

MS&E 246 Final Report

Samuel Hansen, Theo Vadpey, Alex Elkrief, Ben Ertringer

2/23/2017

Contents

Exectutive Summary	1
Exploratory Data Analysis	1
Default Rate vs. Business Type	2
Default Rate by Loan Amount	2
Default Rate by NAICS Code	3
Default Rate by Subprogram Type	4
State GDP vs. Default Rate	5
Modeling Default Probability	5
Binary Response Models	5
Cox Proportional Hazards Model	17
Portfolio Selection	23
Modeling Loss at Default	23
Value-at-Risk	23
Average Value-at-Risk	23
Loss Distributions by Tranche	23

Exectutive Summary

In *MS&E 246: Financial Risk Analytics*, our team analyzed a data set of roughly 150,000 loans backed by the US Small Business Administration (SBA) between 1990 and 2014. In doing so, we aimed to implement and test models of the risk and loss of loan default. This report summarizes our findings from exploratory data analysis, details our approaches to modeling loan default probability and loss, and presents our methods of estimating the loss distributions of tranches backed by a portfolio of loans.

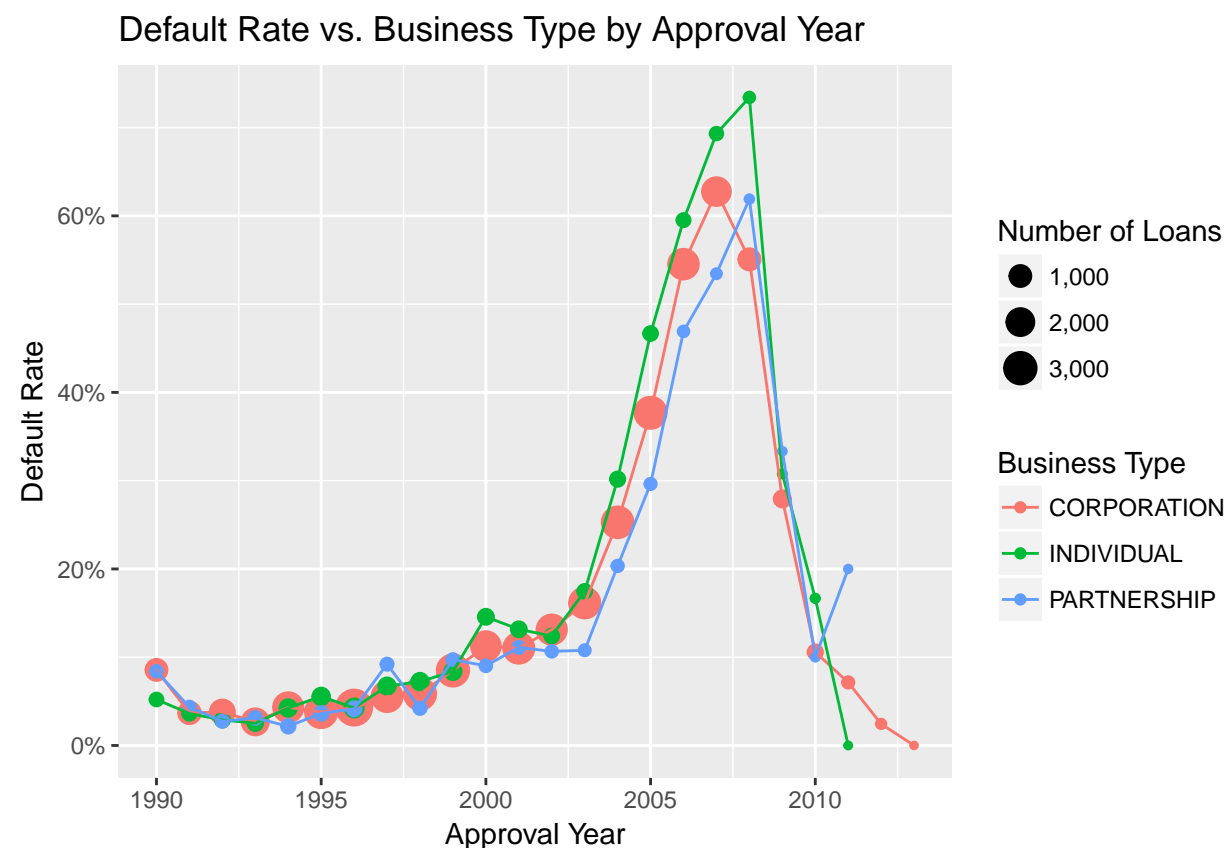
Exploratory Data Analysis

Prior to model building, we explored the data to detect patterns that may provide signal for models of loan default. Because we first aimed to build binary response models of default probability, we excluded “Exempt” loans from our exploratory analysis. Subsequently, we examined the relationship between default rates and the predictor variables, including **Business Type**, **Loan Amount**, **NAICS Code**, and **Loan Term**, among others.

Further, we collected additional predictor variables such as monthly **GDP**, **Crime Rate**, and **Unemployment Rate** by State, as well as macroeconomic predictors such as monthly measures of the **S&P 500**, **Consumer Price Index**, and 14 other volatility market indices (see “Data Cleaning” section for data collection details). We include insights from exploratory analysis of these measures as well.

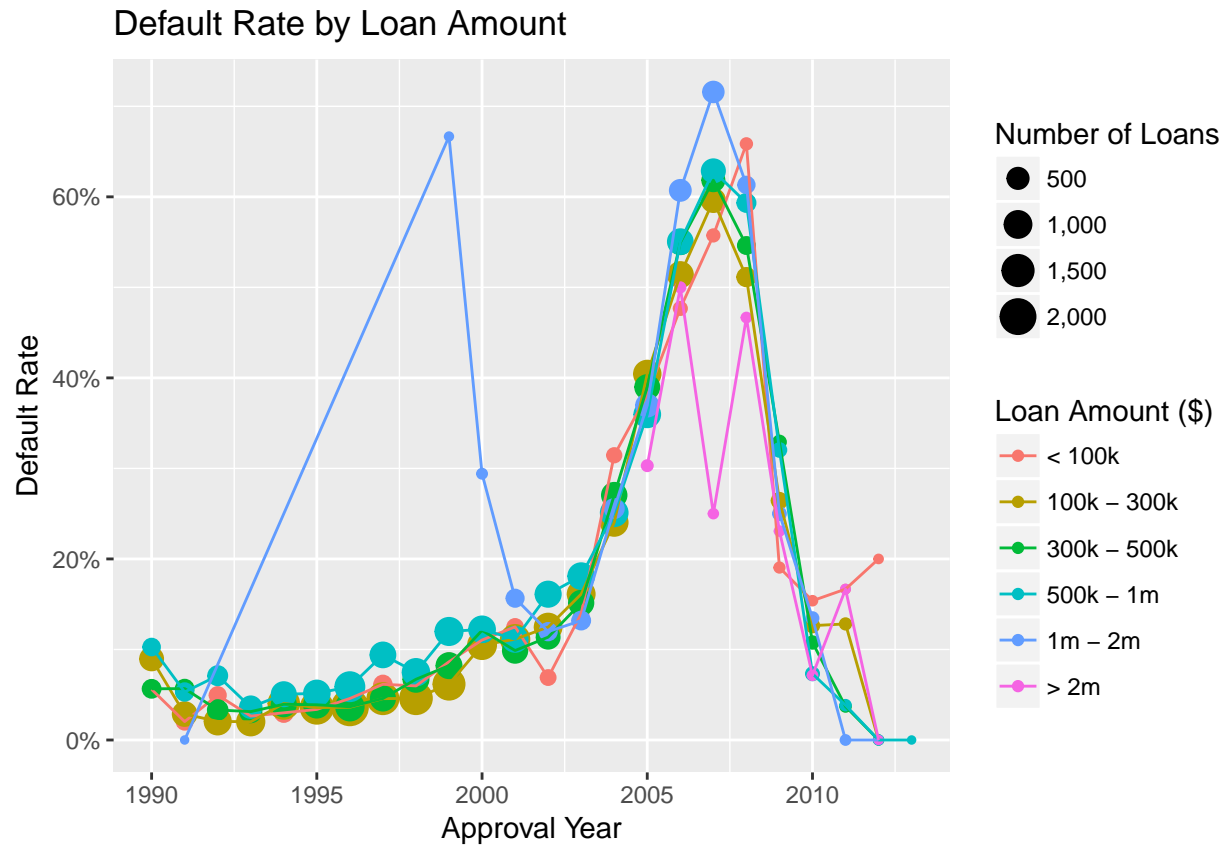
Default Rate vs. Business Type

First, we examined the relationship between default rate and **Business Type** by loan approval year. As shown on the plot below, we observe an interaction effect between these three features, such that default rates spiked for loans that were approved around the Great Recession (approximately 2006- 2009). Further, the different trajectories of the 3 curves implies the “Individual” **Business Type** suffered greater default rates than corporations and partnerships. Although corporations constitute a greater share of the data set, as evidenced by the greater mass in the red circles, they exhibit medium default risk, as compared to the other business types. Taken together, this plot reveals business types were affected differently by the recession, offering useful signal for subsequent modeling.



