, \hat{y}^{(t-1)})$$

$$\Omega(f) = \gamma T + \frac{1}{2}\lambda\|w\|^2$$

*Fomulas source (in this page):*
Tianqi Chen and Carlos Guestrin. 2016.

Plug in

$$\tilde{\mathcal{L}}^{(t)} = \sum_{i=1}^{n}[g_i f_t(\mathbf{x}_i) + \frac{1}{2}h_i f_t^2(\mathbf{x}_i)] + \gamma T + \frac{1}{2}\lambda\sum_{j=1}^{T}w_j^2$$

$$= \sum_{j=1}^{T}[(\sum_{i\in I_j}g_i)w_j + \frac{1}{2}(\sum_{i\in I_j}h_i + \lambda)w_j^2] + \gamma T$$

Solve the quadratic function of $w$

$$w_j^* = -\frac{\sum_{i\in I_j}g_i}{\sum_{i\in I_j}h_i + \lambda}$$

Output, minimized objective function

$$\tilde{\mathcal{L}}^{(t)} = -\frac{1}{2}\sum_{j=1}^{T}\frac{(\sum_{i\in I_j}g_i)^2}{\sum_{i\in I_j}h_i + \lambda} + \gamma T$$

# How XGBoost grows trees

- XGBoost begins with a weak learner

- Loop from 1 to k:
    - Build the first tree
    - Learn the structure and use the minimized objective function
    - To avoid overfitting, learning rate and gamma is added to each tree.
    - Final model = all of the trees are combined additively

- The process stops when:
    - Gain score becomes negative (cannot gain further information from splitting)
    - fixed number of the iteration (k) is reached

- Source: (Chen & Guestrin, 2016)
o For a big data sample:
  - sparsity-aware split finding algorithm to handle the problem of missing values in the data
  - For a big dataset – training speed: fast

o In practice:
  - Python package of XGBoost
  - Parallelization of tree construction
  - Out-of-Core Computing
  - Cache Optimization

# Results

Model parameter:

```python
%%time
xgb_model = xgb.XGBClassifier(objective='binary:logistic',
                              learning_rate=0.05,
                              scale_pos_weight=3,
                              n_estimators=200,
                              max_depth=5,
                              min_child_weight=1,
                              gamma=0.5,
                              colsample_bytree=1,
                              subsample=1,
                              use_label_encoder=False,
                              random_state=42)

xgb_model.fit(x_train, y_train,
          verbose=0,
          early_stopping_rounds=10,
          eval_metric='aucpr',
          eval_set=[(x_train, y_train), (x_test, y_test)])
```

- Development:
  - Initial mode
  - optimized model (tunning hyperparameters)
  - final mode
  - Each model is evaluated with confusion matrix, metrics scores...

- Trade-off: Bias vs Variance

1) CLASSIFICATION REPORT:

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.86 | 0.69 | 0.77 | 45006 |
| 1 | 0.38 | 0.64 | 0.48 | 13542 |
| accuracy | | | 0.68 | 58548 |
| macro avg | 0.62 | 0.66 | 0.62 | 58548 |
| weighted avg | 0.75 | 0.68 | 0.70 | 58548 |

2) CONFUSION MATRIX



- Target: High recall for positive class & high accuracy
- Cost of TP and FN
- Trade-off: precision vs recall score

$$accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$precision = \frac{TP}{TP + FP}; \quad recall = \frac{TP}{TP + FN}$$

# Results

ROC Curve (AUC=0.73)

Precision and Recall Trade-off (AUC PR=0.33)

