# BFSI: CREDIT RISK ASSIGNMENT

By: Shreevatsa Hegde, Shraddha Nerlekar, Shruthi J

# Contents

- Objective
- Background
- Data Analysis
  - Pre Processing of Data
  - EDA
  - Model Building
  - Interpreting the Results
  - Recommendations

# OBJECTIVE

- The objective is to build a statistical model to estimate borrowers' **Loss Given Default (LGD)**

$$LGD = \frac{Loan\ Amount - (Collateral\ value + Sum\ of\ Repayments)}{Loan\_Amount}$$

# BACKGROUND

- Credit risk analytics in the context of the banking sector and model a common metric used for estimating the expected credit loss (ECL)

- ECL method is used for provisioning the capital buffer to protect banks against possible default of the customers.

**Expected credit loss = Exposure at default x Probability of Default x Loss given default**

- The **loss given default (LGD)** is a measure of the amount of loss that a bank is expected to incur in the event of a default by a borrower.
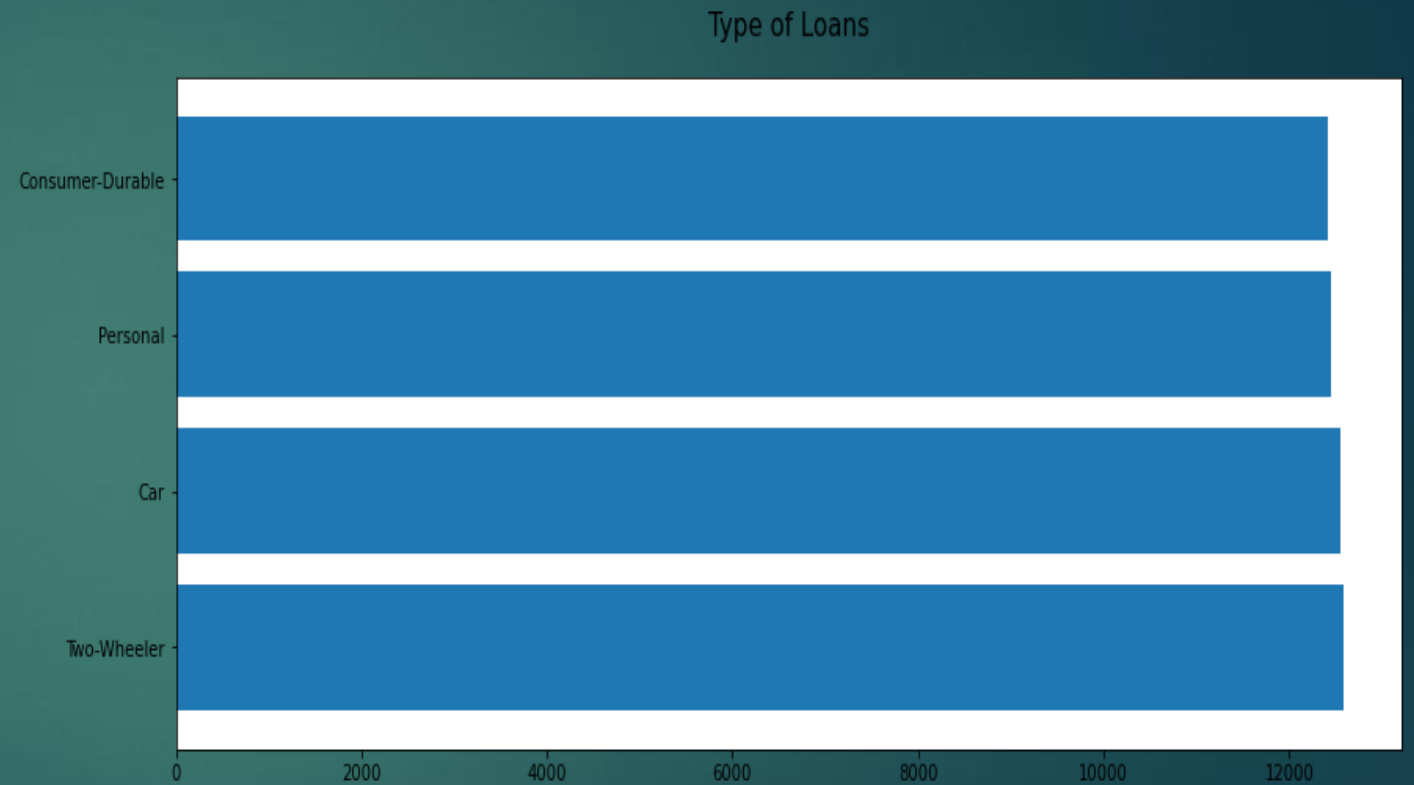
# DATA SOURCES