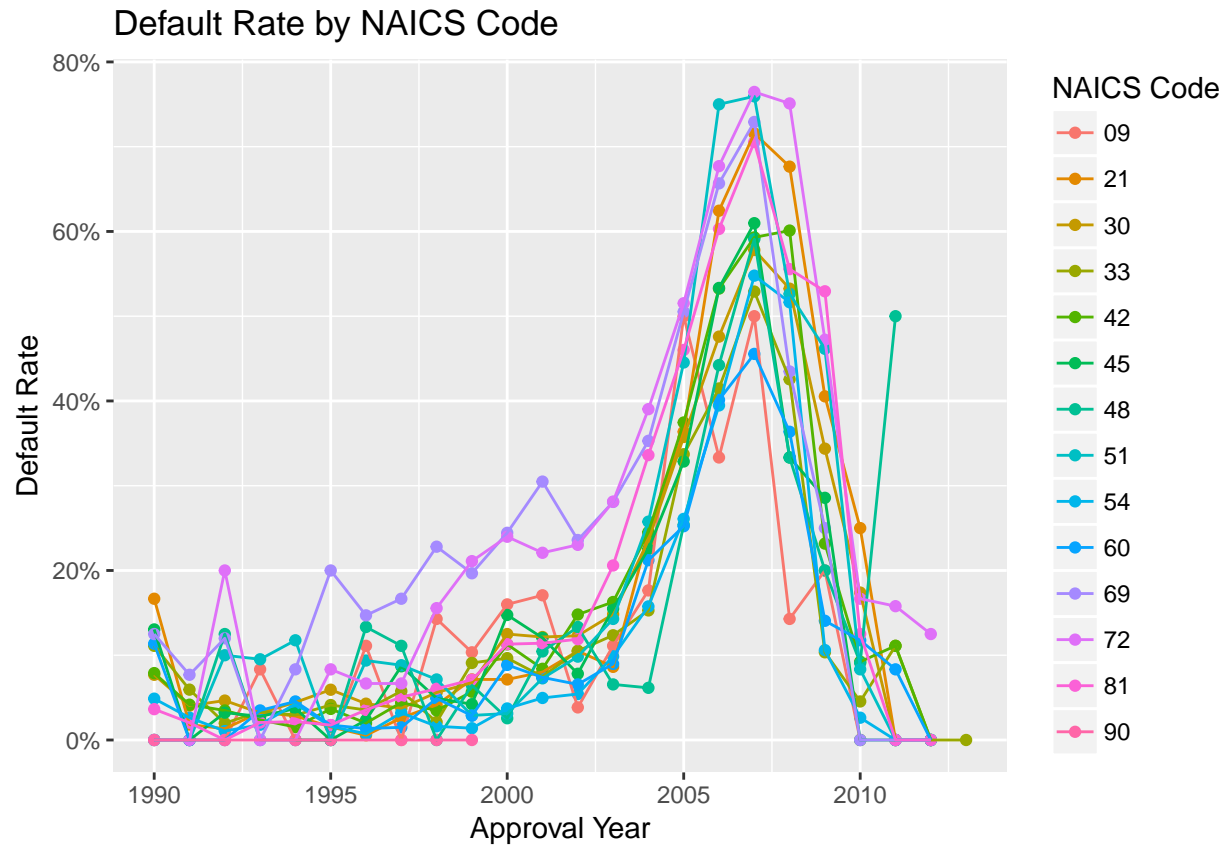
Default Rate by Loan Amount

Second, we examined whether we would observe a similar time-dependent interaction effect between default rate and **Loan Amount**. The plot below reveals that loans of all sizes approved around the Great Recession faced the greatest default rates. However, loans of sizes \$500k-\$1m and \$1m-\$2m appear to have experienced larger default rates over time compared to smaller loans of size \$100k-\$300k and \$300k-\$500k. The spiking behavior of \$1m-\$2m loans in 1999 and of loans greater than \$2m seem to be due to small sample sizes, as depicted by circle diameter. Overall, since loans of different sizes have different default rate patterns over time, we would also expect the **Loan Amount** feature to offer predictive power.



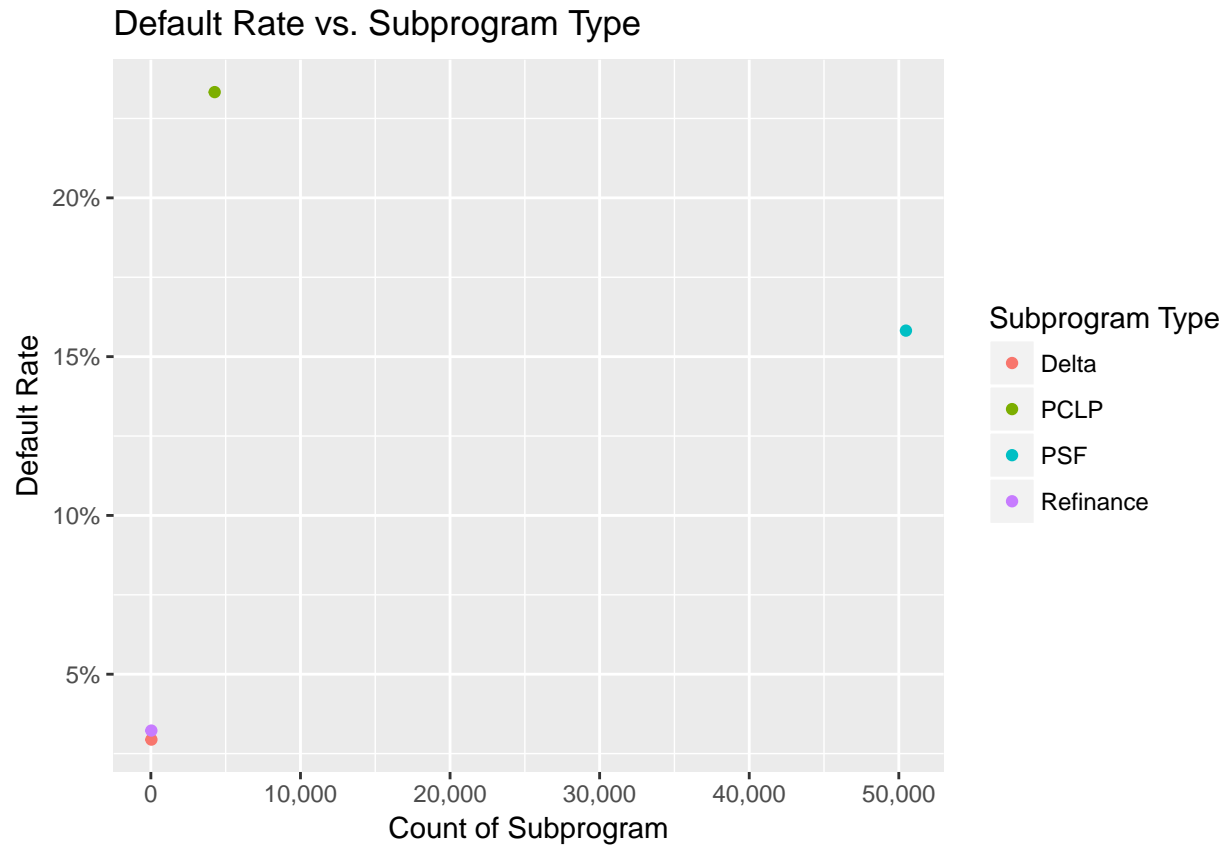
Default Rate by NAICS Code

Third, we hypothesized different economic sectors would exhibit different default rates over time. In turn, we extracted the North American Industry Classification System (NAICS) code for each loan and truncated it to the first two digits, which represents broad industry classes such as “Agriculture” and “Manufacturing.” The following plot shows the default rate for loans of each truncated NAICS code approved in each year between 1990-2014. We observe considerable variance in default rates between sectors; for instance, code 72, corresponding to “Accommodation & Food Services”, has one of the highest default rates even before the recession. However, code 54, corresponding to “Professional, Scientific, and Technical Services,” consistently has one of the lowest default rates. These patterns are consistent with intuition, and underscore the value of including the truncated NAICS code as a predictive feature of defaulting.



Default Rate by Subprogram Type

Fourth, we compared the default rates between different loan subprogram types. The plot below shows the default rates of the different loan subprograms versus their respective counts in the data. We observe that the PSF subprogram is the most common and has medium default risk. However, loans in the Premier Certified Lenders Program (PCLP) are less common, but have higher default risk. This suggests **Subprogram Type** offers useful signal for predicting default risk. Lastly, the loans belonging to the Delta and Refinance subprograms are highly uncommon and have low default risk. In order to reduce to the dimensionality of the feature space, we collapsed these two factor levels into “Other.”



State GDP vs. Default Rate

- [Make plot here](#)

Modeling Default Probability

Building upon our exploratory data analysis, we constructed two types of predictive models of loan default probability: binary response models and the Cox Proportional Hazards model. Here, we present our approach to fitting both model types, including data cleaning, feature engineering, feature selection, hyper-parameter optimization, and evaluation.

Binary Response Models

First, we built binary response models of small-businesses defaulting on loans, which estimate the probability that a given loan *ever* defaults. To do so, we implemented a machine learning pipeline that:

1. Performs feature engineering;
2. Splits the data into train and test sets;
3. Normalizes continuous features;
4. Selects features using recursive feature elimination;
5. Trains binary response predictive models.

Lastly, we evaluated the performance of these models on resampled partitions of the training data, and on a held-out test set in terms of AUC, sensitivity, and calibration.

Feature Engineering

Building on insights derived from exploratory data analysis, we engineered the following features from the raw data:

- **NAICS_code**: truncated to the first two digits of the NAICS code;
- **subprogram**: condensed infrequent factor levels into “other” category;
- **approval_year**: extracted year from loan approval date-time object.
- **SameLendingState**: created flag for whether borrower received loan from in-state;
- **MultiTimeBorrower**: created flag for whether loan recipient is a multi-time borrower;
- **ThirdPartyLender** created flag for whether borrower received third party aide.

In effect, these features represent dimensionality reduction of factors with many levels. For instance, there are 1,239 unique NAICS six-digit NAICS codes in the raw data, yet only 25 unique 2-digit codes. Although we lose fine-grained detail by truncating the NAICS code, we aimed to optimize our models by reducing variance introduced by high dimensionality. After applying such dimension reductions, we eliminated extraneous variables, such as the Borrower’s Zip Code and the Project’s State, which were used to engineer features.

In addition to constructing features from the raw data, we also incorporated data from external sources, including monthly State-based measures of crime rate, GDP, and unemployment rate. We also joined in time-varying risk factors, including monthly snapshots of the **S&P 500**, **Consumer Price Index**, and 14 other volatility market indices.

