# L&T Loan Default Prediction

Aleksa Radojičić

# Overview

- **Goal**: Predict the probability of Indian borrower defaulting on a vehicle loan in the first EMI (Equated Monthly Installments) on the due date ([competition](#))

- **Problem type**: Supervised binary classification (target = default on first EMI)

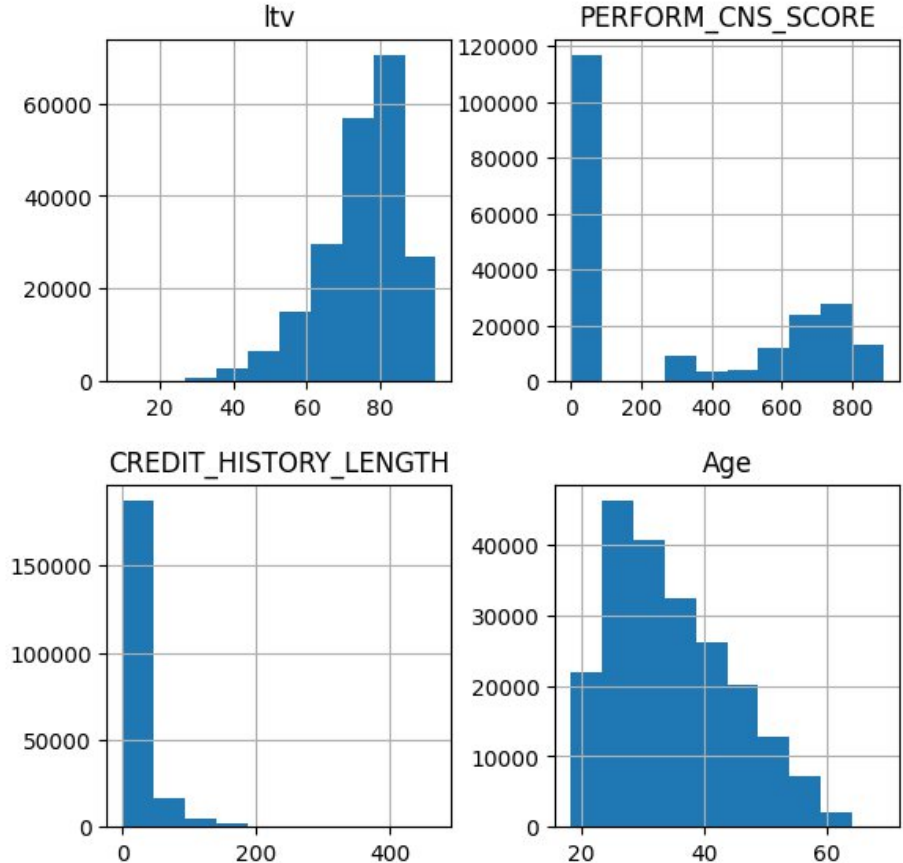- **Evaluation metric**: AUC (Area under the Curve)

# Data Description

- Train set rows: 233,154
- Test (submission) set rows: 112,392

- 39 features:
  `PERFORM_CNS.SCORE.DESCRIPTION`,
  `asset_cost`,
  `disbursed_amount`,
  `State_ID`,
  `PRI.NO.OF.ACCTS`,
  `NO.OF_INQUIRIES`, ...