# Explanation of the model

Decision Tree

int_rate<12.0599995

no

yes, missing

int_rate<8.02499962

yes, missing

no

int_rate<7.2750001

int_rate<10

yes, missing

FICO_average<764.5

no

FICO_average<694.5

yes, missing

num_actv_rev_tl<6.5

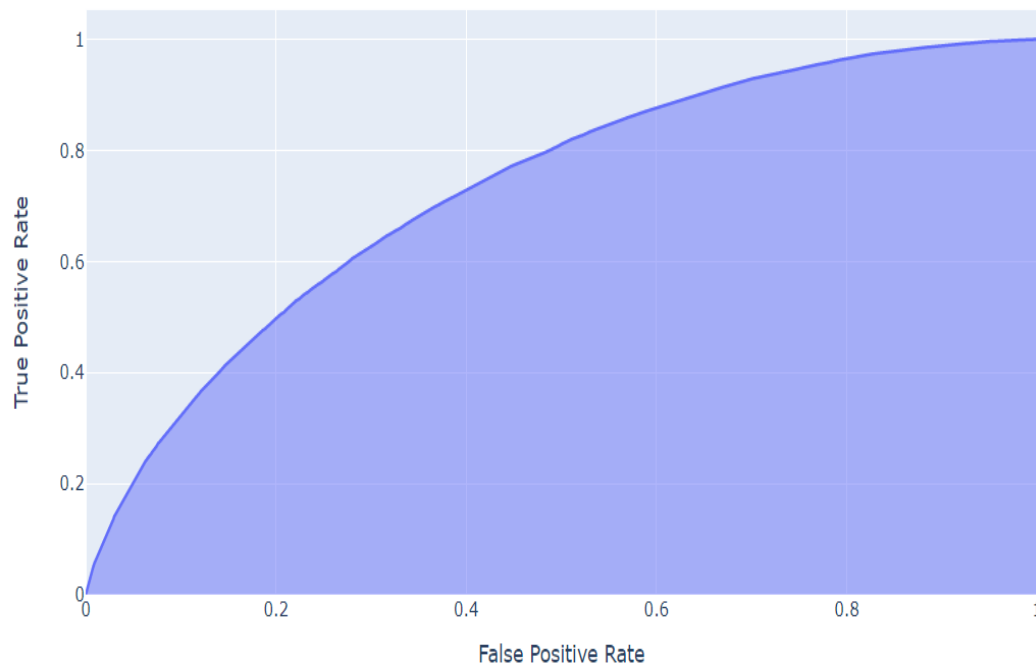