- **BEN**: Fill in where the data came from and any other important info

Data Splitting

We randomly partitioned the data into 70% training and 30% test sets. This approach does not implement a time-based split, but rather a random sampling of observations over the entire 1990-2014 window. We adopted this splitting approach because we were interested in capturing the signal of the Great Recession within our models. Further, we did not create a validation set because we performed feature selection and hyper-parameter optimization using cross-validation on the training set.

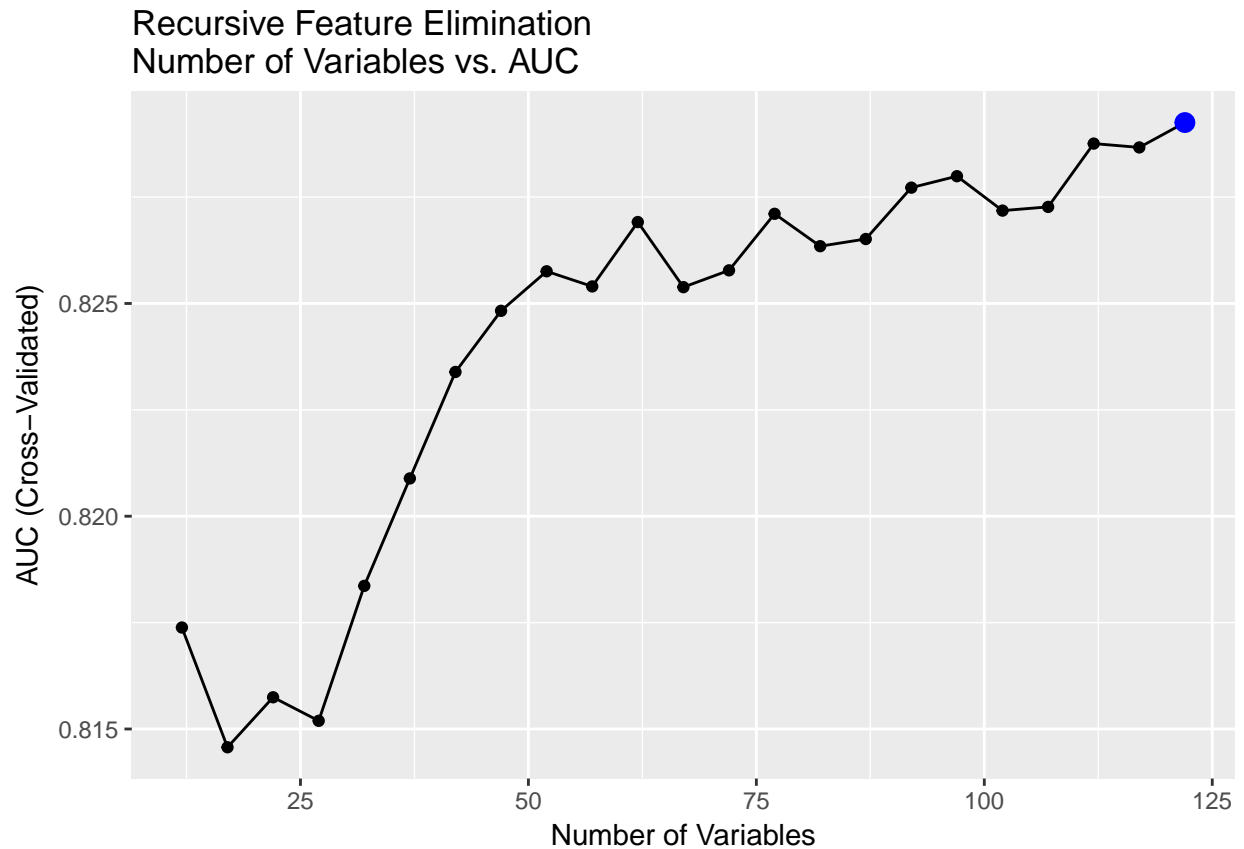
Data Preprocessing

After engineering features and joining in external data sources, we applied several preprocessing steps to our main data frame. First, we centered and scaled the continuous predictors to apply regularization techniques during the modeling phase. Doing so adjusted for variables being on different scales; for example, **Gross Approval** varies in dollar amounts from \$30,000 to \$4,000,000, whereas **Term in Months** ranges from 1 to 389. Second, we applied a filter to remove features with near zero variance to eliminate predictors that do not offer meaningful signal.

Feature Selection

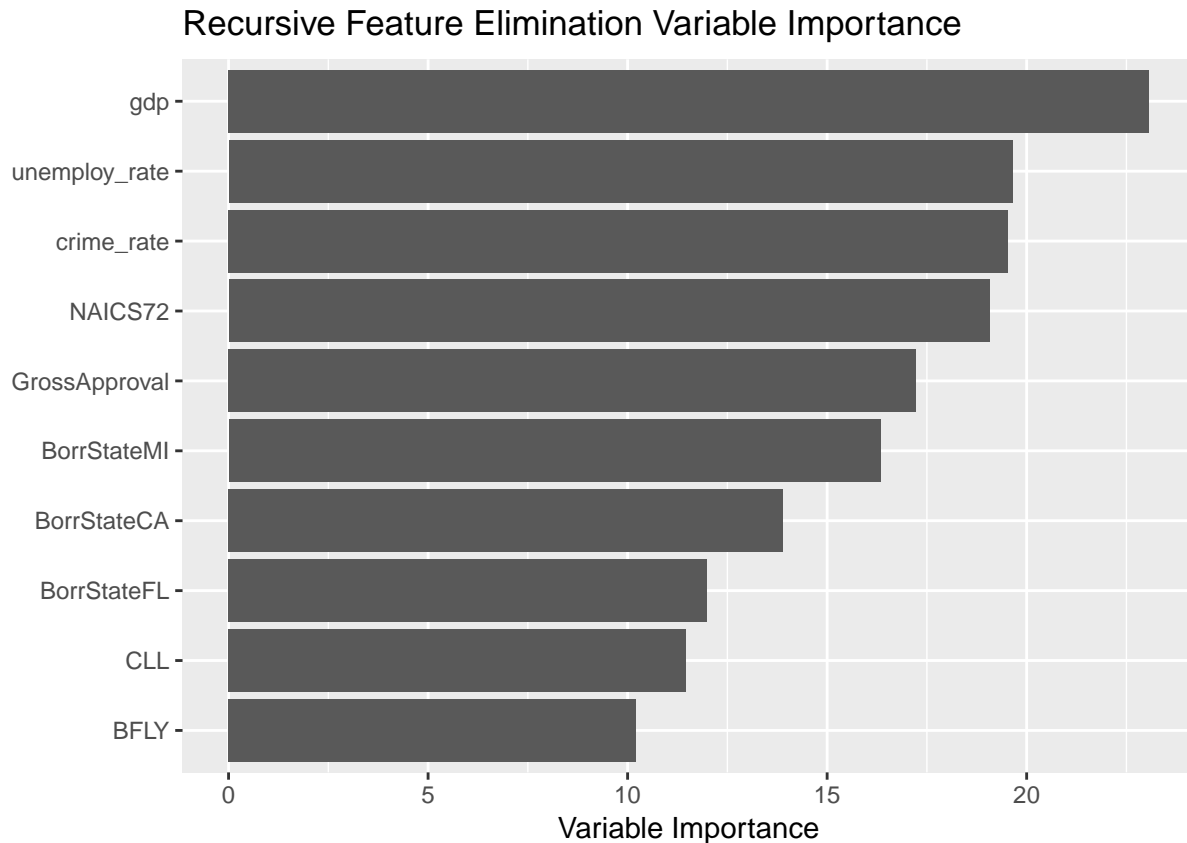
To perform feature selection, we used recursive feature elimination with 10-fold cross-validation. This method uses random forests to iteratively remove variables with low variable importance, as measured by mean increase in out-of-bag area-under-the-curve (AUC). In other words, variables that do not contribute to significant improvements in AUC are eliminated. We performed a grid search over the number of potential features to determine how many features to include. Note that factors were converted to separate dummy variables using a one-hot encoder.

The following plot shows that recursive feature selection chose 122 variables because AUC is maximized (see plot below). In effect, all variables were kept because they offered predictive power regarding loan defaults.



The importances of the top 10 selected features are shown in the plot below. We observe that State GDP, a monthly time-dependent risk factor, is the most important feature, meaning it led to the greatest average increase in AUC across cross-validation iterations. State unemployment rate and crime rate are also highly important, suggesting local time-dependent risk factors are the most predictive of whether a loan defaults.

The importance of NAICS code 72, corresponding to “Accommodation & Food Services”, is consistent with our exploratory data analysis finding that the sector is especially risk prone. Borrower States such as Michigan, California, and Florida also offer predictive power regarding defaulting. Lastly, the importances of the Collar Index (CLL) and Iron Butterfly Index (BFLY) imply market volatility measures also improve the discrimination of loan defaults.