- Label: `loan_default`

Split 10% of original train set for testing:
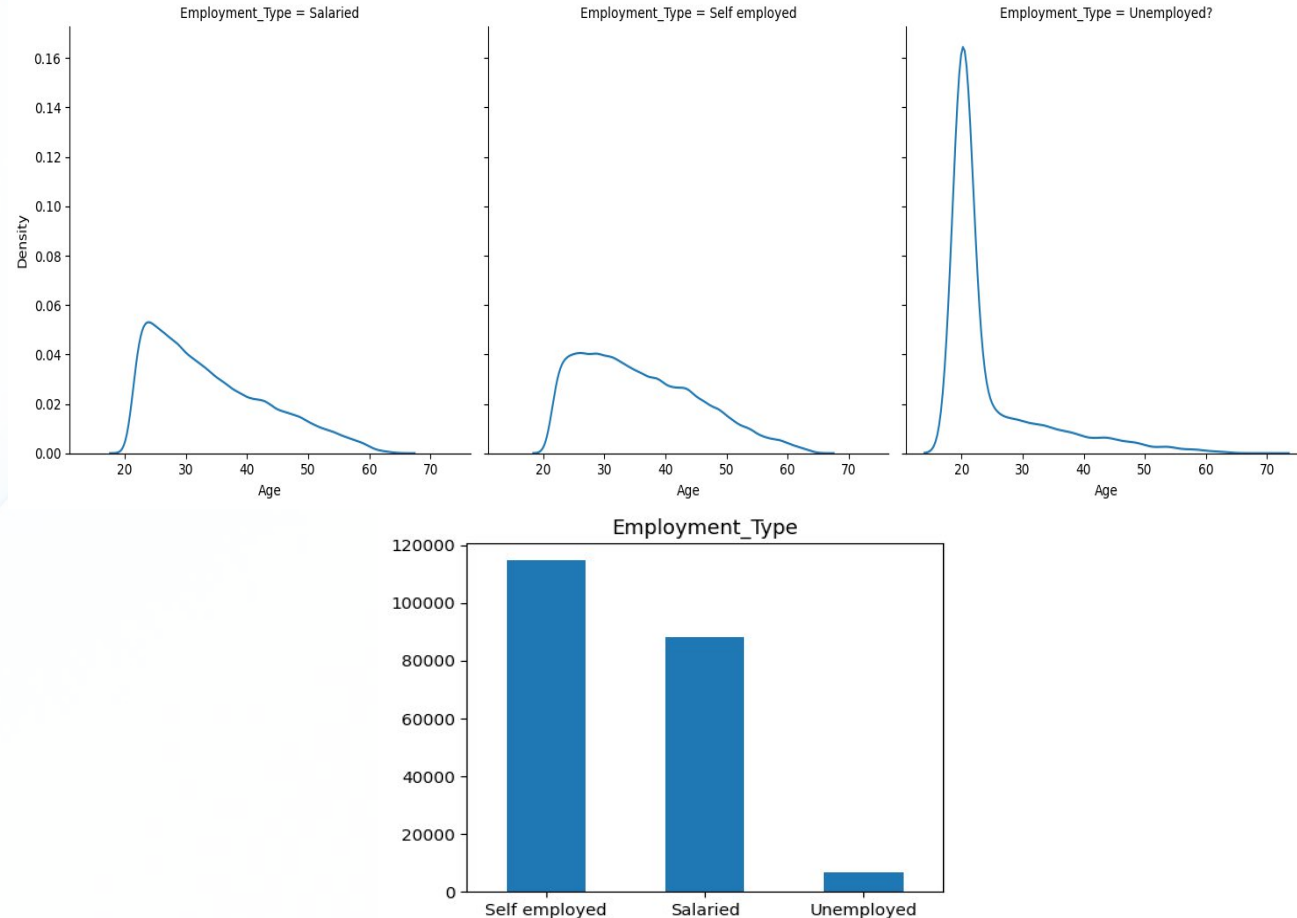- Train set: 209,835
- Test set: 23,316

# EDA

- Entire EDA performed on training data subset.

- Mean LTV (loan to value): 76.81.

- ~50.15% do not have bureau history available, meaning no CNS Score either.

- Median age: 26.

- For more than 50% of the population this was the first loan.

# Preprocessing

- NA values were only in `Employment_Type` column. After examining relationship of `Age` to `Employment_Type`, there is a high number of young people in NaNs, considerably more than in other groups. Therefore, NaNs are filled with 'Unemployed' value.

- Removed rows where `PERFORM_CNS_SCORE_DESCRIPTION` = 'Not Scored: More than 50 active Accounts found' (3 entries).