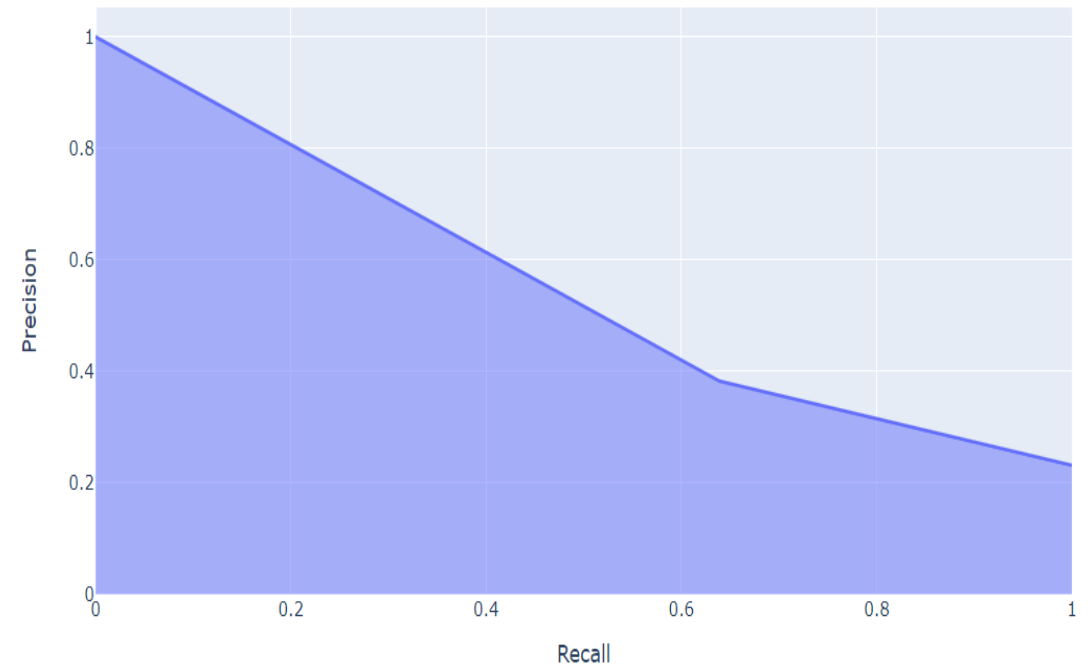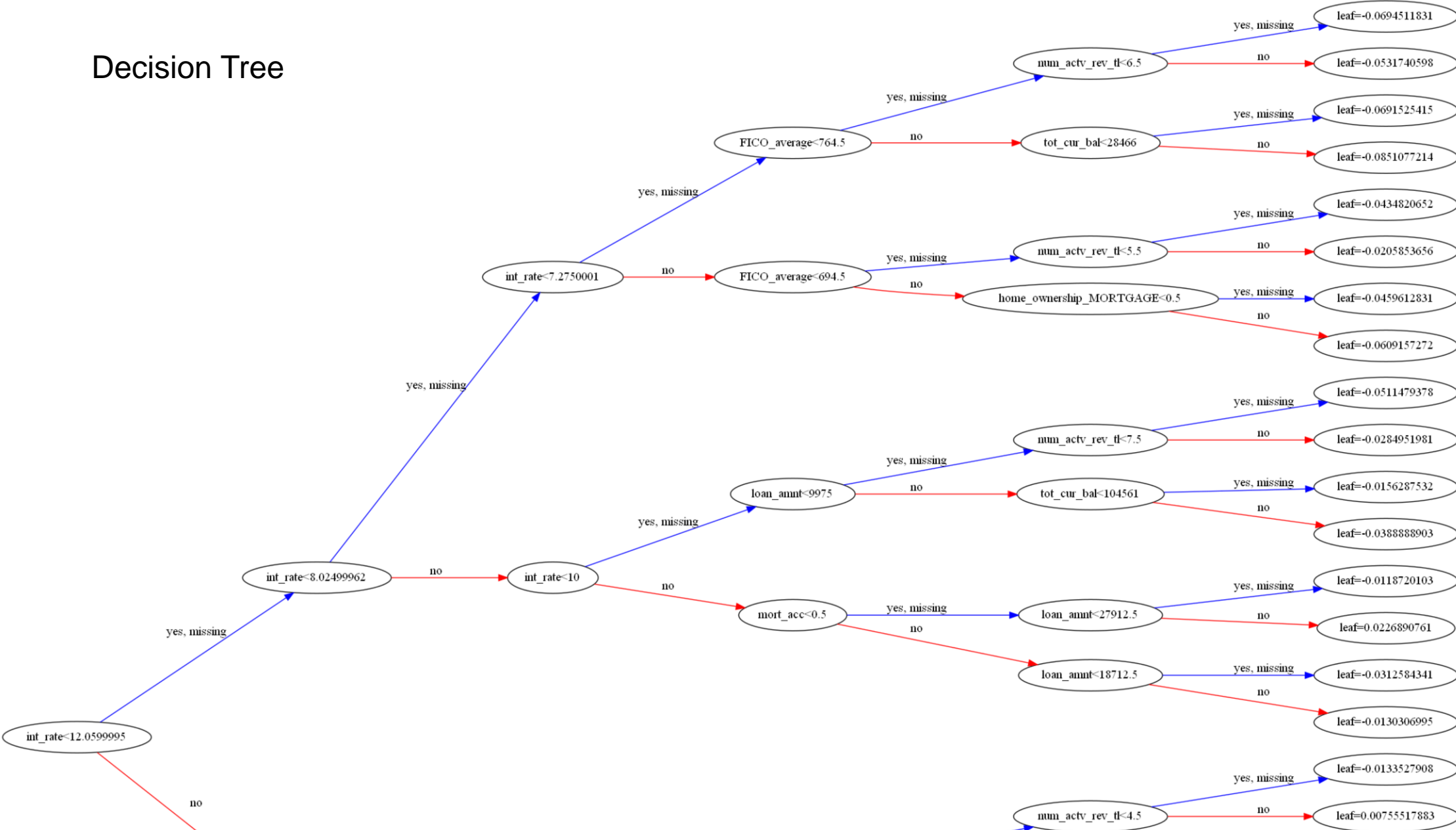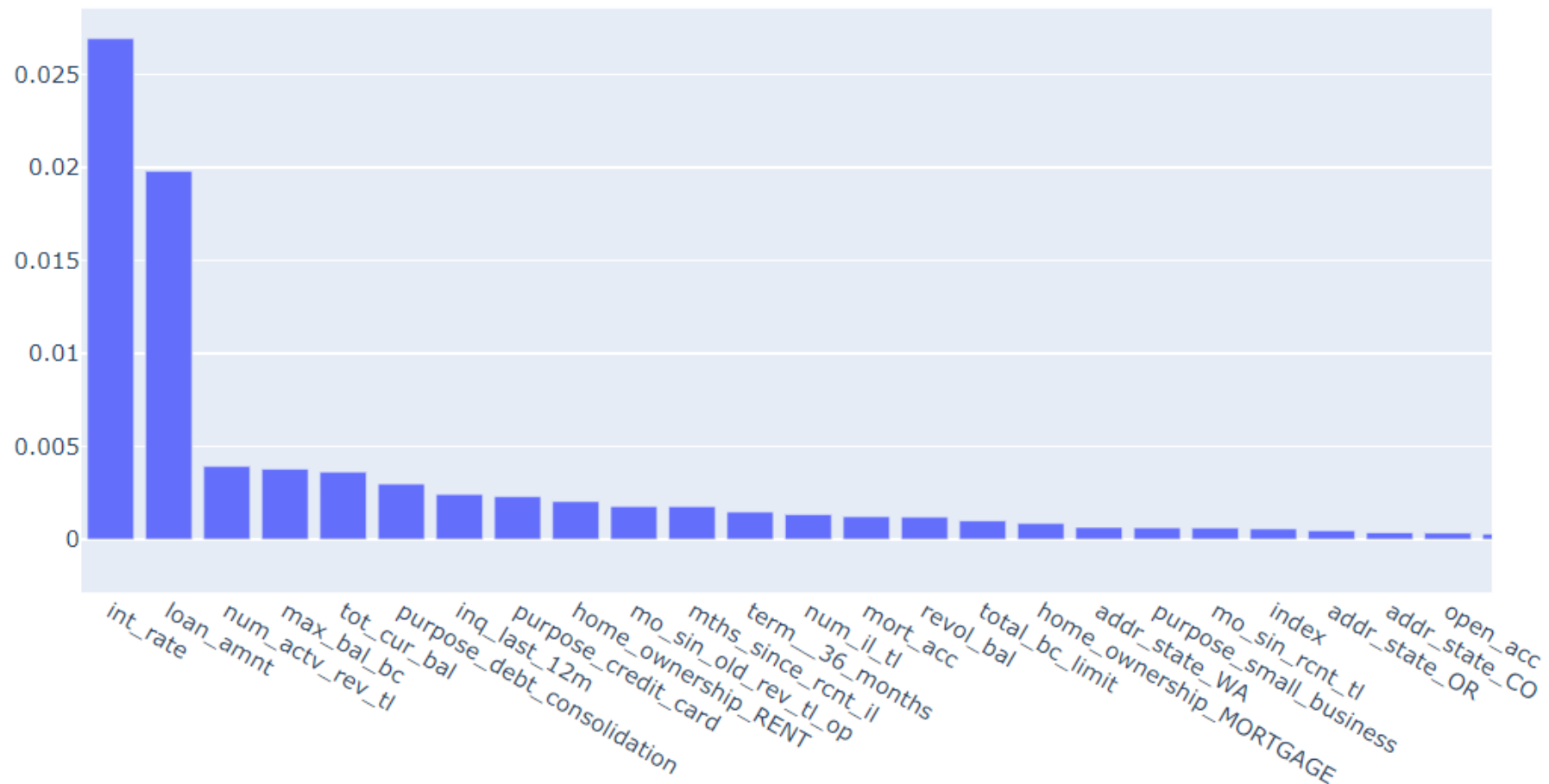yes, missing → leaf=-0.0694511831