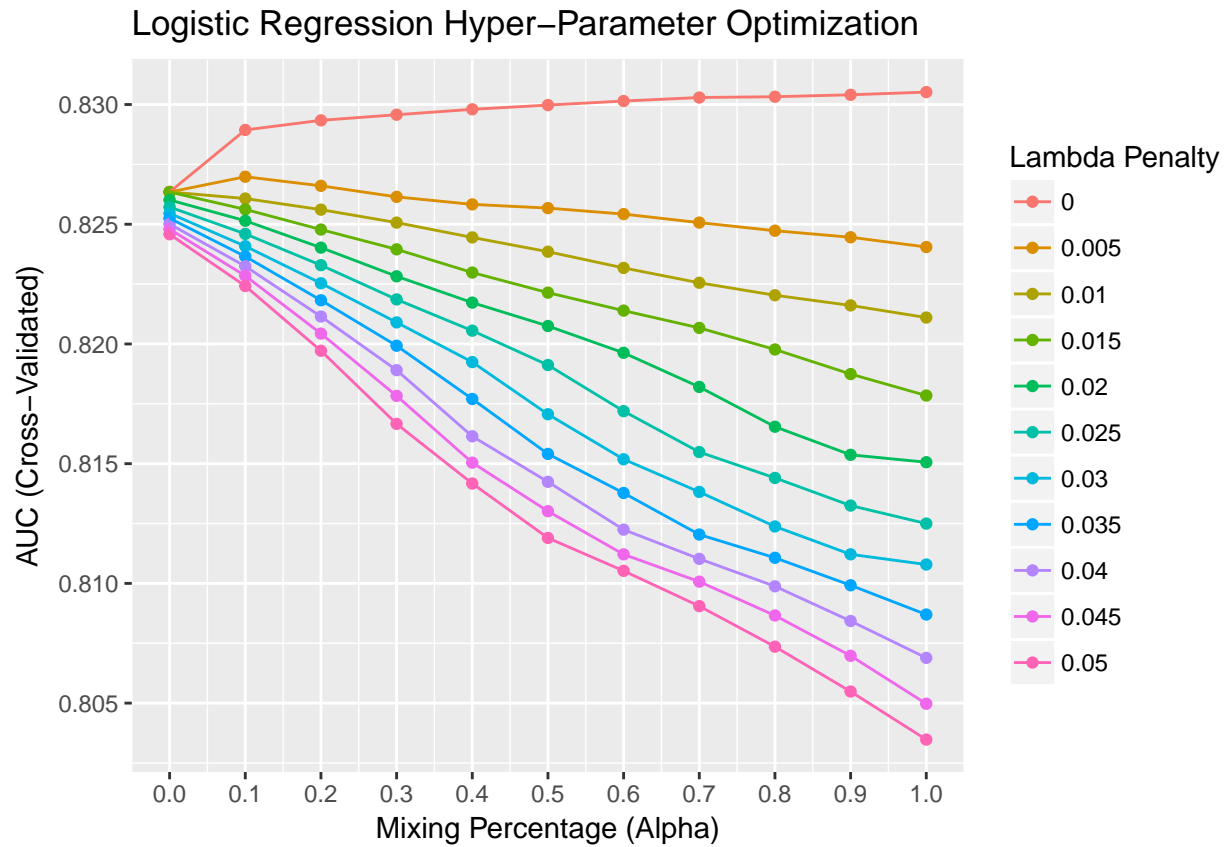


Model Fitting

Using these selected features, we fit models predicting the binary outcome of whether a small business defaults on a loan. We constructed linear and nonlinear models, including a logistic regression model with the elastic net penalty, a random forest classifier, and a gradient boosting machine classifier. To tune hyper-parameters, we used 10-fold cross-validation with the one standard-error rule, which selects parameters that obtain the highest cross-validated AUC within one standard error of the maximum. For each model type, we performed a grid search over the hyper-parameters to ensure optimal selection.

Logistic Regression with Elastic Net

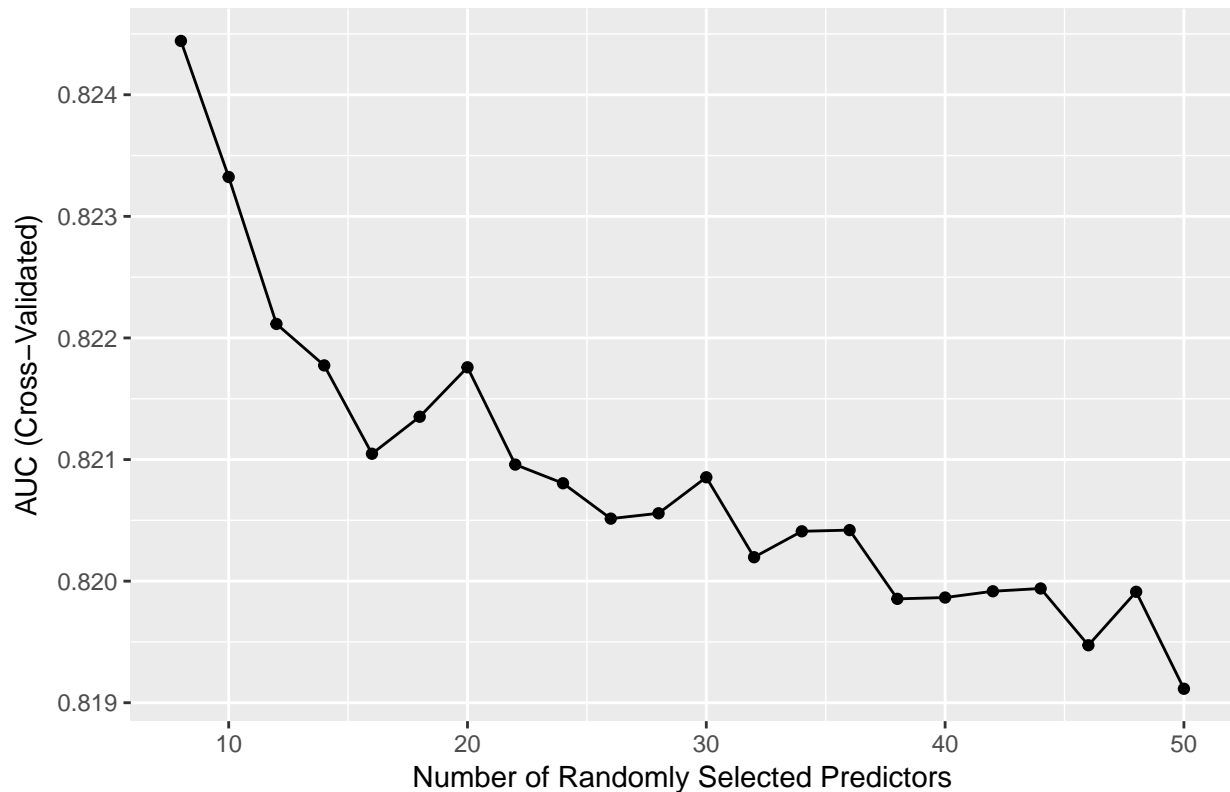
First, AUC was used to select the optimal logistic regression model with an elastic net penalty using the one standard-error rule. As shown in the plot below, the final values used for the model were $\alpha = 0.1$ and $\lambda = 0$, indicated by the spike in the red curve at $\alpha = 0.1$. This implies the optimal model used the ridge penalty more than the LASSO penalty with minimal regularization.



Random Forest Classifier

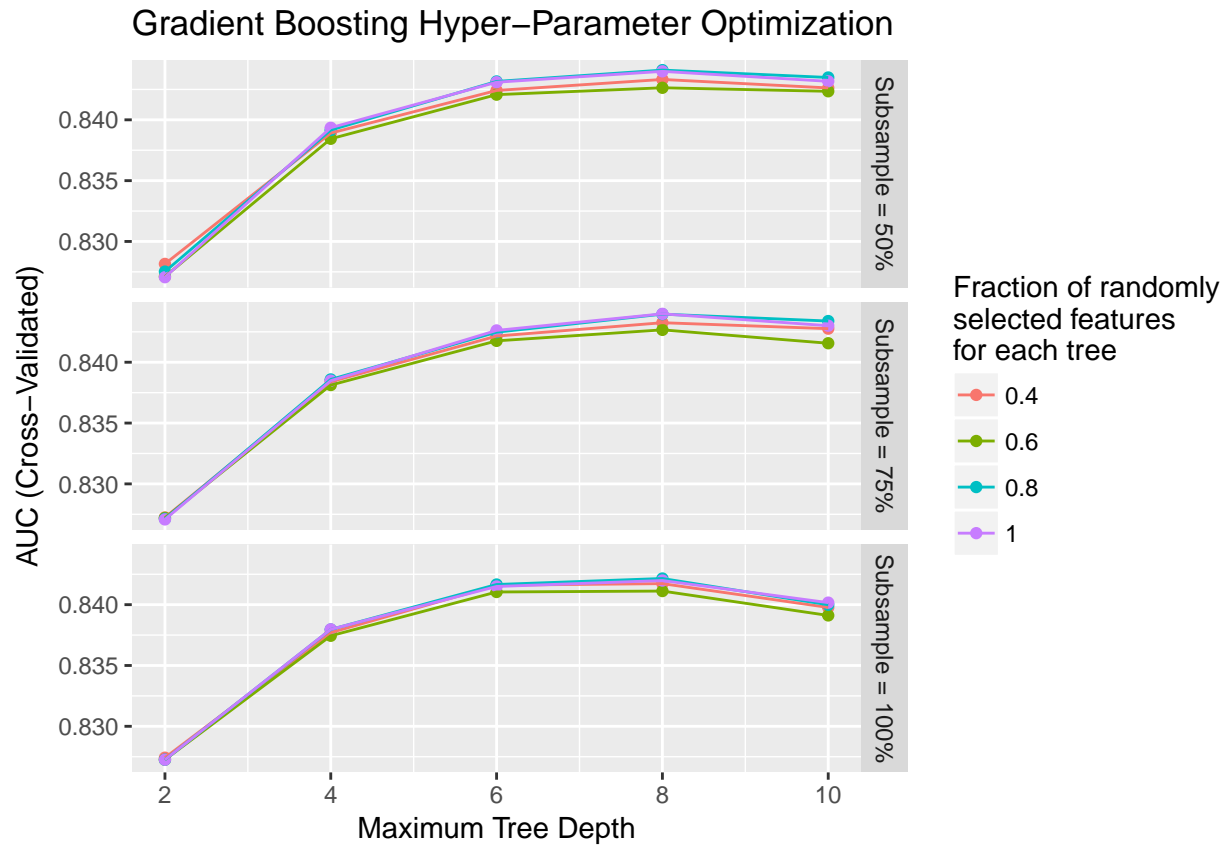
Second, AUC was used to select the optimal random forest model, which selected `mtry = 8` as the best parameter. This means 8 random predictors were chosen to build each tree of the random forest. The plot below shows steadily declining AUC as the number of randomly chosen predictors increases, indicating that the optimal model is sparsest.