- Removed columns for high correlation with others (>0.9): [`PRI_CURRENT_BALANCE`, `SEC_CURRENT_BALANCE`, `PRI_DISBURSED_AMOUNT`, `SEC_DISBURSED_AMOUNT`]
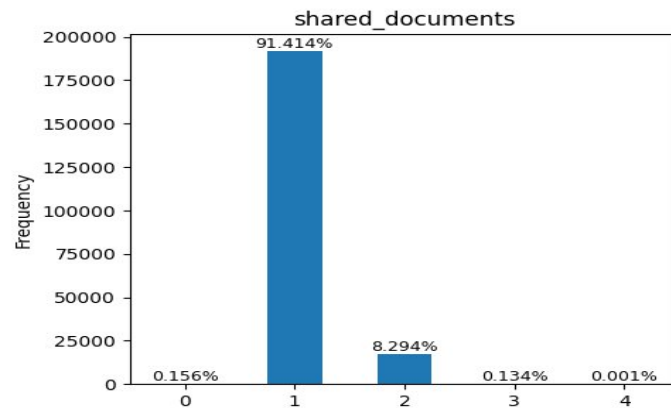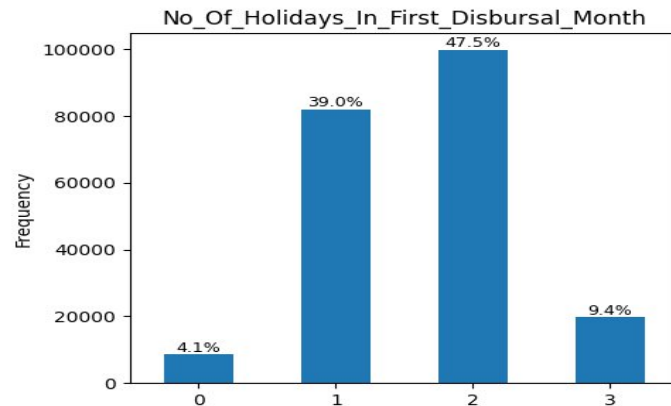
# Feature Engineering – Numerical Features

- Combined primary and secondary columns into total aggregates (prefix 'TOT').

- Add other numerical columns:
  [`PRI_OVERDUE_TO_ACTIVE_ACCTS_RATIO`,
  `SEC_OVERDUE_TO_ACTIVE_ACCTS_RATIO`,
  `TOT_OVERDUE_TO_ACTIVE_ACCTS_RATIO`,
  `PRI_ACTIVE_ACCTS_RATIO`,
  `SEC_ACTIVE_ACCTS_RATIO`,
  `TOT_ACTIVE_ACCTS_RATIO`,
  `DELINQUENT_TO_NEW_ACCTS_RATIO`,
  `DELINQUENT_TO_ALL_ACCTS_RATIO`,
  `PERFORM_CNS_SCORE_NORMALIZED_BY_ltv`,
  `Invalid_Age_First_Loan`].

# Feature Engineering – Categorical Features

- Added categorical columns:
  `Invalid_Age_First_Loan`
  `No_Of_Holidays_In_First_Disbursal_Month`,
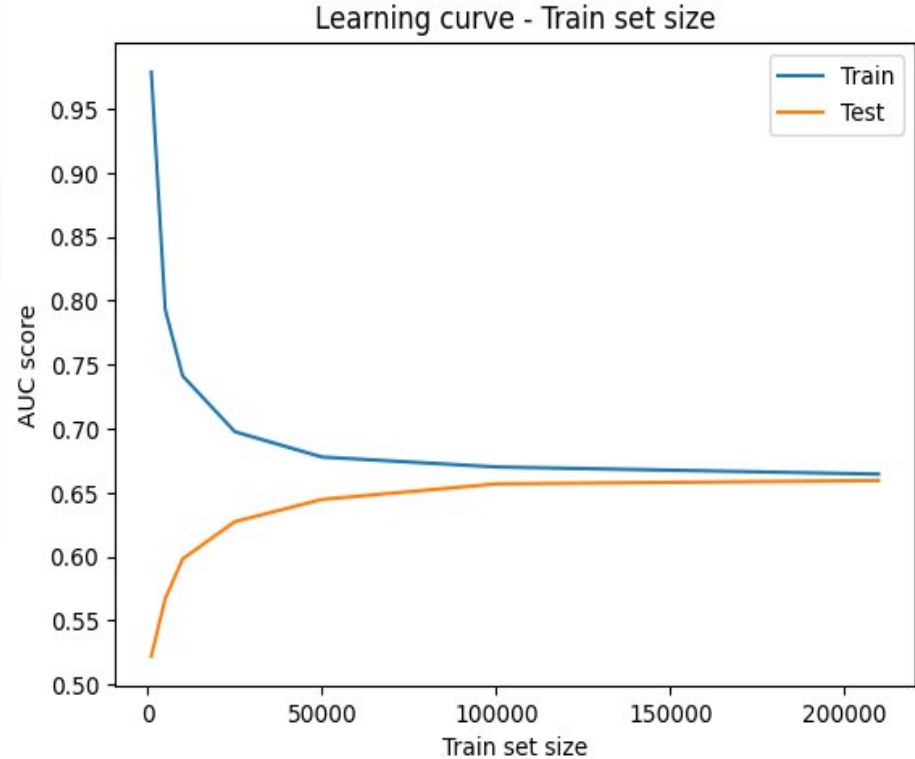  `shared_documents`,
  `Month_of_Birth`



