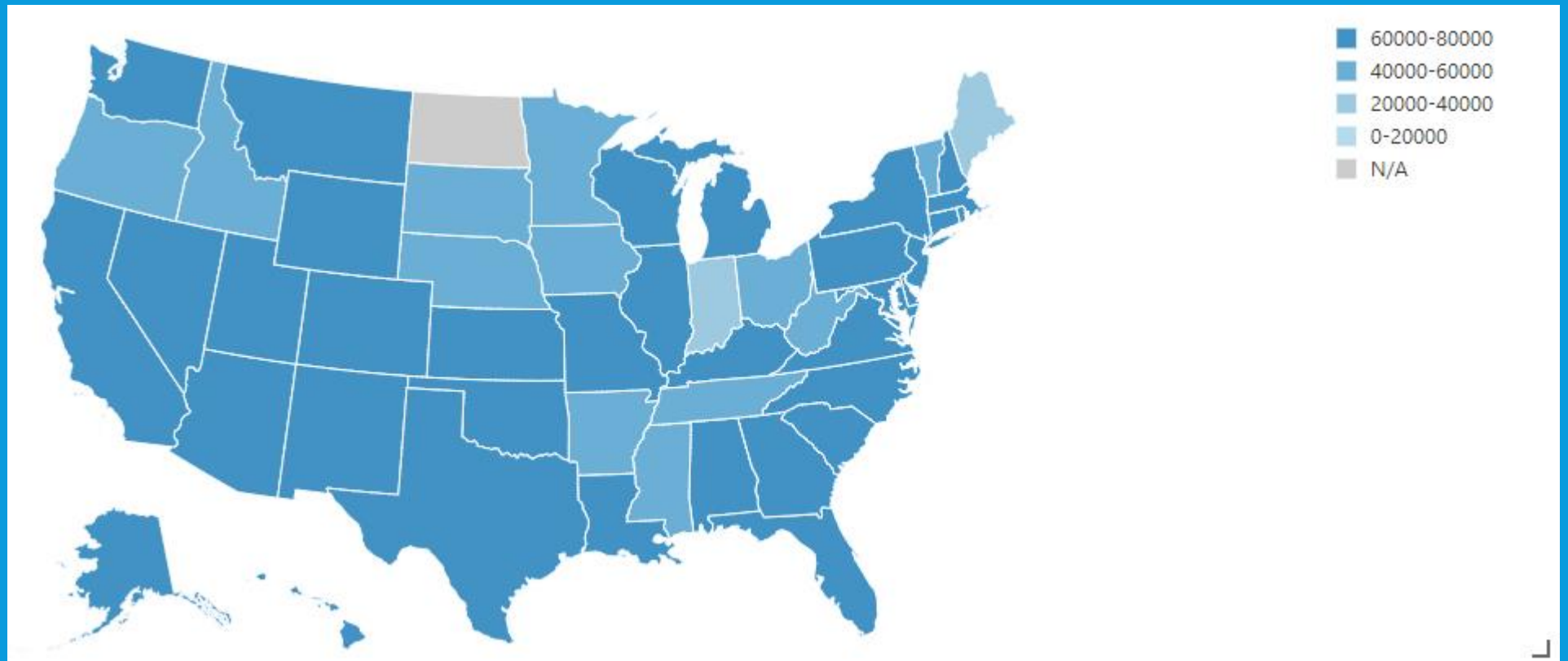# Evaluating Risk for Loan Approvals on LendingClub

The self-reported annual income provided by the borrower during registration.