Random Forest Hyper-Parameter Optimization

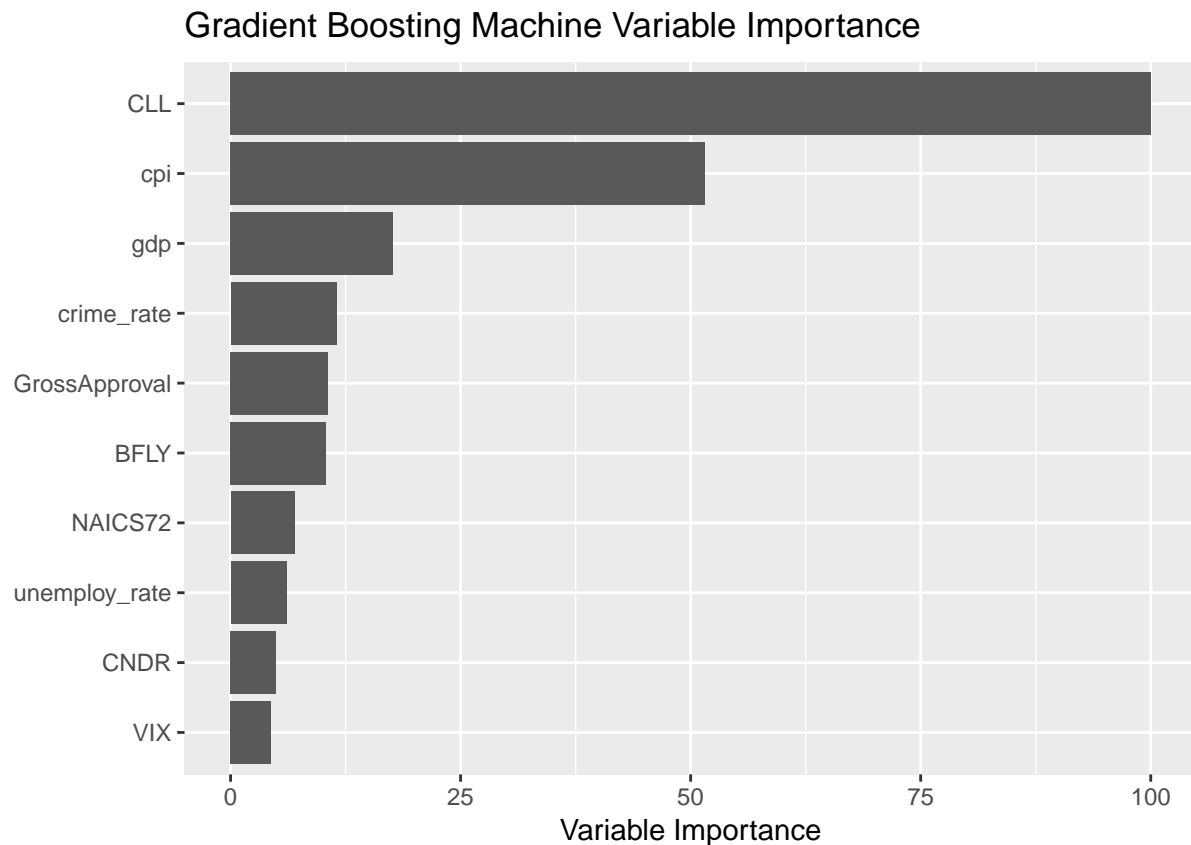


Gradient Boosting Machine Classifier

Third, AUC was similarly used to select the optimal gradient boosting machine (GBM) model. The final values used for the model were `nrounds = 100`, `max_depth = 6`, `eta = 0.03`, `gamma = 0`, `colsample_bytree = 0.4`, `min_child_weight = 1` and `subsample = 0.5`. This means that the tuning procedure utilized a learning rate of 0.03 and a minimum loss reduction of 0, resulting in the optimal model with 100 trees of maximum depth 6 that subsamples 50% of the observations and 40% of the features for each tree. This combination of optimal hyper-parameters is shown by the spike of the red curve in the first subplot at the maximum tree depth of 6.



Examining the variable importance of the final GBM model, we observe the most important feature for predicting defaults is the Collar Index (CLL), which is “designed to provide investors with insights as to how one might protect an investment in S&P 500 stocks against steep market declines” (CBOE). Other important features include the national consumer price index (CPI), State GDP, crime, and unemployment rates, loan amount, and Chicago Board Options Exchange (CBOE) indices including the Butterfly Index (BFLY), the Iron Condor Index (CNDR), and the Volatility index (VIX). Such variables are “important” because they lead to the greatest improvements to cross-validated AUC across boosting iterations.



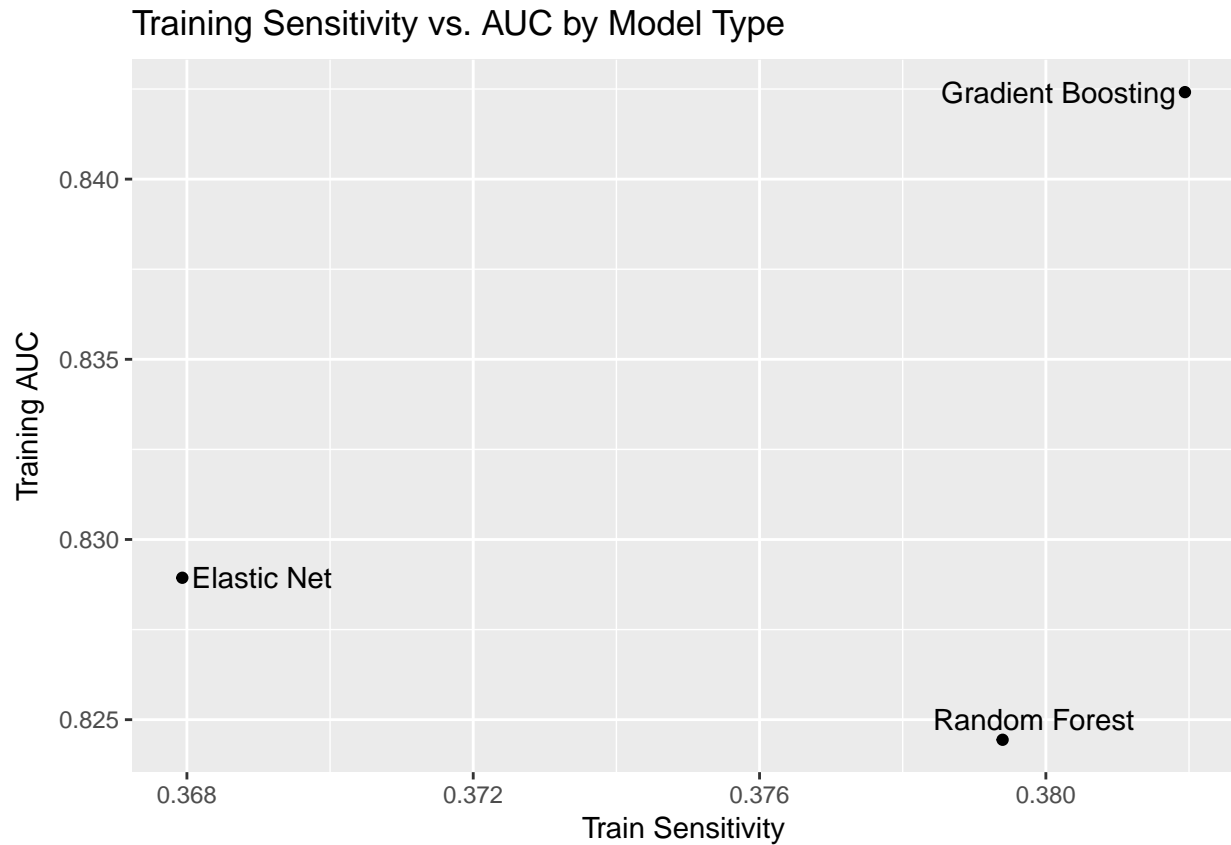
Model Evaluation

After we optimized the hyper-parameters of our models, we evaluated the models using in-sample and out-of-sample metrics, including AUC, sensitivity, ROC curves, and calibration. To do so, we used these “best” models to predict loan defaults in the training and test sets.

In-Sample Evaluation

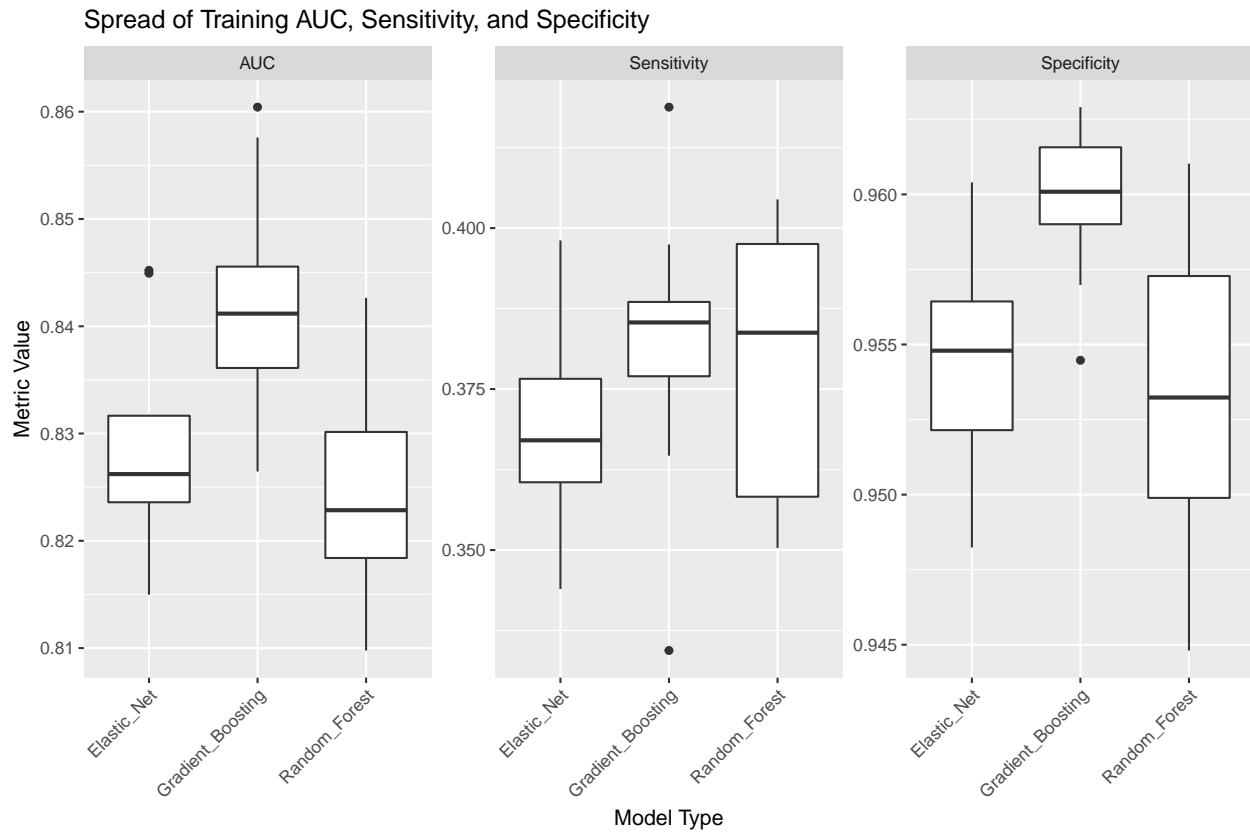
Training AUC and Sensitivity of Best Models