No_Of_Holidays_In_First_Disbursal_Month



shared_documents

Modelling

# Learning Curve – Train Set Size

Train set sizes:
[1,000; 5,000; 10,000; 25,000; 50,000;
100,000; 209,835]

**<u>Conclusion</u>:**
Adding more data pass 200,000 set size
won't improve test AUC (high bias), implying
stronger model is required for better
performance.
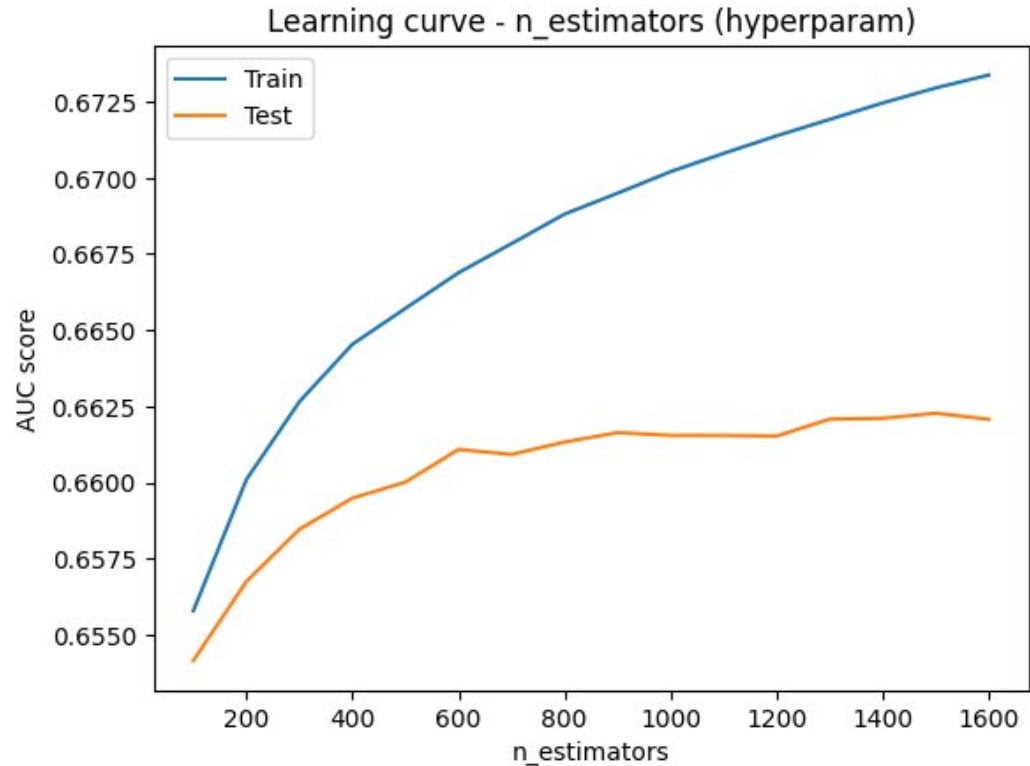


Learning curve - Train set size

# Learning Curve – `n_estimators`

Train AdaBoost model with varying values of hyperparameter `n_estimators` from 100 to 1600 to estimate gain.

**Conclusion**:
Increasing the number of estimators beyond 600 yields only marginal improvements in test AUC, suggesting that further gains may require tuning other hyperparameters, such as the base tree depth or enhancing feature engineering.

# Hyperparameters Tuning pt. 1 – BayesSearchCV

- Total features after feature engineering: 44

- Pipeline:

1. Preprocessing pipeline

2. Feature selection using `mutual_info_classif`

3. AdaBoost model

- BayesSearch stratified 5-fold cross validation with 30 iterations, with hyperparameter space:

```python
param_space: dict[str, Any] = {
    f'{fs_name}__k': space.Integer(low=13, high=44),
    f'{classifier_name}__learning_rate': space.Real(low=10**-3, high=10**0, prior="log-uniform"),
    f'{classifier_name}__n_estimators': space.Integer(low=900, high=1700),
    f'{classifier_name}__minority_class_weight': space.Real(low=1, high=5)
}
```