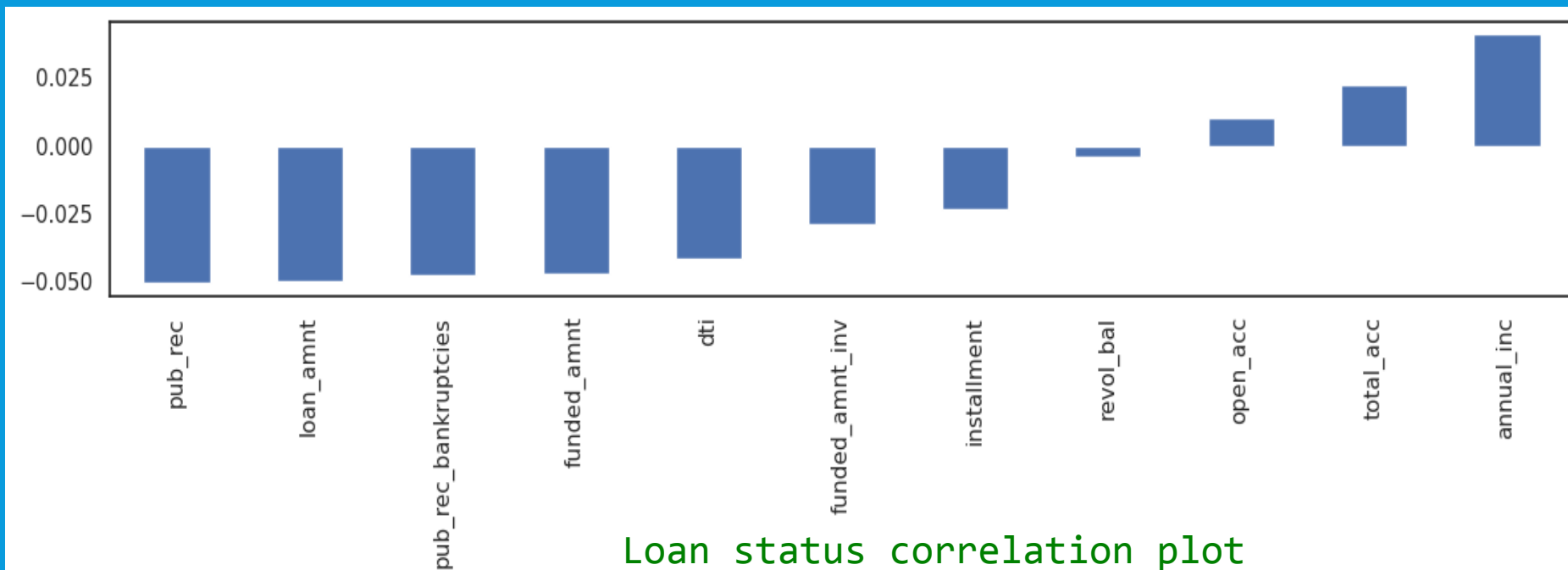The number of open credit lines in the borrower's credit file.
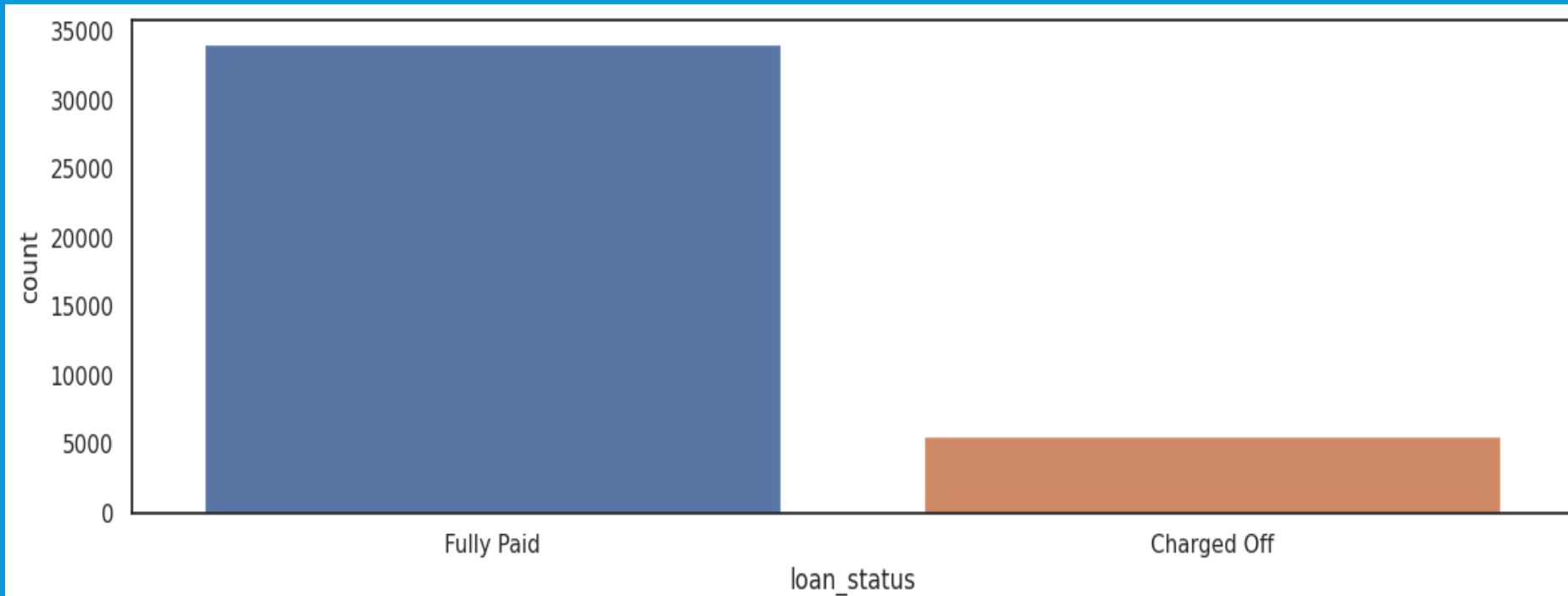
The total number of credit lines currently in the borrower's credit file

Total credit revolving balance

**Debt-to-income ratio**

**Loan amount granted**

**Number of derogatory public records**

Loan status correlation plot

**High DTI translates into higher default rates**

A ratio calculated using the borrower's total monthly debt payments



**Higher the installment amount, higher the default rate**

The monthly payment owed by the borrower if the loan originates.

**60 months term loan are defaulted more often**

The number of payments on the loan. Values are in months and can be either 36 or 60.