no → leaf=-0.0531740598

no

tot_cur_bal<28466

yes, missing → leaf=-0.0691525415

no → leaf=-0.0851077214

yes, missing

num_actv_rev_tl<5.5

yes, missing → leaf=-0.0434820652

no → leaf=-0.0205853656

no

home_ownership_MORTGAGE<0.5

yes, missing → leaf=-0.0459612831

no → leaf=-0.0609157272

yes, missing

loan_amnt<9975

yes, missing

num_actv_rev_tl<7.5

yes, missing → leaf=-0.0511479378

no → leaf=-0.0284951981

no

tot_cur_bal<104561

yes, missing → leaf=-0.0156287532

no → leaf=-0.0388888903

no

mort_acc<0.5

yes, missing

loan_amnt<27912.5

yes, missing → leaf=-0.0118720103

no → leaf=0.0226890761

no

loan_amnt<18712.5

yes, missing → leaf=-0.0312584341

no → leaf=-0.0130306995

num_actv_rev_tl<4.5

yes, missing → leaf=-0.0133527908
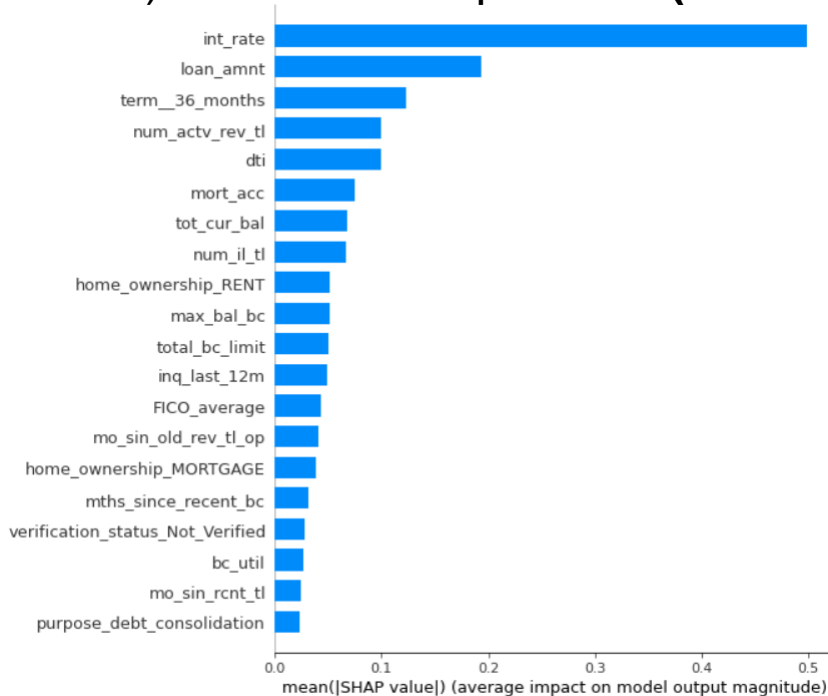
no → leaf=0.0075517883
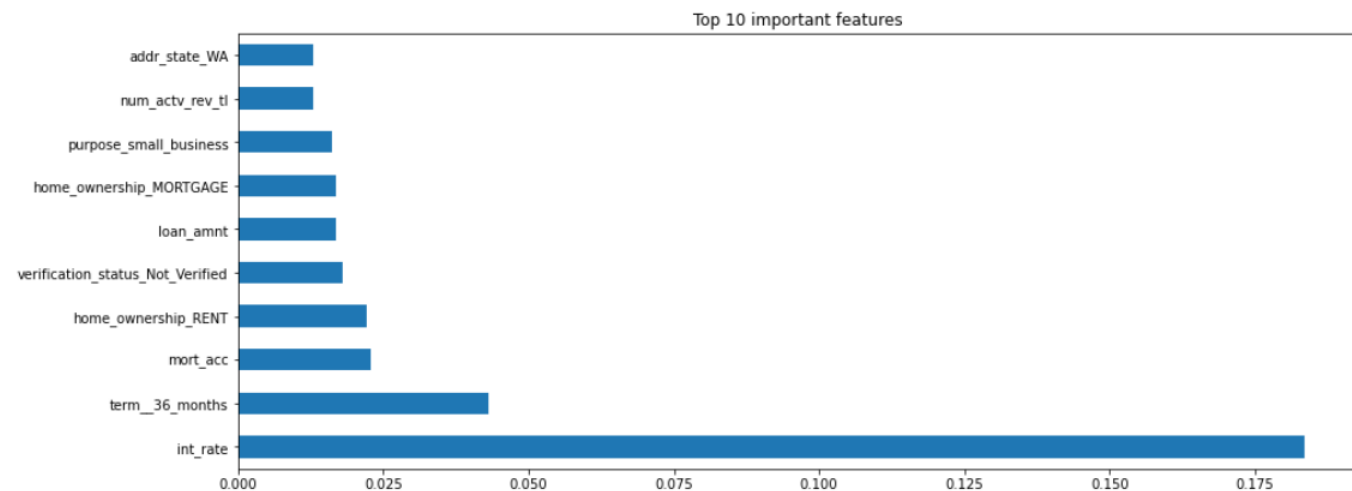
Permutation feature importance
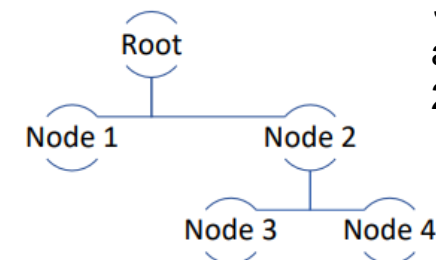
# Explanation of the model