The following plot compares averaged **training** area under the ROC curve and sensitivity across the model types with optimized parameters. We observe that the gradient boosting machine classifier has the highest AUC and sensitivity, whereas the logistic regression model with the elastic net penalty performs the worst.



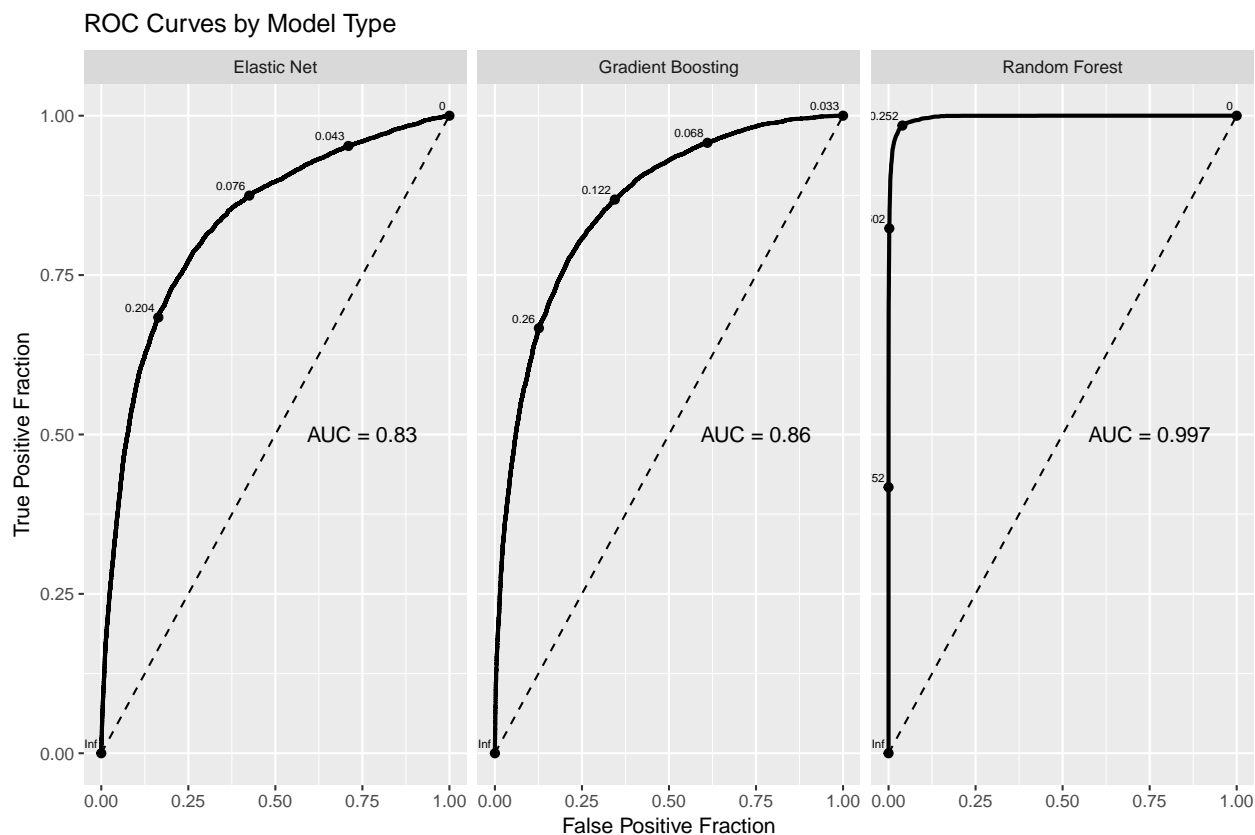
Distribution of Resampled Training AUC, Sensitivity, and Specificity

To examine the spread of **training** area under the ROC curve, sensitivity, and specificity across model types, we leverage the resampled data generated during the cross-validation of model fitting to plot their respective distributions. In the following plot, we observe that the GBM classifier has the highest median AUC, sensitivity, and specificity, as well as the smallest spread. Although the random forest classifier has comparable sensitivity, it exhibits enormous variance compared to the other models, suggesting it is prone to overfitting. For this reason, the logistic regression classifier (a linear model) outperforms the random forest classifier (a non-linear model) in terms of AUC and specificity.



Training ROC Curves

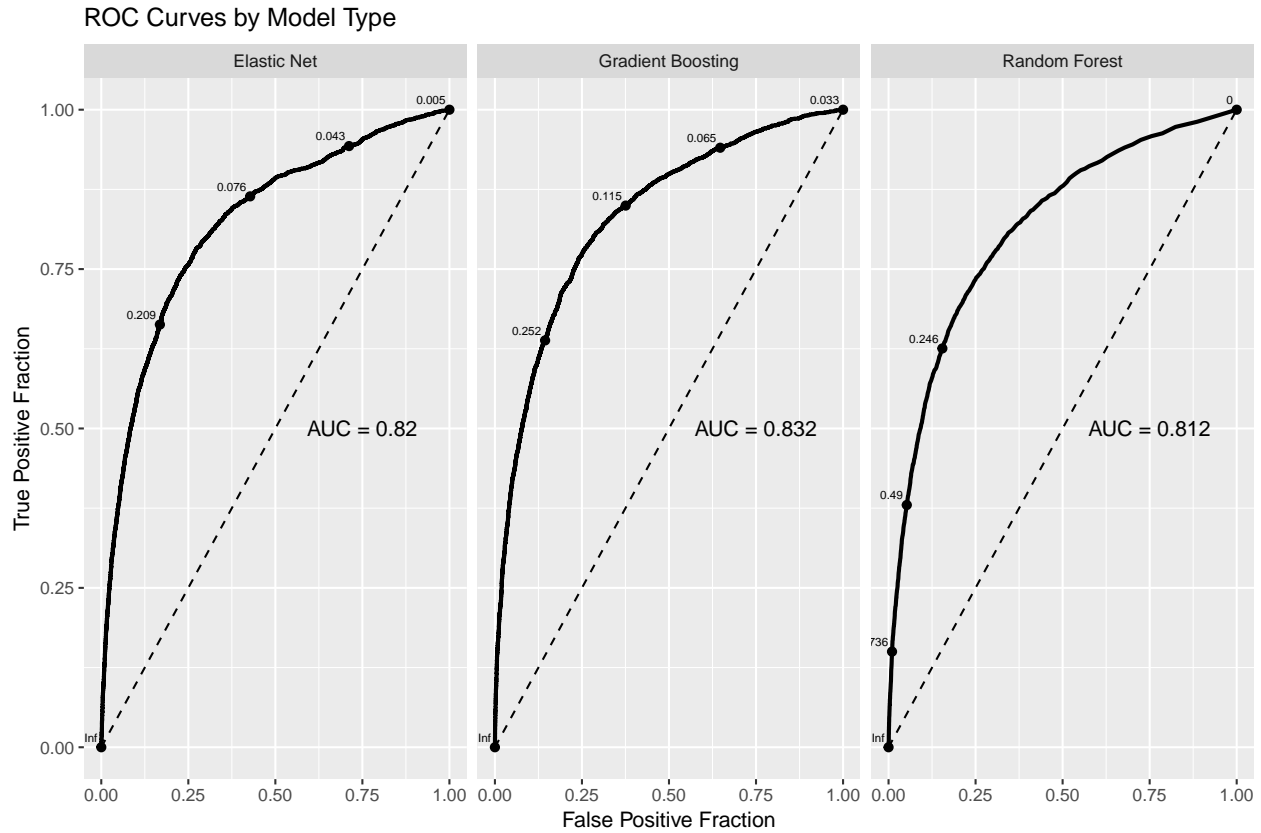
Lastly, we can examine the **training** ROC curves by model type. We observe that the random forest model has a near-perfect ROC curve, which also implies it is overfitting to the training data. The GBM model again performs worse than the random forest model on the training data, but likely because it is avoiding overfitting. The logistic regression model with the elastic net penalty performs the worst.