- Used 3 Data sets for model Building
  - The main_loan_base data set contains information about loan accounts and other relevant information for the corresponding borrowers.
  - The repayment_base data set contains information about the repayments received by the banks in the form of EMIs or through other collection efforts.\
  - The monthly_balance_base contains the information pertaining to the monthly balance statements in the borrower's accounts.

# PRE PROCESSING OF DATA

- For each data set converted Data types if necessary

- Null values are handled using deletion and imputation techniques. As well duplicate values are removed from data sets.

- Merging the data sets and created target variable(LGD)

- Exploratory Data Analysis has been performed

- Variable Transformation
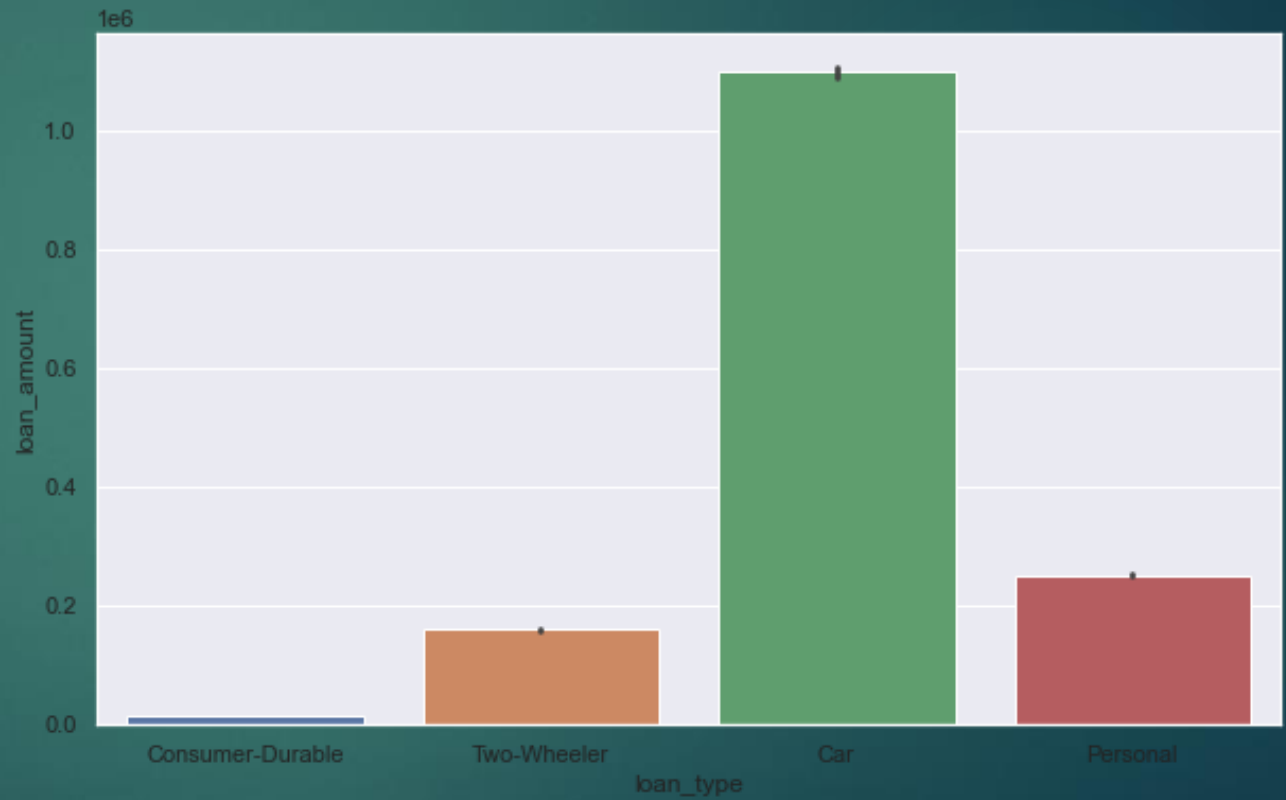
- Dummy Encoding