Feature Importance
a) Based on Shap values **(Lundberg & Lee, 2017**)          b) Based on Gain values



$$\mathcal{L}_{split} = \frac{1}{2}\left[\frac{\left(\sum_{i\in I_L} g_i\right)^2}{\sum_{i\in I_L} h_i + \lambda} + \frac{\left(\sum_{i\in I_R} g_i\right)^2}{\sum_{i\in I_R} h_i + \lambda} - \frac{\left(\sum_{i\in I} g_i\right)^2}{\sum_{i\in I} h_i + \lambda}\right] - \gamma$$
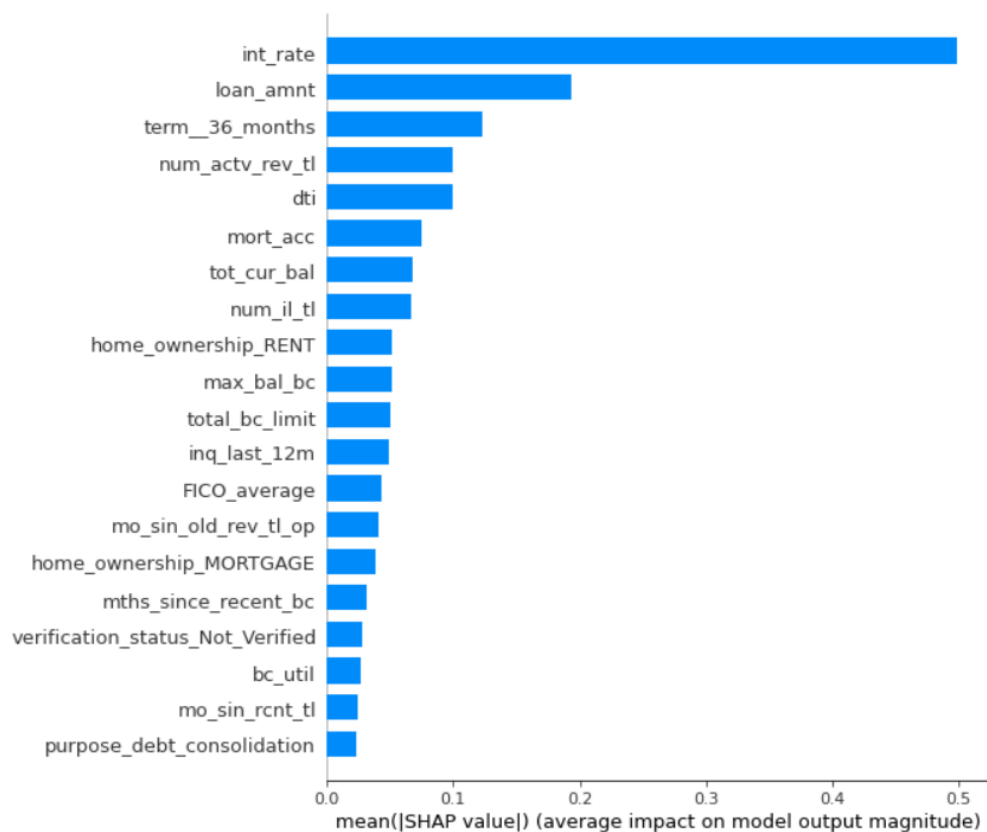
*Source:* Tianqi Chen
and Carlos Guestrin.
2016.



