# Building an Ensemble Model to predict personal loan interest rates

## Data Preparation

Amount Requested: Continuous, Income: Continuous, Rate: Continuous, Months: Continuous

IL $ ILL: Grouped into IL (Illinois)

TEX $ TX: Grouped into TX(Texas)

WI $ WISC: Grouped into WI (Wisconsin)

Imputations

The missing values in column „Utilization" are imputed by calculating Credit balance/ Credit limit. As the credit balance is 0 for all the missing values in „Utilization", these missing values are imputed with 0.00%. Missing values in column "Month" are imputed with Multivariate Normal Imputation.

Monthly Pmt, Amt Funded and Standing variables have not been used for modeling the data.

Transformations

Z ln - Amt Requested, z ln - Pmts to date, z ln - Income, z ln - Open LoC, zln - Total LoC, z sqrt - Balance, z sqrt - Delinquent, z - Utilization, z – Months, z -D/I

## Exploratory Data Analysis

The variables Amt Requested, Pmts to date, Open LoC, Total LoC and Income are right skewed. Rate and FICO Score can be inversely proportional. Rate is observed to be high when there is zero Utilization. The interest rates for small businesses and consolidated debts are observed to be high. FICO Score has an inverse relation with Delinquencies. Greater the FICO score, lesser are the number of delinquencies. Income, State and Employment in years may not have a direct impact on the interest rates. The minimum interest rate is 5.42 % for applicants who have a FICO score of greater than 781. The maximum interest rate is 24.8 % for applicants who have a FICO score between 660-680.
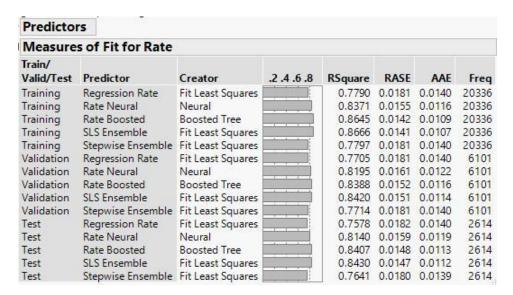
**Fitting the model**

The chosen model for predictions uses Ensemble method (Regression/Standard least squares).

The 4 Predictor variables used in the final model are output variables of 4 base models of

regression, k-NN, Boosted tree and Neural net, which are denoted as Regression Rate, Rate

Neural, Rate Boosted, Rate kNN in the model. Output variable is Rate.

**Boosted Tree model & Neural net model specifications** -



A screenshot of all generated models along with their training and validation R-square values is

given below -



**Standard Least squares Ensemble** : Validation R square - 0.842, Training R squared - 0.8666