**Verified accounts defaulted in large number**

Indicates if income was verified by LC, not verified, or if the income source was verified

LC assigned loan grade

**The 'E' and 'F' grade loans are defaulted more often.**



LC assigned loan subgrade

**'F' and 'G' sub-grades don't get paid back that often**

# Better to avoid small business loans followed by educational loans

The CA, NY, TX states had high number of applications and high default rate

# Machine Learning based Loan Default Predictions

# Modeling Binary Classifiers

Five binary classifiers have been modeled namely, Linear SVC, Logistic Regression, Gaussian NB, Random Forest Classier, Gradient Boosting Classifier and XGBClassifier

Chosen Recall, Precision, and F1-score as evaluation metrics.

The precision is the measure of how accurate the classifier's prediction of a specific class.

The Recall is the measure of the classifier's ability to identify a class.

Resampling (Oversampling)

This technique is used to upsample the minority class of an imbalanced dataset using replacement. This technique is called oversampling.

Synthetic Minority Oversampling Technique (SMOTE)

SMOTE is another technique to oversample the minority class. It looks into minority class instances and uses k nearest neighbor to pick a random nearest neighbor, and a synthetic instance is created randomly in feature space.

# Classification Metrics without application of Data Imbalance handling Techniques

**Linear SVC**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0.0 | 0.33 | 0.00 | 0.00 | 534 |
| 1.0 | 0.87 | 1.00 | 0.93 | 3440 |
| accuracy |  |  | 0.87 | 3974 |
| macro avg | 0.60 | 0.50 | 0.47 | 3974 |
| weighted avg | 0.79 | 0.87 | 0.80 | 3974 |

**RandomForest Classifier**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0.0 | 0.33 | 0.10 | 0.16 | 534 |
| 1.0 | 0.87 | 0.97 | 0.92 | 3440 |
| accuracy |  |  | 0.85 | 3974 |
| macro avg | 0.60 | 0.54 | 0.54 | 3974 |
| weighted avg | 0.80 | 0.85 | 0.82 | 3974 |

**Logistic Regression**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0.0 | 0.25 | 0.00 | 0.01 | 534 |
| 1.0 | 0.87 | 1.00 | 0.93 | 3440 |
| accuracy |  |  | 0.86 | 3974 |
| macro avg | 0.56 | 0.50 | 0.47 | 3974 |
| weighted avg | 0.78 | 0.86 | 0.80 | 3974 |

**GradientBoosting Classifier**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0.0 | 0.00 | 0.00 | 0.00 | 534 |
| 1.0 | 0.87 | 1.00 | 0.93 | 3440 |
| accuracy |  |  | 0.87 | 3974 |
| macro avg | 0.43 | 0.50 | 0.46 | 3974 |
| weighted avg | 0.75 | 0.87 | 0.80 | 3974 |

**Gaussian NB**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0.0 | 0.25 | 0.32 | **0.28** | 534 |
| 1.0 | 0.89 | 0.85 | **0.87** | 3440 |
| accuracy |  |  | 0.78 | 3974 |
| macro avg | 0.57 | 0.59 | 0.58 | 3974 |
| weighted avg | 0.80 | 0.78 | 0.79 | 3974 |

**XGBClassifier**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0.0 | 0.38 | 0.06 | 0.11 | 534 |
| 1.0 | 0.87 | 0.98 | 0.92 | 3440 |
| accuracy |  |  | 0.86 | 3974 |
| macro avg | 0.62 | 0.52 | 0.52 | 3974 |
| weighted avg | 0.80 | 0.86 | 0.81 | 3974 |

# Classification Metrics with application of Data Imbalance handling Techniques-Resampling (Oversampling)

## Linear SVC

|            | precision | recall | f1-score | support |
|------------|-----------|--------|----------|---------|
| 0.0        | 0.56      | 0.68   | 0.62     | 3367    |
| 1.0        | 0.61      | 0.48   | 0.54     | 3450    |
| accuracy   |           |        | 0.58     | 6817    |
| macro avg  | 0.59      | 0.58   | 0.58     | 6817    |
| weighted avg | 0.59    | 0.58   | 0.58     | 6817    |

## RandomForest Classifier

|            | precision | recall | f1-score | support |
|------------|-----------|--------|----------|---------|
| 0.0        | 0.80      | 0.94   | 0.86     | 3367    |
| 1.0        | 0.93      | 0.76   | 0.84     | 3450    |
| accuracy   |           |        | 0.85     | 6817    |
| macro avg  | 0.86      | 0.85   | 0.85     | 6817    |
| weighted avg | 0.87    | 0.85   | 0.85     | 6817    |

## Logistic Regression

|            | precision | recall | f1-score | support |
|------------|-----------|--------|----------|---------|
| 0.0        | 0.63      | 0.57   | 0.60     | 3367    |
| 1.0        | 0.61      | 0.67   | 0.64     | 3450    |
| accuracy   |           |        | 0.62     | 6817    |
| macro avg  | 0.62      | 0.62   | 0.62     | 6817    |
| weighted avg | 0.62    | 0.62   | 0.62     | 6817    |