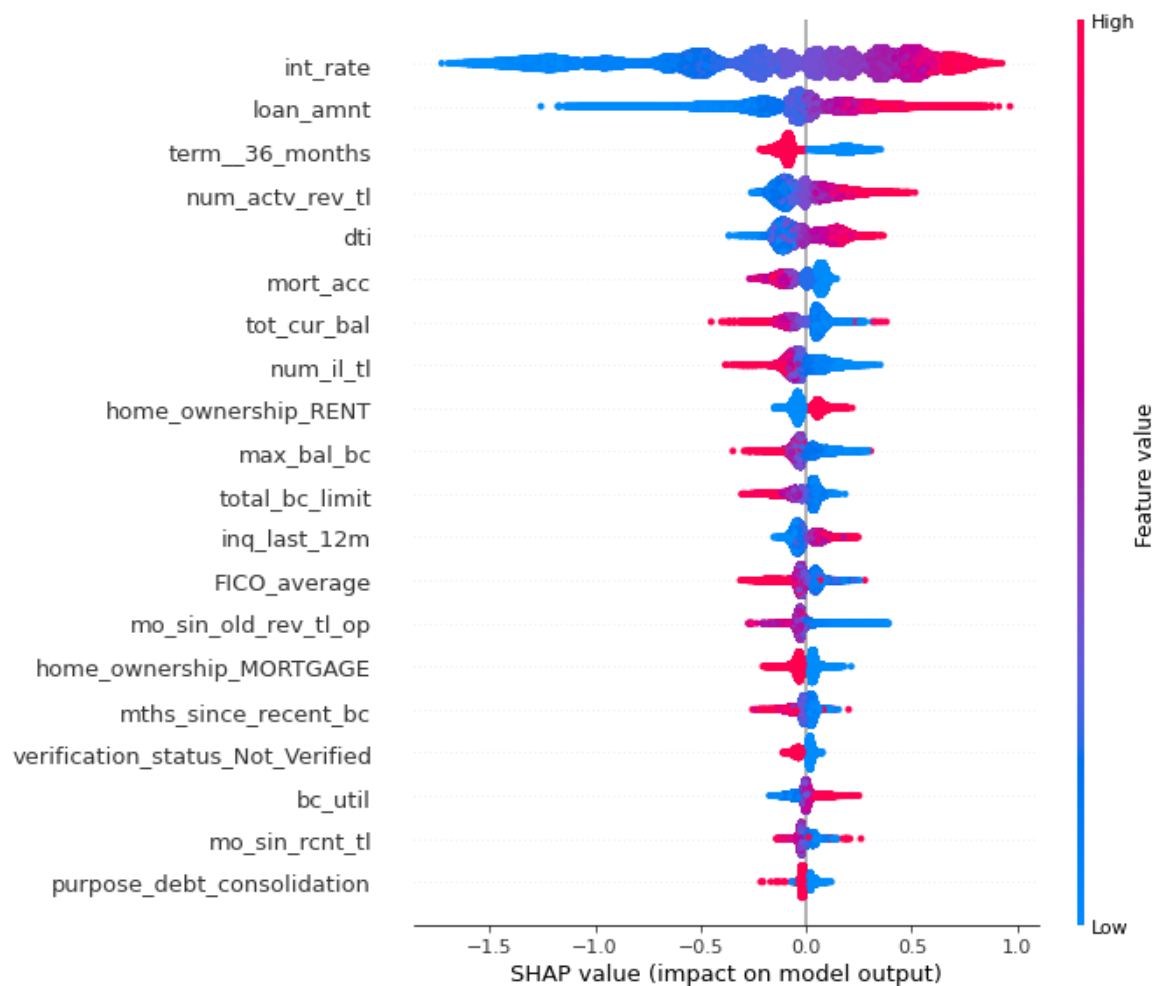**Figure 2.1**: An example tree showing a node split

