# Home Credit Default Risk

Phase 3, Group 9

Rakesh Jothishankar, Joel Klein, Suriyadeepan Narayanasamy, & Ashley Thornton

INDIANA UNIVERSITY

# Contents

INDIANA UNIVERSITY

# Overview

1. Goal: Improve method of approving or declining loan applications.

2. Data: Home Credit data from Kaggle.

3. Methods: Logistic regression, XGBoost, Light GBM, Random Forest, SVM, & Neural Network.
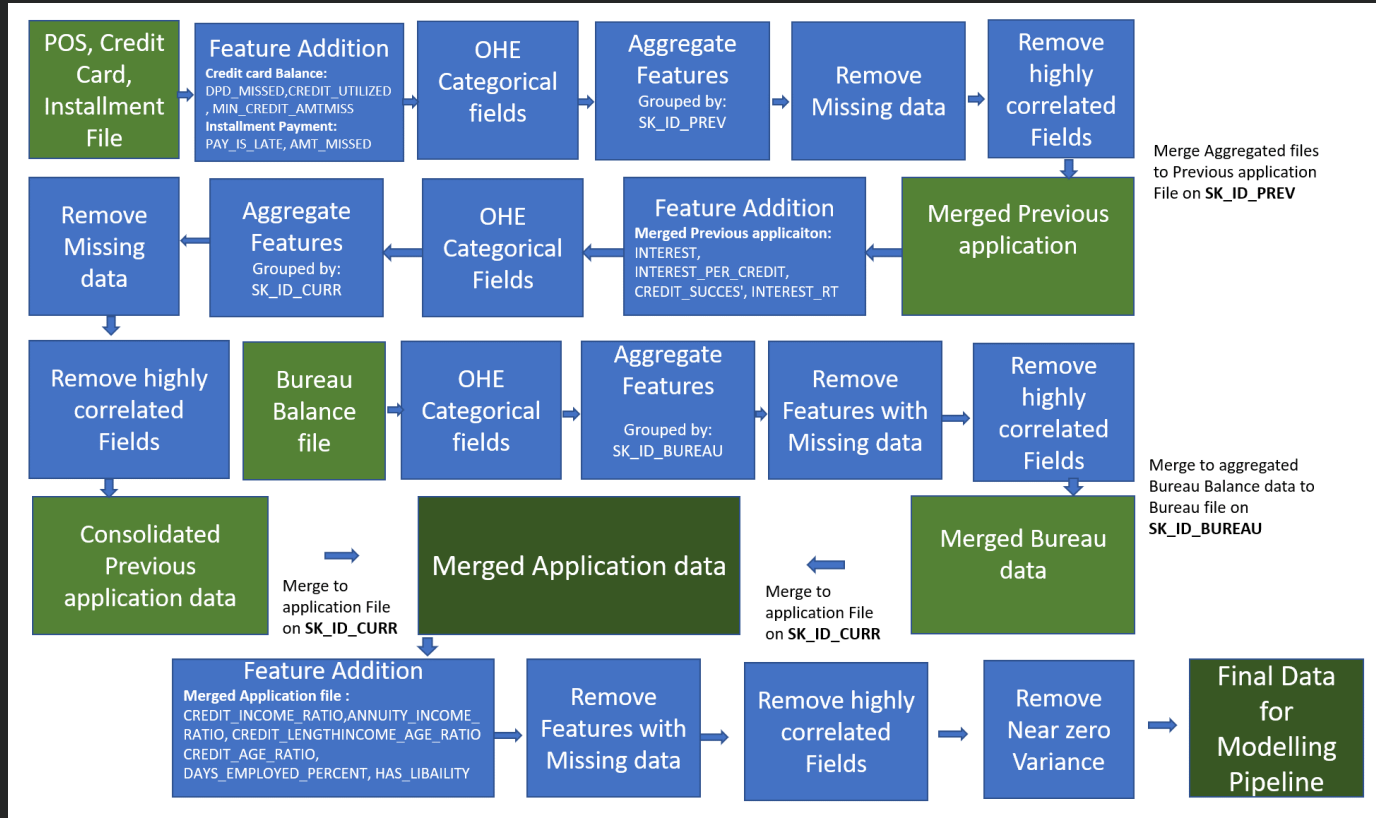
# Data Prep

1. POS cash balance, installment payment, and credit card balance files get rolled up to the previous applications file joined by SK_ID_PREV.