# Hyperparameters Tuning pt. 2 – GridSearchCV

- GridSearch stratified 5-fold cross validation for the split criterion hyperparameter: ["gini", "entropy"]

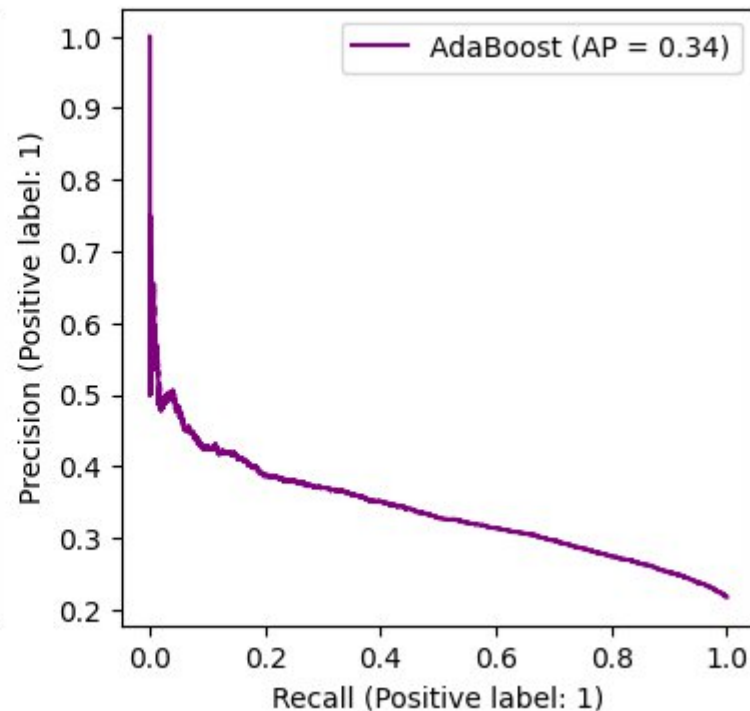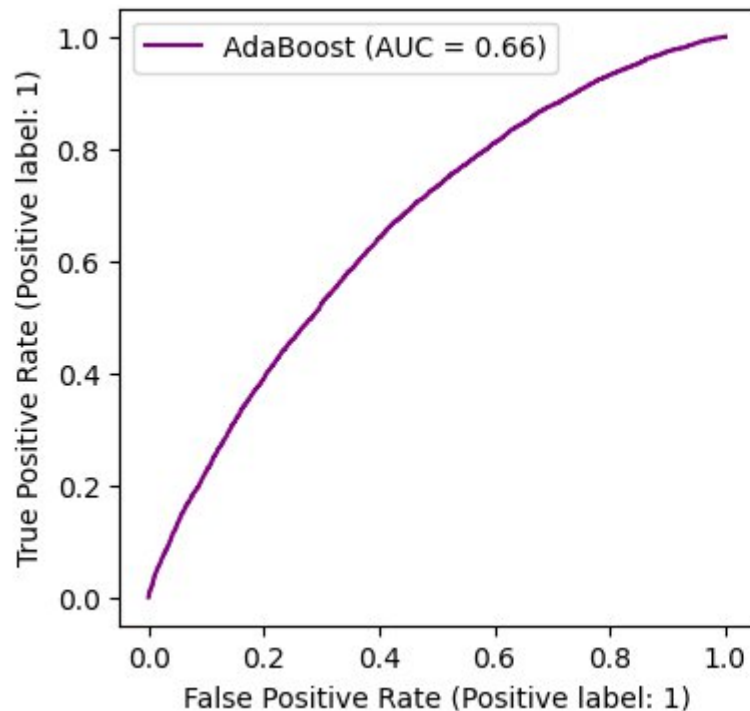- Refit the best model on the entire train data.

# Results

## Optimal Hyperparameters

| Hyperparameter | Value |
|---|---|
| AdaBoost__learning_rate | 0.8026 |
| AdaBoost__minority_class_weight | 5.0 |
| AdaBoost__n_estimators | 1700 |
| AdaBoost__estimator__criterion | gini |
| mutualinfoclassif__k | 42 |

## Evaluation

| Dataset | AUC |
|---|---|
| CV | 65.888% |
| Train Set | 67.144% |
| **Test Set** | **66.302%** |
| Entire Set | 67.069% |

# ROC & PR curve

# Feature Importance



AdaBoost: Feature importances