- Scaling using Standard Scaler

# EDA

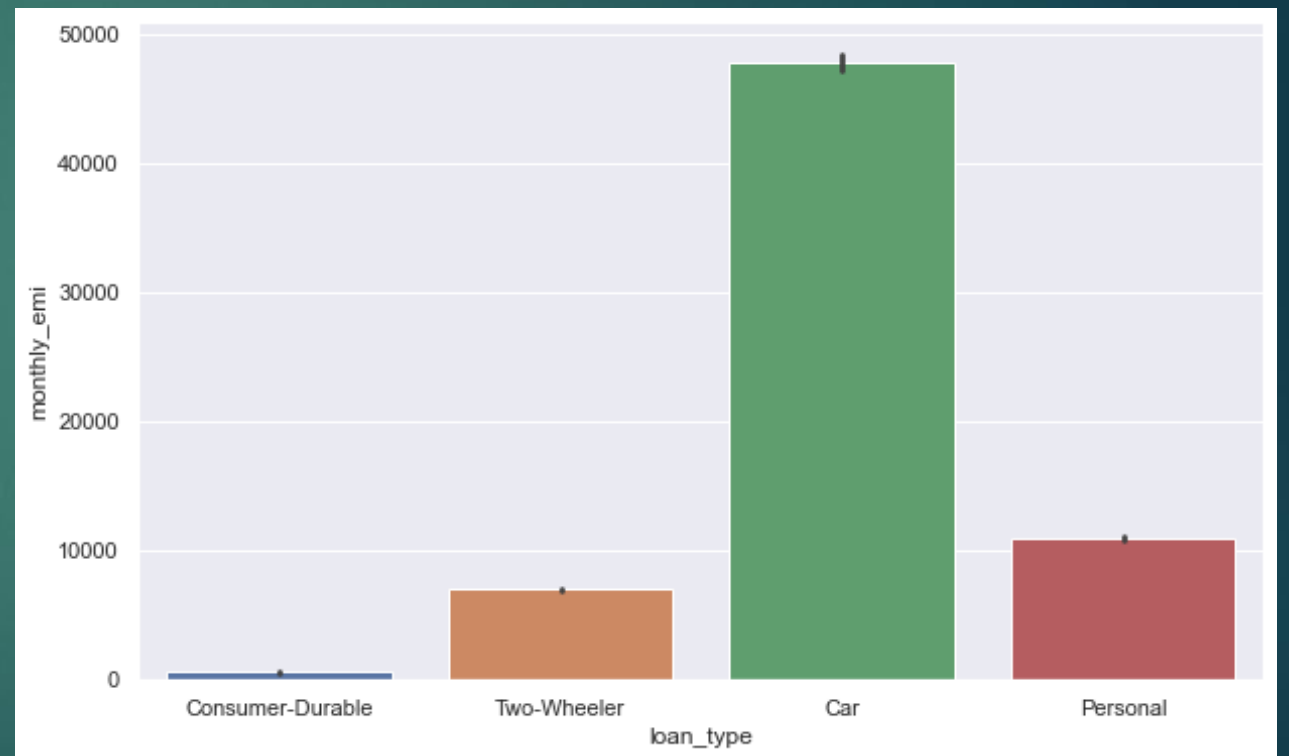Number of loans in Two-wheeler is higher than all others.