# Explanation of the model
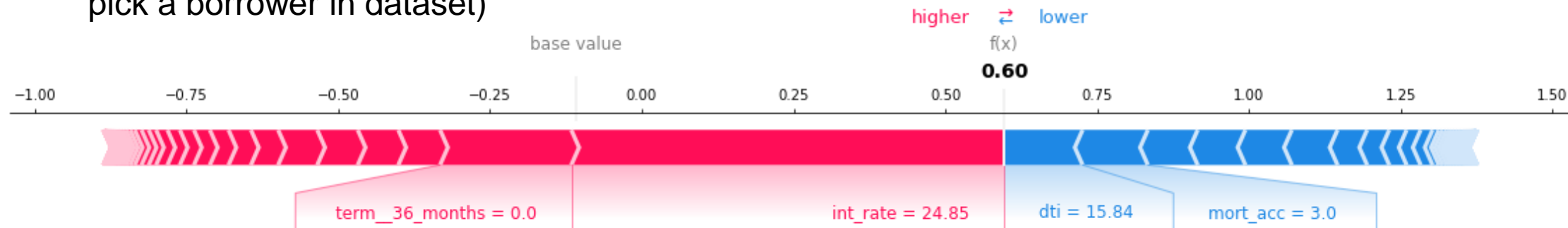
## Feature Importance based on Shap values

### Bar chart

### Density Scatter plot

OTTO VON GUERICKE
UNIVERSITÄT
MAGDEBURG

IFM

LEHRSTUHL FÜR
INNOVATIONS- UND
FINANZMANAGEMENT

Financial
Technology
ego.-INKUBATOR

# Explanation of the model

Local explanation using SHAP values (Lundberg & Lee, 2017) of the 9th borrower (randomly pick a borrower in dataset)



- Ouput value: 0.6 – convert to probability is 0.645 - Charged Off
- Base value
- Color (Blue/Red)
- Manitude of each feature

$$\frac{1}{1 + e^{-y^{\wedge}}}$$

Source: Daniel & Martin (2021). Speech and Language Processing. https://web.stanford.edu/~jurafsky/slp3/5.pdf P.3)

# Improvement

- Increase model performance by:

  o Training on larger subsample of data ( due to hardware

  capacity, only 700000 rows are used)

  o Further EDA with informative features – more distinguishing

  features

  o Tunning hyperparameters

- End-users: need knowledge of Machine Learning to read the

dashboard

THANK YOU!

Q&A