Out-of-Sample Evaluation

Test ROC Curves

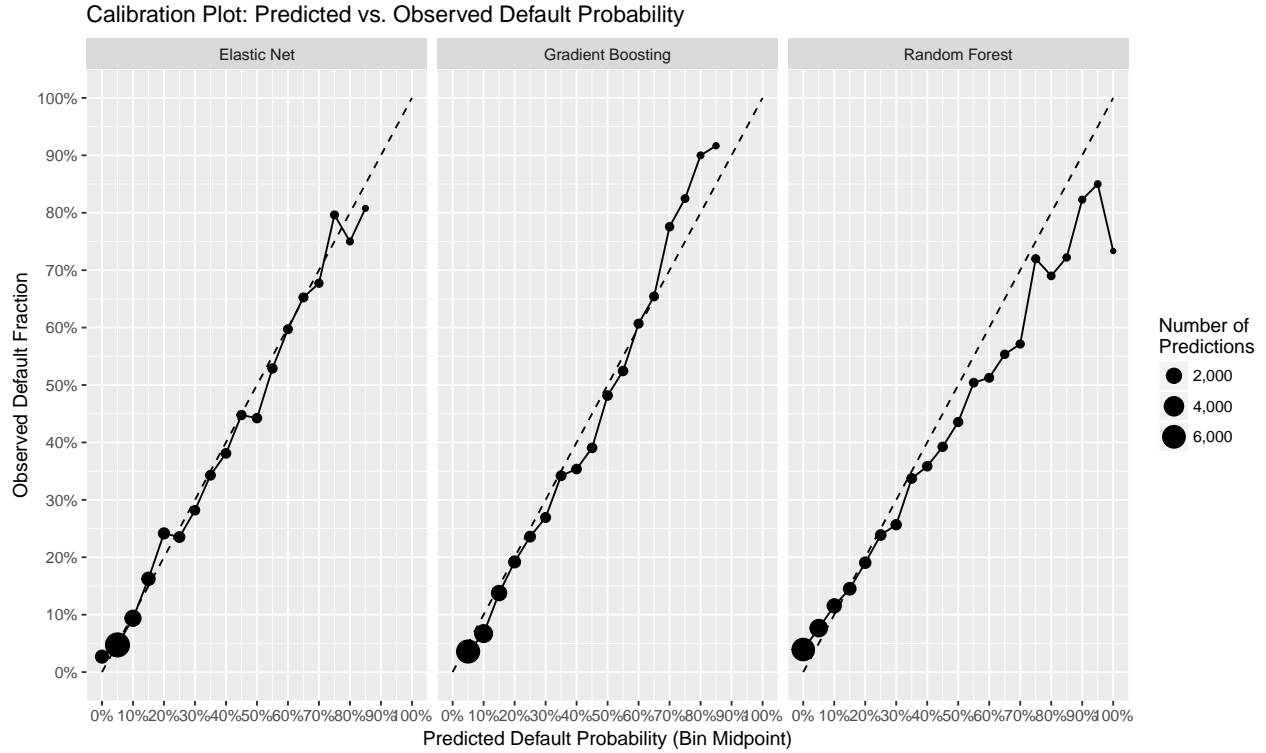
We evaluated our best models on a held-out test set representing 30% of the original data. The ROC curves below reveal that the GBM model performed the best on the test set, followed by the logistic regression model, and finally, the random forest classifier. The weak performance of the random forest classifier is likely due to overfitting on the training set. Nevertheless, all models achieve good performance over “random guessing” baselines.



Test Calibration Plots

Lastly, we evaluated the calibration of the our models' predicted probabilities of loan default against the observed fraction of defaults in the data. A point on the dashed line means that the model's predicted probability of default matched the empirical default rate. Points to the right of the line mean the model overestimated the default probability, whereas points to the left mean the model underestimated the default probability.

The GBM model achieves the best calibration because its points follow the dashed line most closely. The logistic regression model with the elastic net penalty achieves comparable performance; however, the random forest classifier tends to overestimate default probabilities. Again, this weaker performance is likely due to overfitting.



The overfitting of the random forest classifier may be due to the fact that too many features were randomly selected to build trees at each iteration. Our hyper-parameter optimization approach performed a grid search over possible values of `mtry`, representing the number of features randomly chosen to build each tree in the forest. However, our grid may have not been large enough, since the minimum value of `mtry` was chosen. However, computational resources limited our ability to refit models over a larger search space.

Moreover, the gradient boosting machine classifier demonstrated the best performance on the test set in terms of AUC and calibration.

Cox Proportional Hazards Model

Motivation:

Survival analysis gives more detailed information about how the default risk of a loan varies over time. With binary classification, we estimated the probability that a given loan *ever* defaults. With a hazard model, we are able to estimate the probability that a loan defaults between any two points of time in its life.

Model Choice:

There exist many specialized Cox models that assume a particular form of the baseline hazard function. The Cox Proportional Hazards Model does not have this requirement. We can see this in the following description of the partial maximum likelihood procedure used to estimate the parameters of the Cox PH model:

The form of the cox model is:

$$h(t) = h_0(t) \exp(\beta^T X)$$

Suppose there are r observed death times in the data (all distinct), and that t_j is a death time in the set of possible death times: $R = \{t_1, t_2, \dots, t_r\}$.

Then the conditional probability that an individual dies at time t_j given t_j is a time of death in the set R :

$$\begin{aligned} & \frac{P(\text{individual with feature vector } X^{(j)} \text{ dies at } t_j)}{P(\text{one death at } t_j)} \\ &= \frac{P(T = t_j | X^{(j)}, T \geq t_j)}{P(T = t_j | X^{(k_0)}, T \geq t_j) \cup P(T = t_j | X^{(k_1)}, T \geq t_j) \cup \dots P(T = t_j | X^{(k_q)}, T \geq t_j)} \end{aligned}$$

Where k_0, \dots, k_q correspond to the indices of observations with event times greater than or equal to t_j . Since the probabilities in the denominator are *assumed to be conditionally independent*, the denominator can be expressed as a sum of probabilities. Converting the above to continuous time, we get:

$$\begin{aligned} &= \frac{\lim_{\delta \rightarrow 0} \frac{P(T < t_j + \delta | X^{(j)}, T \geq t_j)}{\delta}}{\sum_{i=k_0}^{k_q} \lim_{\delta \rightarrow 0} \frac{P(T < t_j + \delta | X^{(i)}, T \geq t_j)}{\delta}} \\ &= \frac{h_j(t_j)}{\sum_{i=k_0}^{k_q} h_i(t_j)} = \frac{h_0(t_j) \exp(\beta^T X^{(j)})}{\sum_{i=k_0}^{k_q} h_0(t_j) \exp(\beta^T X^{(i)})} = \frac{\exp(\beta^T X^{(j)})}{\sum_{i=k_0}^{k_q} \exp(\beta^T X^{(i)})} \end{aligned}$$

And we can see that the contribution of any observation to the likelihood function will not be dependent on $h_0(t)$. \square

Additional Modifications to the Data

Roughly 95% of loans in the training data set had a term of 20 years. We decided that considering loans with the same term was more appropriate for this analysis (84,949 loans).

Within the training data, about 86% of loans were right censored (term did not expire in window, and did not default), about 7% of loans were paid off (term expired in window), and about 7% of loans defaulted within the window (figure 1).

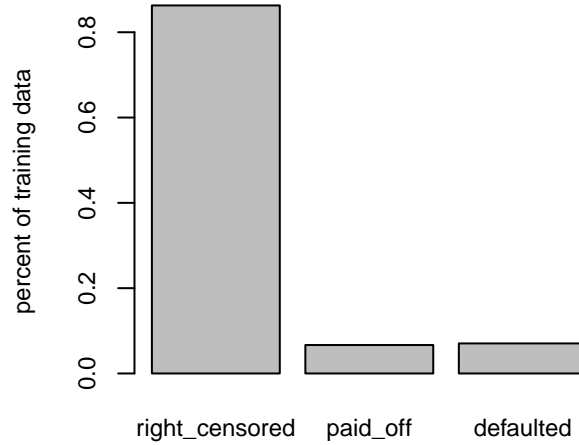


Figure 1: Loans in training data by status

Polynomial terms up to *degree five* were added for all numeric variables. Our intention was to include these features to capture non-linearities in these variables, and conduct feature selection during model fitting (through regularization).

All numeric variables were centered to 0, and scaled by standard deviation.

Missing values were set to 0 and an missing value indicator feature was added for each original variable. Including expanded categorical variables, polynomials, and missing value dummies, the data had 201 features.