Type of Loans

# EDA

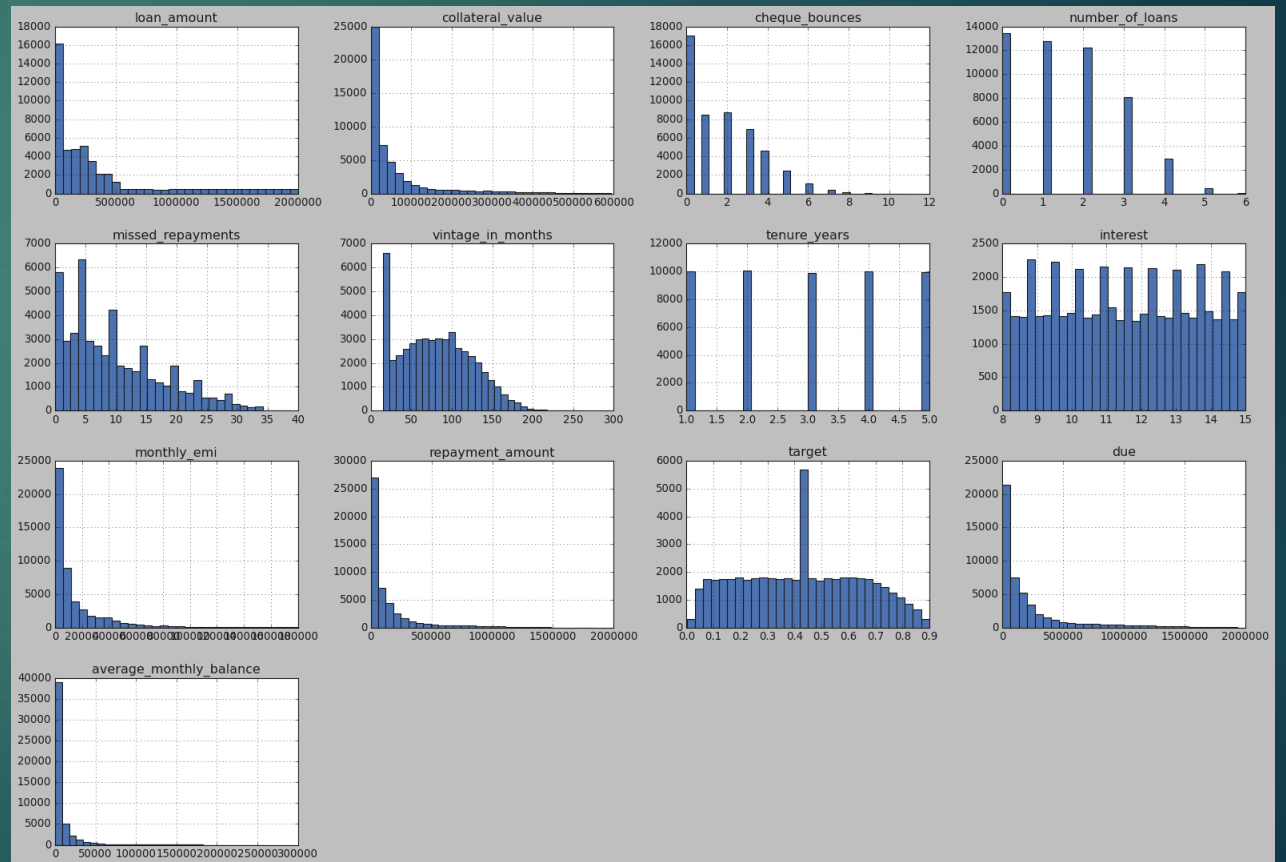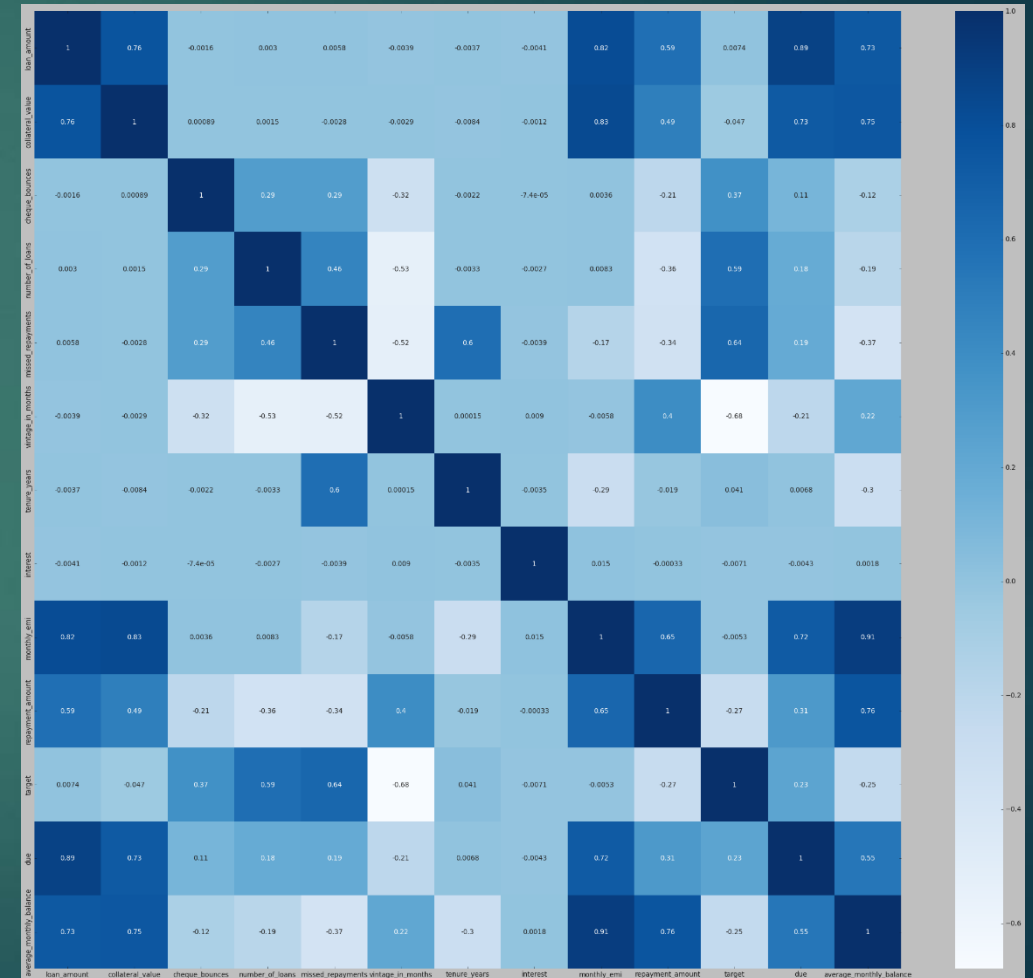But, the loan amount of car loan is the highest.