## GradientBoosting Classifier

|            | precision | recall | f1-score | support |
|------------|-----------|--------|----------|---------|
| 0.0        | 0.62      | 0.66   | 0.64     | 3367    |
| 1.0        | 0.65      | 0.61   | 0.63     | 3450    |
| accuracy   |           |        | 0.63     | 6817    |
| macro avg  | 0.63      | 0.63   | 0.63     | 6817    |
| weighted avg | 0.63    | 0.63   | 0.63     | 6817    |

## Gaussian NB

|            | precision | recall | f1-score | support |
|------------|-----------|--------|----------|---------|
| 0.0        | 0.62      | 0.59   | 0.60     | 3367    |
| 1.0        | 0.61      | 0.64   | 0.63     | 3450    |
| accuracy   |           |        | 0.61     | 6817    |
| macro avg  | 0.61      | 0.61   | 0.61     | 6817    |
| weighted avg | 0.61    | 0.61   | 0.61     | 6817    |

## XGBClassifier

|            | precision | recall | f1-score | support |
|------------|-----------|--------|----------|---------|
| 0.0        | 0.74      | 0.82   | 0.78     | 3367    |
| 1.0        | 0.80      | 0.72   | 0.76     | 3450    |
| accuracy   |           |        | 0.77     | 6817    |
| macro avg  | 0.77      | 0.77   | 0.77     | 6817    |
| weighted avg | 0.77    | 0.77   | 0.77     | 6817    |

# Classification Metrics with application of Data Imbalance handling Techniques-SMOTE

**Linear SVC**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0.0 | 0.51 | 1.00 | 0.68 | 3479 |
| 1.0 | 0.22 | 0.00 | 0.00 | 3338 |
| accuracy |  |  | 0.51 | 6817 |
| macro avg | 0.37 | 0.50 | 0.34 | 6817 |
| weighted avg | 0.37 | 0.51 | 0.35 | 6817 |

**RandomForest Classifier**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0.0 | 0.94 | 0.85 | 0.89 | 3479 |
| 1.0 | 0.85 | 0.94 | 0.90 | 3338 |
| accuracy |  |  | 0.89 | 6817 |
| macro avg | 0.90 | 0.89 | 0.89 | 6817 |
| weighted avg | 0.90 | 0.89 | 0.89 | 6817 |

**Logistic Regression**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0.0 | 0.64 | 0.57 | 0.60 | 3479 |
| 1.0 | 0.59 | 0.66 | 0.63 | 3338 |
| accuracy |  |  | 0.61 | 6817 |
| macro avg | 0.62 | 0.61 | 0.61 | 6817 |
| weighted avg | 0.62 | 0.61 | 0.61 | 6817 |

**GradientBoosting Classifier**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0.0 | 0.85 | 0.79 | 0.82 | 3479 |
| 1.0 | 0.80 | 0.86 | 0.83 | 3338 |
| accuracy |  |  | 0.82 | 6817 |
| macro avg | 0.83 | 0.82 | 0.82 | 6817 |
| weighted avg | 0.83 | 0.82 | 0.82 | 6817 |

**Gaussian NB**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0.0 | 0.61 | 0.74 | 0.67 | 3479 |
| 1.0 | 0.65 | 0.50 | 0.57 | 3338 |
| accuracy |  |  | 0.62 | 6817 |
| macro avg | 0.63 | 0.62 | 0.62 | 6817 |
| weighted avg | 0.63 | 0.62 | 0.62 | 6817 |

**XGBClassifier**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0.0 | 0.98 | 0.84 | 0.90 | 3479 |
| 1.0 | 0.85 | 0.98 | 0.91 | 3338 |
| accuracy |  |  | 0.91 | 6817 |
| macro avg | 0.92 | 0.91 | 0.91 | 6817 |
| weighted avg | 0.92 | 0.91 | 0.91 | 6817 |

# Machine Learning based Loan Default Predictions

# Modeling Binary Classifiers after Data Augmentation

# New Dataset-Classification Metrics without application of Data Imbalance handling Techniques

**Linear SVC**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.82 | 0.86 | 0.84 | 9674 |
| 1 | 0.96 | 0.94 | 0.95 | 32378 |
| accuracy |  |  | 0.93 | 42052 |
| macro avg | 0.89 | 0.90 | 0.90 | 42052 |
| weighted avg | 0.93 | 0.93 | 0.93 | 42052 |

**RandomForest Classifier**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.84 | 0.87 | 0.86 | 9674 |
| 1 | 0.96 | 0.95 | 0.96 | 32378 |
| accuracy |  |  | 0.93 | 42052 |
| macro avg | 0.90 | 0.91 | 0.91 | 42052 |
| weighted avg | 0.93 | 0.93 | 0.93 | 42052 |

**Logistic Regression**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.87 | 0.83 | 0.85 | 9674 |
| 1 | 0.95 | 0.96 | 0.96 | 32378 |
| accuracy |  |  | 0.93 | 42052 |
| macro avg | 0.91 | 0.90 | 0.91 | 42052 |
| weighted avg | 0.93 | 0.93 | 0.93 | 42052 |

**GradientBoosting Classifier**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.84 | 0.88 | 0.86 | 9674 |
| 1 | 0.96 | 0.95 | 0.96 | 32378 |
| accuracy |  |  | 0.94 | 42052 |
| macro avg | 0.90 | 0.92 | 0.91 | 42052 |
| weighted avg | 0.94 | 0.94 | 0.94 | 42052 |