2. Bureau balance file gets rolled up to the bureau file joined by SK_ID_BUREAU.

3. Joined previous application data and joined bureau data get rolled up to the applications file on SK_ID_CURR.

4. Features with a large amount of missing data or highly correlated to other features were removed.

# Data Prep

# Feature List

**Feature Types**

| Loan (6) | Surrounding DPD (4) | Previous Application (36) |
|---|---|---|
| Date (5) | Document forms (20) | Previous Monthly POS/Cash Balance (6) |
| Contact info (6) | Credit bureau inquiries (6) | Previous Loan Installment Payments (6) |
| Family (3) | Demographics (3) | Previous Monthly Credit Card Loan Balance (21) |
| Region (9) | Occupation (2) | Bureau Previous Credits (15) |
| External (3) | Process Time (2) | Bureau Previous Credits Monthly Balance (2) |
| Housing (48) | Other Assets (3) | |

# Exploratory Data Analysis

Features with Most Missing Data



Target Frequency Distribution

# Data Handling Pipeline

1. Imputed missing numeric values with median.

2. Standardized numeric features.

3. Imputed missing categorical values with "Unknown."

4. OHE categorical features.

5. Feature engineered new features.

6. Removed near zero variance features.

7. Removed features with zero importance from previous model (for some test runs).

**New Engineered Features**:
- Late payment
- Amount missed
- Credit utilized
- Min credit amount missed
- Interest
- Interest per credit
- Credit success
- Interest rate
- Credit to income ratio
- Annuity to income ratio
- Credit length
- Income to age ratio
- Credit to age ratio
- Percent of days employed
- Liability

# Sampling Method

# Modeling Pipeline

# Artificial Neural Network Visualization

# Best Performing Model: Ensemble

**Artificial Neural Network**

-Batch size: 10,000
-Epochs: 15
-Learning rate: 0.001

**+**

**XGBoost**

-NZV features removed
-Learning rate: 0.1
-Max depth: 3
-Trees: 300

**+**

**Light GBM**

-Feature selection
-Boosting type: dart
-Learning rate: 0.005
-Max depth: 2
-Trees: 7000