# EDA

Monthly EMI also car loan is much higher compared to other loans.

Plotted histograms for the numerical columns to understand the distribution of data

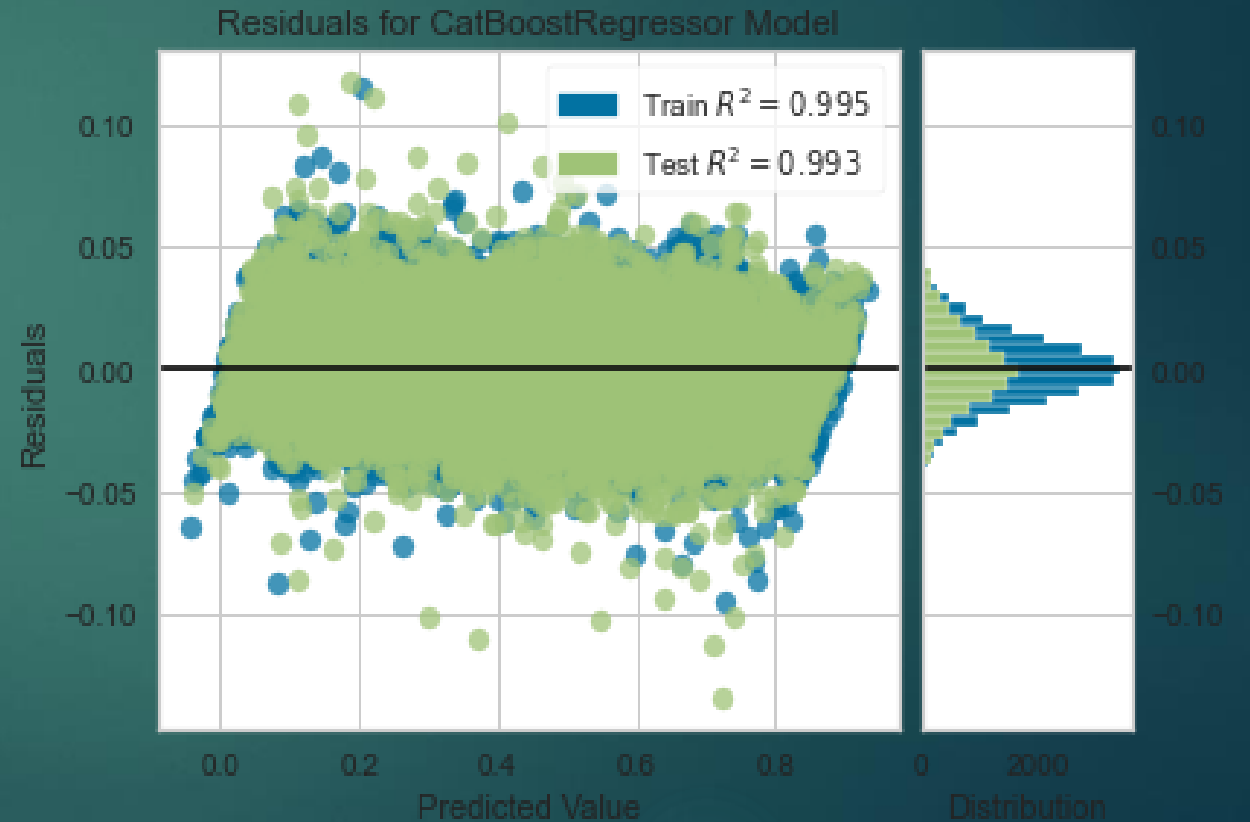Created Heatmap to understand the correlation between the variables