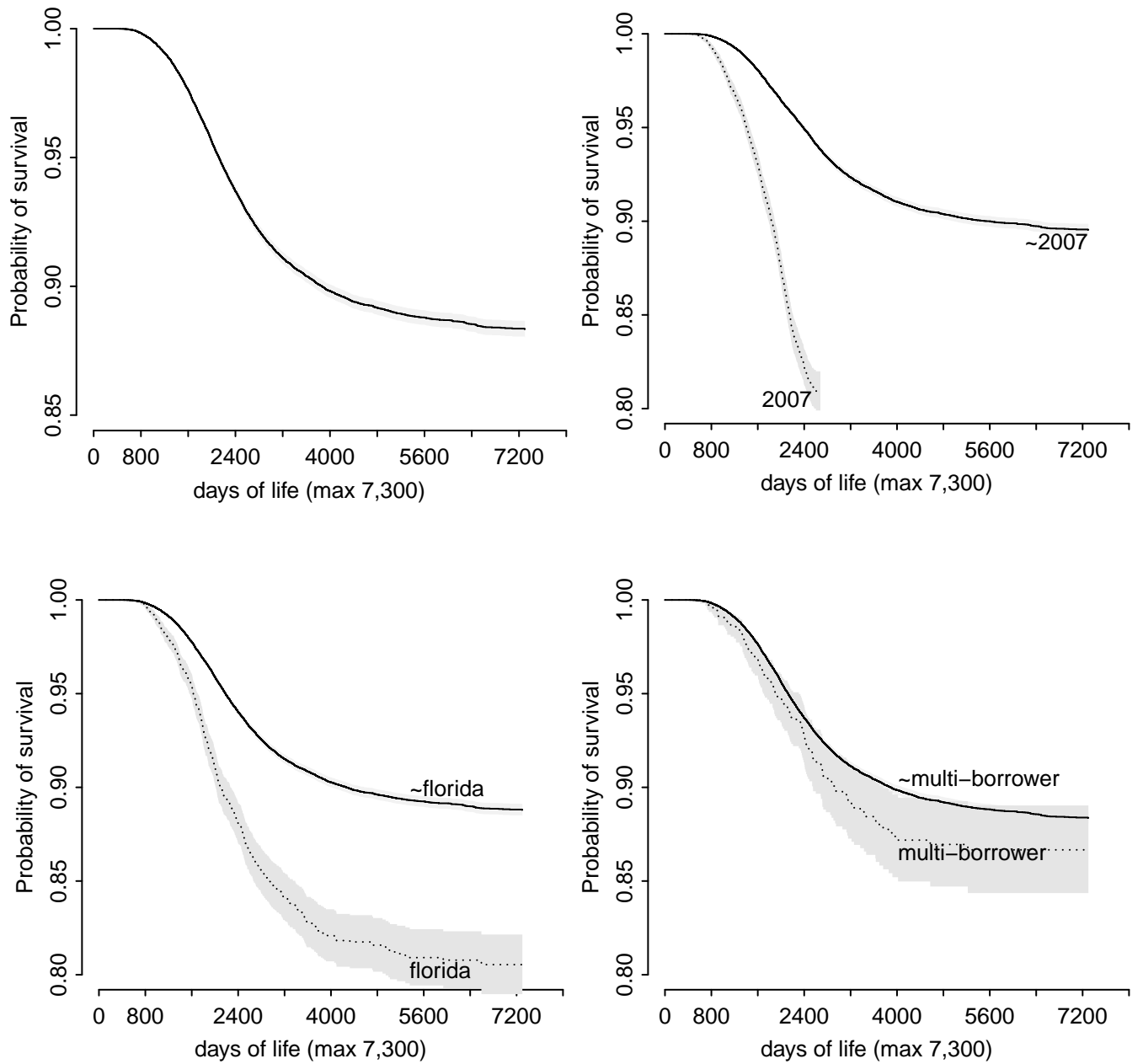
Kaplan-Meier Survival Curve

A Kaplan-Meier curve is a non-parametric estimate of the survival function, $S(t) = P(T > t)$, defined as:

$$\hat{S}(t) = \prod_{t_i \leq t} \left[1 - \frac{d_i}{n_i}\right]$$

Where $\{t_1, \dots, t_r\}$ are the death times of observations in the data, $\{d_1, \dots, d_r\}$ are the number of deaths that occur at those times, and $\{n_1, \dots, n_r\}$ are the number of observations remaining in the at-risk population just before those times.

For expository purposes the following plots show the estimated survival function conditioned on select categorical variables such as a particular year, state, or status, as well as the general survival curve for our loan population. Note that the survival curve was significantly steeper for loans conditioned on these variables (a higher probability of default at all times).



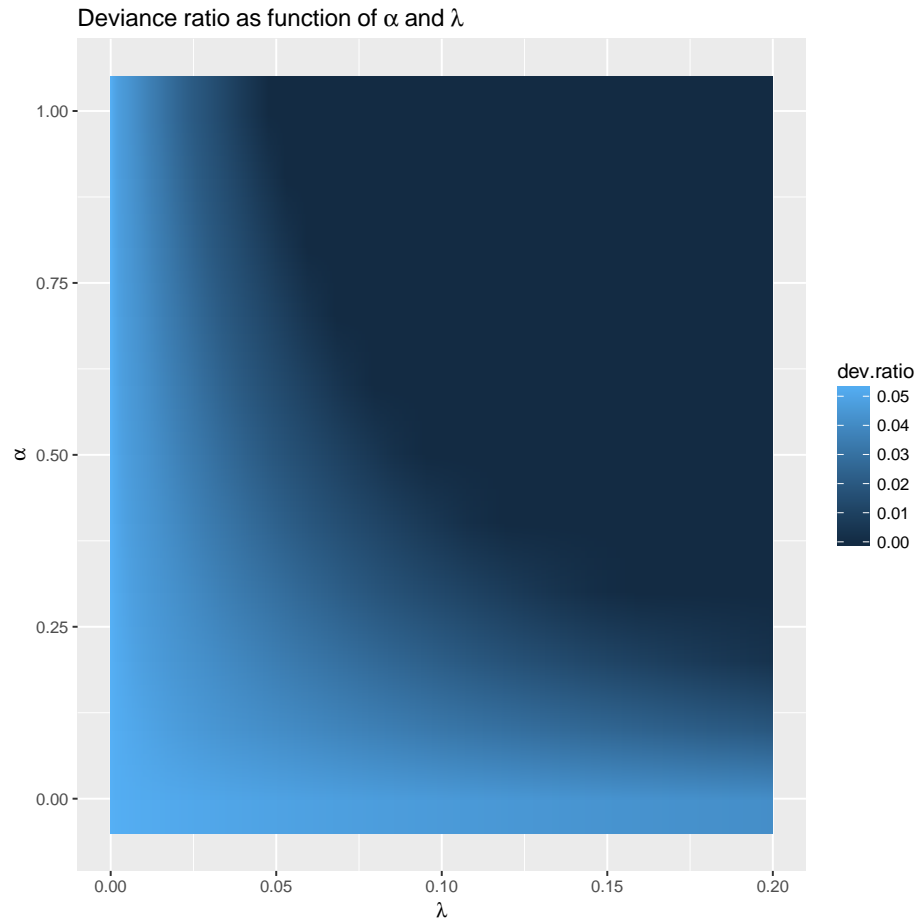
Penalized Cox Proportional Hazards Model

For the purpose of feature selection, we fit a series of penalized Cox models to the training data.

We used an elastic net penalty– a penalty term that is a linear combination of the l_1 and l_2 penalties.

$$\lambda[(1 - \alpha)||\beta||_2 + \alpha||\beta||_1]$$

We fit models varying α and λ in the penalty– we selected the model with the largest evaluated value of the likelihood function.

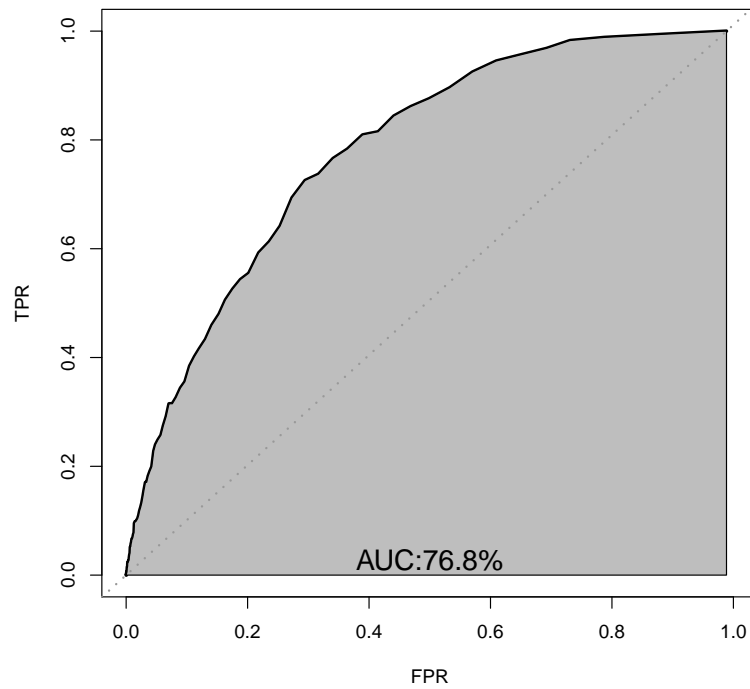


The best model, in terms had a value of λ very close to 0, and α very close to 0 (the ridge penalty). Ninety-seven variables of the original 201 had non-zero coefficients.

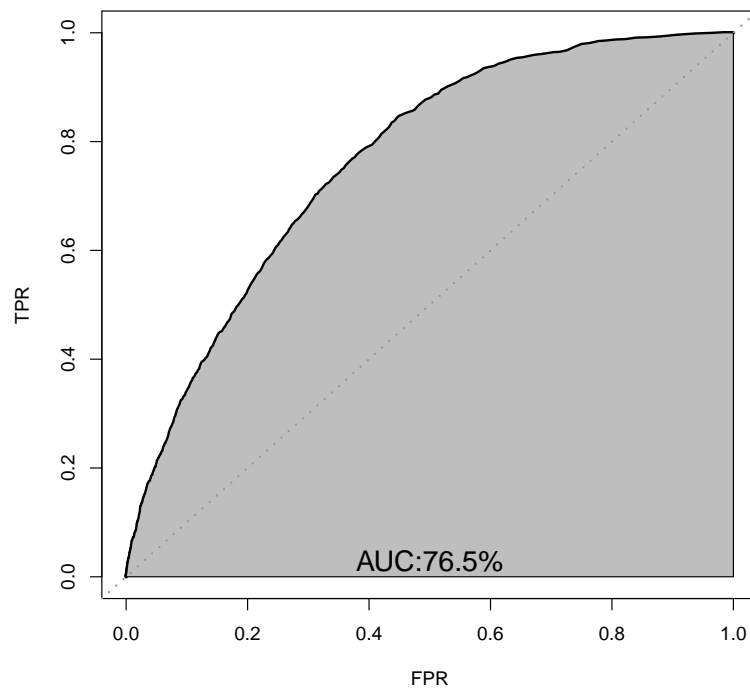
One Year and Five Year Prediction of default (out of sample)

The below figures show the out of sample performance of the one and five year probabilities estimated by the Cox model:

ROC curve for 1 year ahead default predictions



ROC curve for 5 year ahead default predictions



Portfolio Selection

For the next part of the project, we considered a portfolio of 500 loans selected from the withheld test data set. These loans met the following criteria:

1. Loans that had not defaulted as of 02-01-2010.
2. Loans that were approved before 02-01-2010.
3. Loans less than 15 years old.

These conditions were to ensure that the 500 loans in question were active as of the portfolio construction date, which we determined to be 02-01-2010. The 15 year age limit was so that estimation of 5 year ahead default probabilities would be valid.

Modeling Loss at Default

Value-at-Risk

- ALEX

Average Value-at-Risk

- ALEX

Loss Distributions by Tranche

- BEN