**Gaussian NB**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.58 | 0.74 | 0.65 | 9674 |
| 1 | 0.91 | 0.84 | 0.87 | 32378 |
| accuracy |  |  | 0.81 | 42052 |
| macro avg | 0.74 | 0.79 | 0.76 | 42052 |
| weighted avg | 0.84 | 0.81 | 0.82 | 42052 |

**XGBClassifier**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.86 | 0.87 | 0.86 | 9674 |
| 1 | 0.96 | 0.96 | 0.96 | 32378 |
| accuracy |  |  | 0.94 | 42052 |
| macro avg | 0.91 | 0.92 | 0.91 | 42052 |
| weighted avg | 0.94 | 0.94 | 0.94 | 42052 |

# New Dataset-Classification Metrics with application of Data Imbalance handling Techniques-Resampling (Oversampling)

**Linear SVC**

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.90      | 0.89   | 0.90     | 32315   |
| 1            | 0.89      | 0.90   | 0.90     | 32462   |
| accuracy     |           |        | 0.90     | 64777   |
| macro avg    | 0.90      | 0.90   | 0.90     | 64777   |
| weighted avg | 0.90      | 0.90   | 0.90     | 64777   |

**RandomForest Classifier**

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.93      | 0.98   | 0.95     | 32315   |
| 1            | 0.98      | 0.93   | 0.95     | 32462   |
| accuracy     |           |        | 0.95     | 64777   |
| macro avg    | 0.95      | 0.95   | 0.95     | 64777   |
| weighted avg | 0.95      | 0.95   | 0.95     | 64777   |

**Logistic Regression**

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.93      | 0.92   | 0.92     | 32315   |
| 1            | 0.92      | 0.93   | 0.93     | 32462   |
| accuracy     |           |        | 0.92     | 64777   |
| macro avg    | 0.92      | 0.92   | 0.92     | 64777   |
| weighted avg | 0.92      | 0.92   | 0.92     | 64777   |

**GradientBoosting Classifier**

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.92      | 0.93   | 0.92     | 32315   |
| 1            | 0.93      | 0.92   | 0.92     | 32462   |
| accuracy     |           |        | 0.92     | 64777   |
| macro avg    | 0.92      | 0.92   | 0.92     | 64777   |
| weighted avg | 0.92      | 0.92   | 0.92     | 64777   |

**Gaussian NB**

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.79      | 0.92   | 0.85     | 32315   |
| 1            | 0.91      | 0.75   | 0.82     | 32462   |
| accuracy     |           |        | 0.84     | 64777   |
| macro avg    | 0.85      | 0.84   | 0.84     | 64777   |
| weighted avg | 0.85      | 0.84   | 0.84     | 64777   |

**XGBClassifier**

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.93      | 0.94   | 0.93     | 32315   |
| 1            | 0.94      | 0.93   | 0.93     | 32462   |
| accuracy     |           |        | 0.93     | 64777   |
| macro avg    | 0.93      | 0.93   | 0.93     | 64777   |
| weighted avg | 0.93      | 0.93   | 0.93     | 64777   |

# New Dataset-Classification Metrics with application of Data Imbalance handling Techniques-SMOTE

**Linear SVC**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.88 | 0.96 | 0.92 | 32432 |
| 1 | 0.96 | 0.86 | 0.91 | 32345 |
| accuracy |  |  | 0.91 | 64777 |
| macro avg | 0.92 | 0.91 | 0.91 | 64777 |
| weighted avg | 0.92 | 0.91 | 0.91 | 64777 |

**RandomForest Classifier**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.94 | 0.97 | 0.95 | 32432 |
| 1 | 0.97 | 0.94 | 0.95 | 32345 |
| accuracy |  |  | 0.95 | 64777 |
| macro avg | 0.95 | 0.95 | 0.95 | 64777 |
| weighted avg | 0.95 | 0.95 | 0.95 | 64777 |

**Logistic Regression**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.94 | 0.94 | 0.94 | 32432 |
| 1 | 0.94 | 0.94 | 0.94 | 32345 |
| accuracy |  |  | 0.94 | 64777 |
| macro avg | 0.94 | 0.94 | 0.94 | 64777 |
| weighted avg | 0.94 | 0.94 | 0.94 | 64777 |

**GradientBoosting Classifier**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.94 | 0.95 | 0.94 | 32432 |
| 1 | 0.95 | 0.93 | 0.94 | 32345 |
| accuracy |  |  | 0.94 | 64777 |
| macro avg | 0.94 | 0.94 | 0.94 | 64777 |
| weighted avg | 0.94 | 0.94 | 0.94 | 64777 |

**Gaussian NB**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.87 | 0.91 | 0.89 | 32432 |
| 1 | 0.91 | 0.86 | 0.88 | 32345 |
| accuracy |  |  | 0.89 | 64777 |
| macro avg | 0.89 | 0.89 | 0.89 | 64777 |
| weighted avg | 0.89 | 0.89 | 0.89 | 64777 |