- used Power transformation to make numerical variables Normally distributed

- Dropped unnecessary columns for model building

- Used One-hot encoding technique and created dummy variable for necessary categorical variables.

# MODEL BUILDING

- Used various models like Multiple Linear Regression, Random Forest Regressor, Gradient Boosting Regressor, XGBoost Regressor, Adaboost Regressor, ElasticNet :Hybrid Regularized Model, LightGBM for model building.

- Used R Squared as a performance metrics.

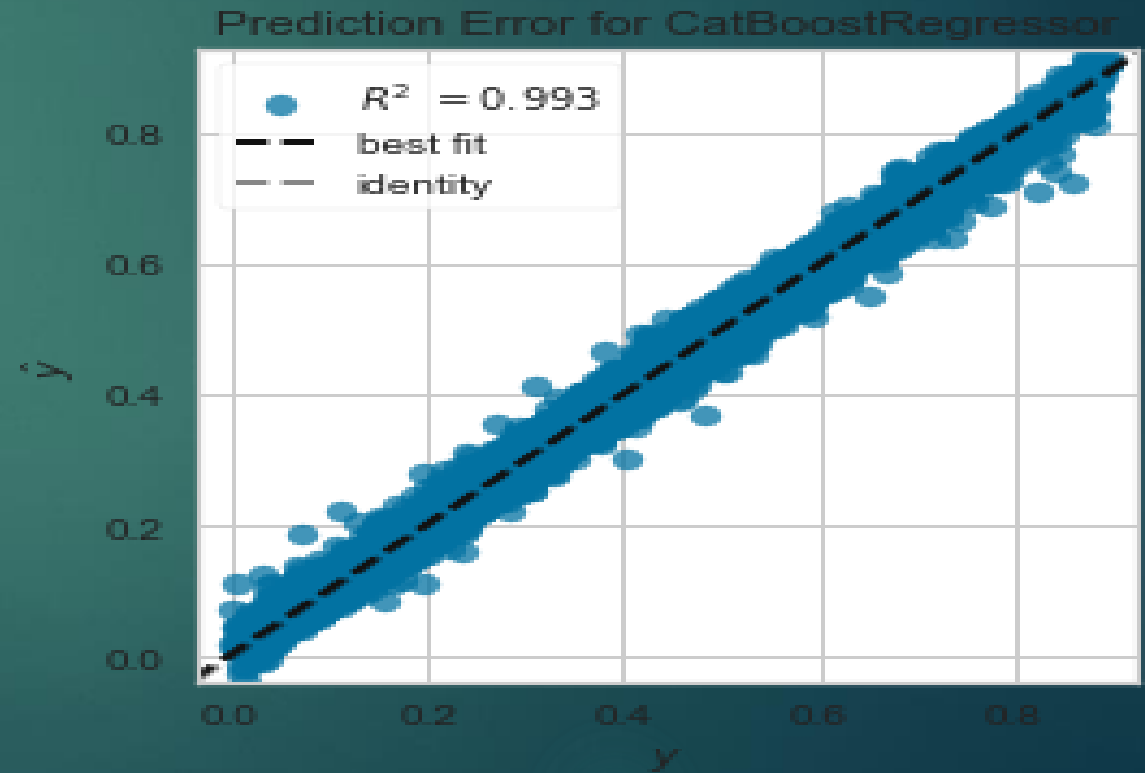- XGBoost has given us 99.5% R squared on test data across the models.