# Results

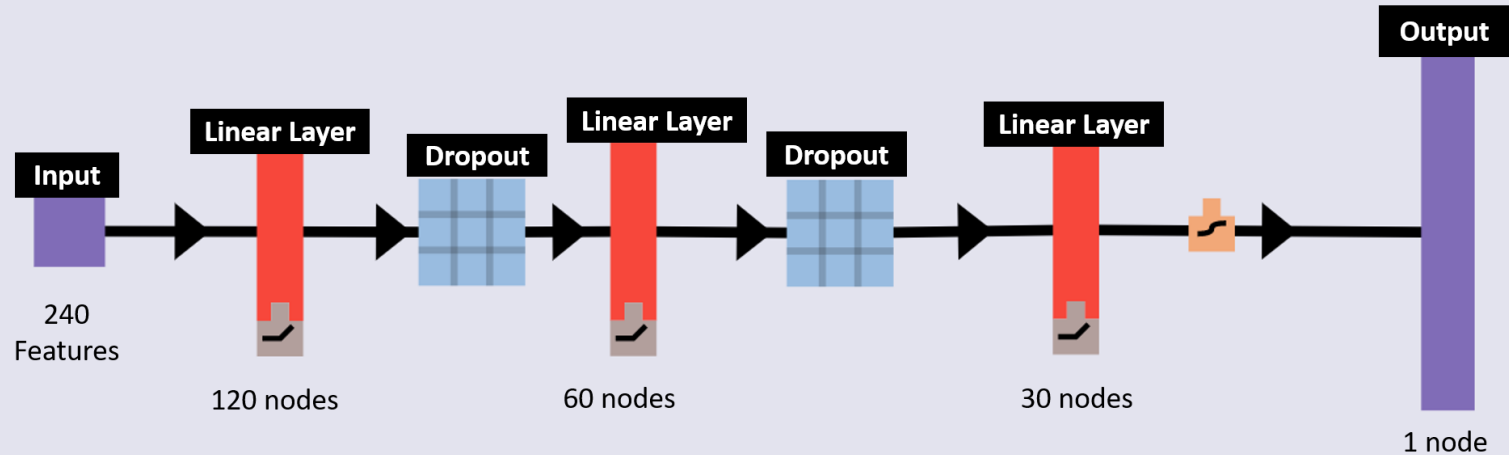| Model | Experiment | Train Accuracy | Test Accuracy | Train Area under ROC | Test Area under ROC | Best Parameters |
|---|---|---|---|---|---|---|
| Stacked ANN + XGBoost + LightGBM | App, agg prev app & bal, agg bureau & bal data w/ NZV features | | | 78.44 | 78.79 | |
| Stacked ANN + LightGBM | App, agg prev app & bal, agg bureau & bal data w/ NZV features | | | 78.29 | 78.63 | |
| XGBoost | App, agg prev app & bal, agg bureau & bal data w/ feature selection | 91.99 | 92.02 | 78.27 | 78.67 | {'xgb__colsample_bytree': 0.1, 'xgb__learning_rate': 0.1, 'xgb__max_depth': 3, 'xgb__n_estimators': 300} |
| XGBoost | App, agg prev app & bal, agg bureau & bal data w/ NZV features | 91.98 | 92.03 | 78.27 | 78.7 | {'xgb__colsample_bytree': 0.1, 'xgb__learning_rate': 0.1, 'xgb__max_depth': 3, 'xgb__n_estimators': 300} |
| Stacked ANN + XGBoost | App, agg prev app & bal, agg bureau & bal data w/ NZV features | | | 78.24 | 78.62 | |
| LightGBM | App, agg prev app & bal, agg bureau & bal data w/ NZV features | 91.98 | 92.01 | 78.03 | 78.48 | {'lgbm__boosting_type': 'goss', 'lgbm__colsample_bytree': 0.1, 'lgbm__learning_rate': 0.005, 'lgbm__max_depth': 2, 'lgbm__n_estimators': 7000} |
| LightGBM | App, agg prev app & bal, agg bureau & bal data w/ feature selection | 91.99 | 92.01 | 78.02 | 78.48 | {'lgbm__boosting_type': 'goss', 'lgbm__colsample_bytree': 0.1, 'lgbm__learning_rate': 0.005, 'lgbm__max_depth': 2, 'lgbm__n_estimators': 7000} |
| ANN | App, agg prev app & bal, agg bureau & bal data w/ NZV features | 91.94 | 91.89 | 77.17 | 77.6 | {'BATCH_SIZE' = 10000 'EPOCHS' = 15 'LEARNING_RATE' = 0.001} |
| Logistic Regression | App, agg prev app & bal, agg bureau & bal data w/ NZV features | 91.93 | 91.92 | 76.56 | 76.81 | {'logistic__C': 0.01, 'logistic__l1_ratio': 0.2} |
| XGBoost | App, agg prev app & bal, agg bureau & bal data w/ PCA | 91.93 | 91.95 | 75.26 | 75.36 | {'xgb__colsample_bytree': 0.2, 'xgb__learning_rate': 0.1, 'xgb__max_depth': 3, 'xgb__n_estimators': 300} |
| Logistic Regression (Baseline) | All application data features | 91.91 | 91.93 | 74.17 | 74.48 | {'logistic__C': 0.1, 'logistic__l1_ratio': 0.6} |

# Kaggle Submission



| Submission and Description | Private Score | Public Score | Use for Final Score |
|---|---|---|---|
| hcdr_kaggle_submission_phase3_ensemble (4).csv<br>a few seconds ago by Rjothis<br><br>Ensemble - XGB FS... ANN + XGB + DART | 0.78502 | 0.78628 | ☐ |

Place: 3,651 out of 7,176

# Phase 3 Issues

1. Size of data.

2. Sklearn is not optimized for training neural networks.

3. Additional experiments with SVM never completed.

# Conclusion



**Chorus**
– The gradient is the weighted sum of the training data, where the weights are proportional to the error (for each example) !

## Past

Phase 1: Defined project, performed EDA, and ran baseline model.

Phase 2: Improved model through hyperparameter tuning and additional algorithms.

## Present

Phase 3: Neural Network with PyTorch

## Future

Professional root finders!