**XGBClassifier**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.96 | 0.96 | 0.96 | 32432 |
| 1 | 0.96 | 0.95 | 0.96 | 32345 |
| accuracy |  |  | 0.96 | 64777 |
| macro avg | 0.96 | 0.96 | 0.96 | 64777 |
| weighted avg | 0.96 | 0.96 | 0.96 | 64777 |

## ANN Classification Metrics  ATT LC  Dataset  Vs. Augmented Dataset

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0.0 | 0.25 | 0.00 | 0.01 | 534 |
| 1.0 | 0.87 | 1.00 | 0.93 | 3440 |
| accuracy | | | 0.86 | 3974 |
| macro avg | 0.56 | 0.50 | 0.47 | 3974 |
| weighted avg | 0.78 | 0.86 | 0.80 | 3974 |

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.86 | 0.85 | 0.85 | 9674 |
| 1 | 0.96 | 0.96 | 0.96 | 32378 |
| accuracy | | | 0.93 | 42052 |
| macro avg | 0.91 | 0.90 | 0.91 | 42052 |
| weighted avg | 0.93 | 0.93 | 0.93 | 42052 |

# Conclusions

Driving Factors (or driver variables)

**(1) Grade:** Default Rate is high in high risk loan applicants. It is important to thoroughly check high risk loan applications.

**(2) Installment Amount:** Defaulter rate increases as the requested loan installment amount increases.

**(4) Annual Income:** Applicants from Low income group have a greater share of defaulted loans.

**(5) States:** The CA, NY, TX states had high number of applications and high default rate. The plot is represented in Figure 16.

(6) **Purpose:** Better to avoid small business loans followed by educational loans

(7) **DTI**: Higher DTI translates to higher default rates

(8) **Income source verification:** Should be checked thoroughly.