# REGRESSION INTERPRETATION
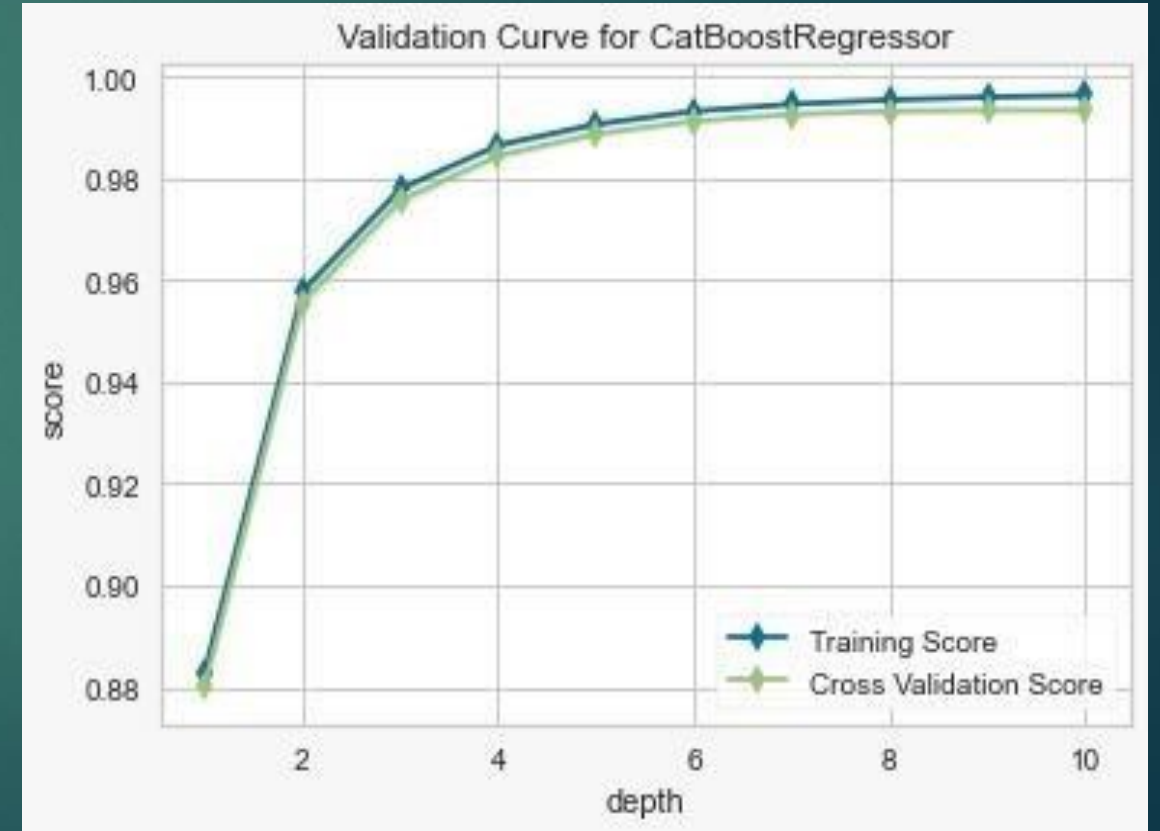


- Residual plot of the finest model

# REGRESSION INTERPRETATION
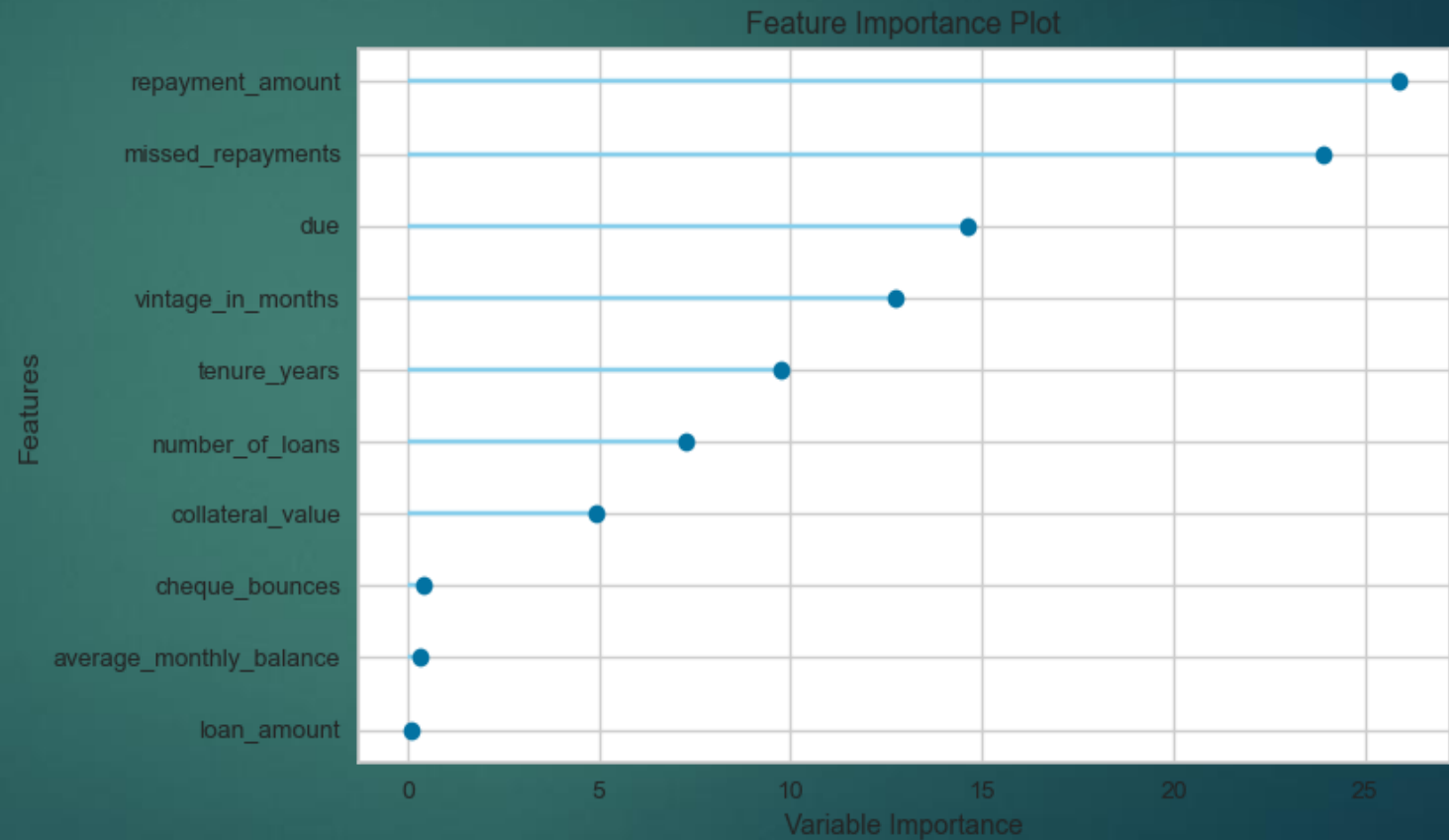
- Best fit line corresponding the prediction error



Prediction Error for CatBoostRegressor

# REGRESSION INTERPRETATION

- Validation Curve

# REGRESSION INTERPRETATION

- Feature Importance


Feature Importance Plot

# RECOMMENDATIONS

1. We should focus more on Car and Two-wheeler loan types
2. Missed Repayment customers with high repayment amount should be highlighted
3. Customer's due factors and tenure are another subset of influencers to predict the Loss Given Default of the customers.

# THANK YOU