(9)  **Binary Classification:** Random Forest and XGBoost provided the best F-Scores and produces good accuracy as well.
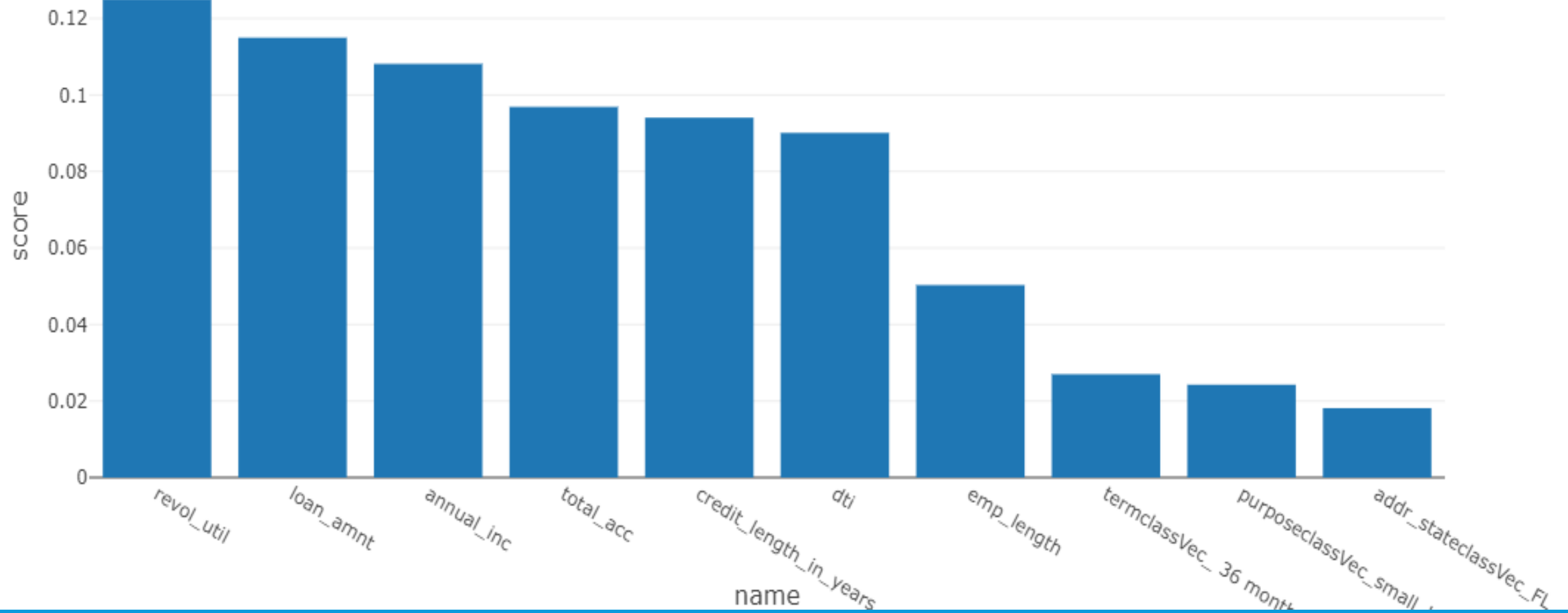
(10) **ANN:** ANN F-score improved on Augmented data.
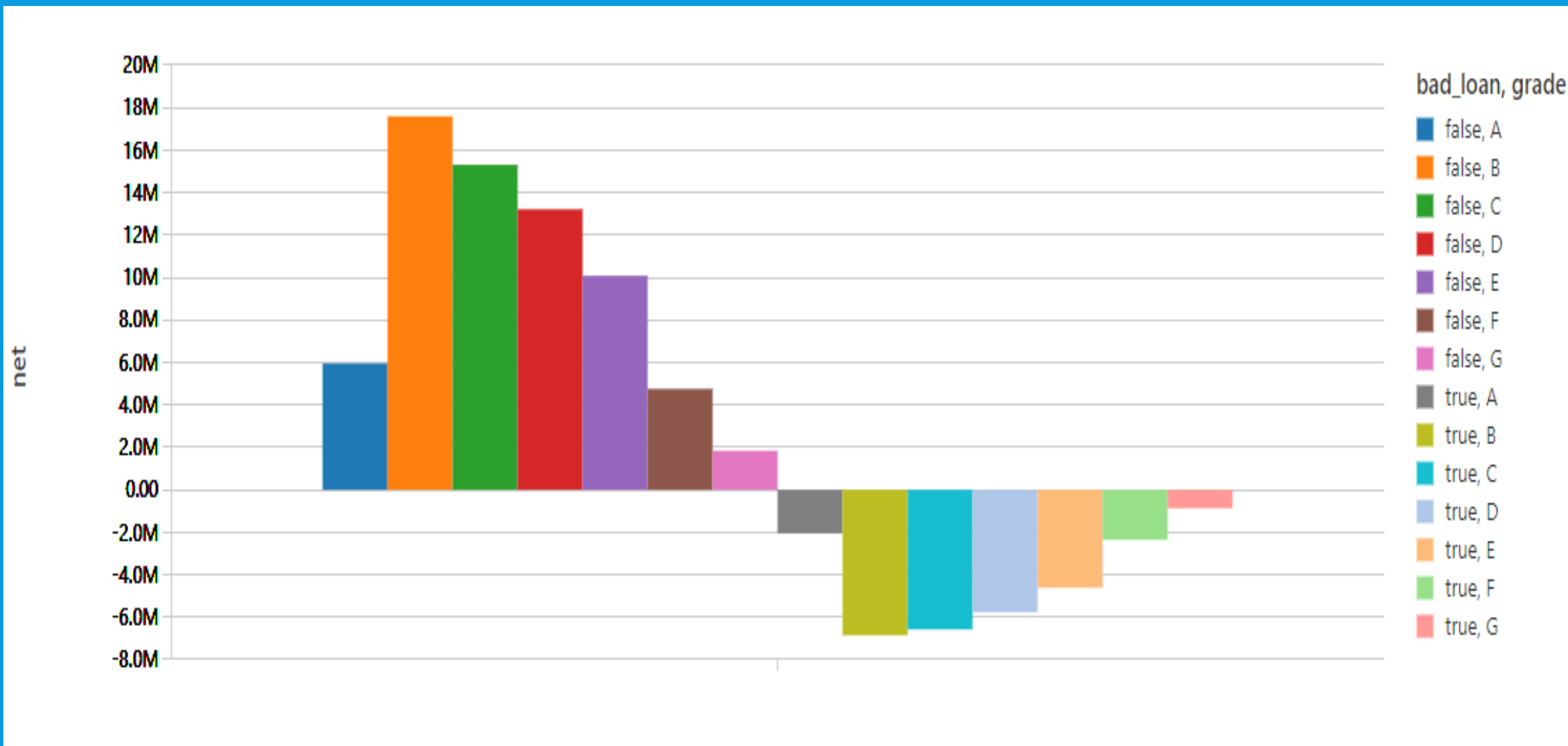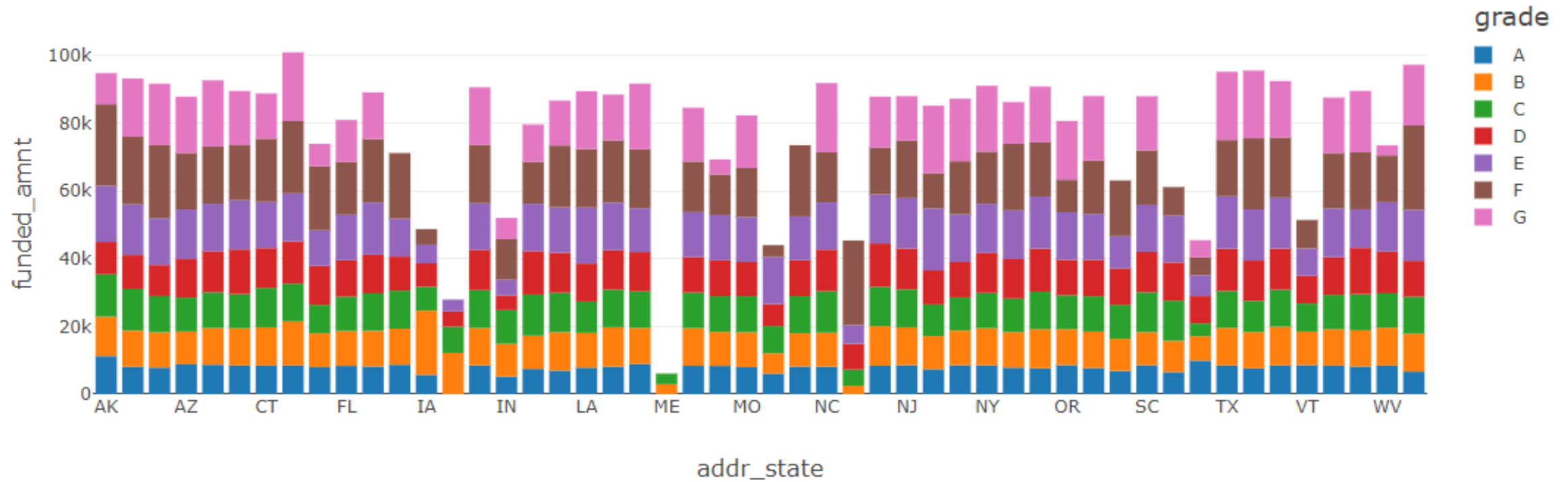
Thank you

# LC Loan Stats

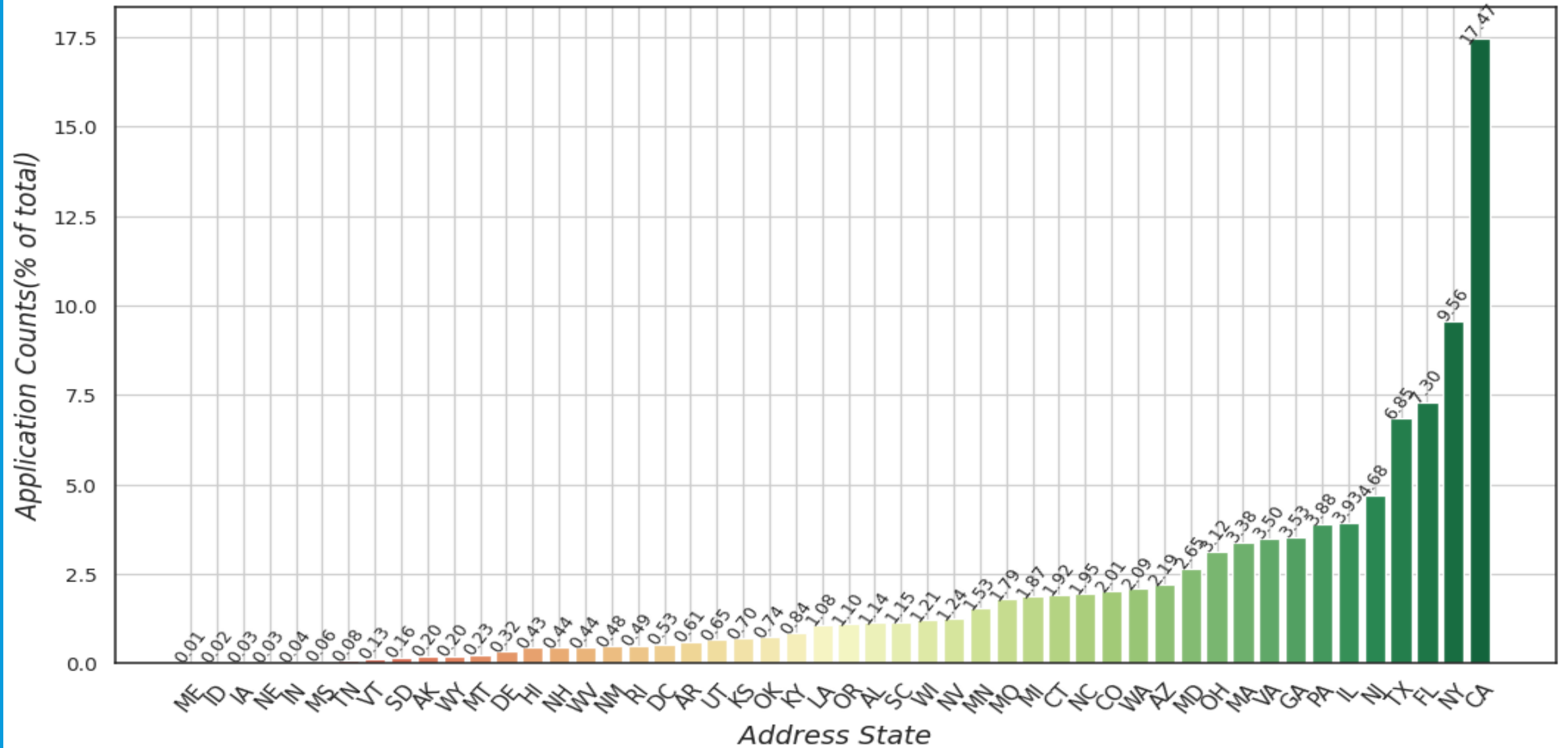# Identify Important Features from Model

Net Profits by Asset Class and Default Status

LC Asset Allocation by Grade

**Address State Analysis(% wise) of Loan Applicants**

**Loan Purpose Analysis**