

# Stat 652 Project

*Blake Shurtz*

*February 6, 2019*

The goal of the project is to predict the **loan\_status** of a given loan. A background introduction to the data set is provided. An exploratory analysis of the target variable yields a highly accurate null model of a subset of the data, which also reduces the remaining levels of the target to two. The dimension of the data is reduced through a walkthrough of its 144 features. Statistical and machine learning models are executed on training data and evaluated on test data. Compared to the [top model accuracy of 0.841](#), the models herein have predictive accuracies of between 96-98%.

## Introduction

Data on loans comes from the Lending Club website at <https://www.lendingclub.com/info/download-data.action>. Data for the years 2012-2014 were downloaded, merged and cleaned. For a record of this process, see the rPub [here](#).

Table 1: Dimension of Data

Dimension	Value
Rows	423812
Columns	145

## Exploratory Analysis

Table 2: Target Levels

Level	Freq	Description
Charged Off	70613	<i>The original creditor has given up on being repaid according to the terms of the loan</i>
Current	14585	<i>The loan is open and is being re-paid</i>
Default	1	<i>The loan won't be paid</i>
Fully Paid	337704	<i>The load has been paid</i>
In Grace Period	362	<i>The loan is late, &lt; 13 days</i>
Late (16-30 days)	126	<i>The next payment is past due by 13-30 days</i>
Late (31-120 days)	417	<i>The loan is late by 31-120 days</i>

The two main categories as represented by their frequency are *Fully Paid* and *Charged-Off*.

The third category by size is *Current*. We'll return to category this momentarily.

The last 3 categories, *In Grace Period*, *Late (13-30 days)*, and *Late (31-120 days)* are chronologically ordered levels between *Current* and *Charged-Off*. If a *Current* status loan is not repaid by the due date, it enters *In Grace Period*, then the two categories of lateness, then after 120 days the loan is categorized as *Charged-Off* then *Default*.

## “loan\_status = Current” Data Subset

Exploratory analysis indicates that loans classified as *Current* consistently have non-zero values for **out\_prncp**, defined as “Remaining outstanding principal for total amount funded”. In other words, current loans still have outstanding debt repayments.

As it turns out, almost all (14573/14585) *Current* loans have **out\_prncp** greater than 0.

Interestingly, only 903/409223 non-*Current* observations also have a non-zero value for this variable. It turns out that the remaining 903 observations are assigned to other categories of late-payers.

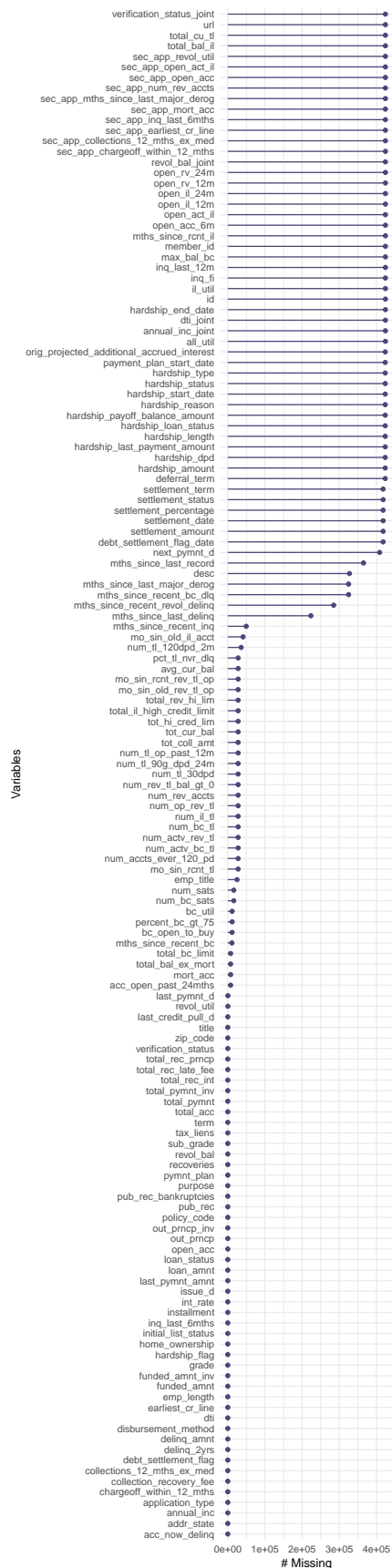
In other words, **out\_prnc** > 0 belongs exclusively to 5 levels of the target. Furthermore, the levels are overwhelmingly represented by the loan class of *Current*.

Therefore, by subsetting the data on the basis of **out\_prncp** > 0 and applying the null model that all observations are *Current*, I can achieve 94% accuracy on the subsetting data (comprising 3% of the total data). This also suggests that I can exclude **out\_prncp** from subsequent models.

In addition, by eliminating 5 of the 7 response categories, I’m left with only two categories to focus on: *Charged-Off* and *Fully Paid*. Thus, the target **loan\_status** becomes a binary variable, which will allow me to apply binary classification methods.

## Missing Data

Many variables are missing most or all of the data. See the figure on the next page. There is a pretty big drop-off between **mths\_since\_last\_delinq** and **mths\_since\_recent\_inq**, so I’ll draw the line there. In all, 58 features were removed.



## Dimension Reduction

Before building the models, I have to reduce the number of features in my data. Otherwise, the ML models will take too long to compile.

### Long Right-Tailed Variables

There are a number of long right-tailed variables in the data set, which I define as a variable where the vast majority of the observations are 0, with a right-tail distribution that comprises around 5%-10% of the remaining observations.

The correlation between each of these features and the target were measured. 11 features with correlations of  $< 0.05$  were removed.

### Highly Correlated Predictors

Just as there are features that are correlated with the target, there are also groups of features that are highly correlated with one another. The dimension of the data can be reduced by removing highly correlated predictors while maintaining the number of unique groups.

Four features can be removed: the **term** is reflected by the **int\_rate**, so I'll remove **term**. **grade** and **sub\_grade** are highly correlated, I'll remove **sub\_grade**. **loan\_amnt/funded\_amnt/funded\_amnt\_inv** are similar, so I'll remove the latter two.

### Date Variables

There are four time series variables in the data set: **issue\_d**, **earliest\_cr\_line**, **last\_credit\_pull\_d** and **last\_pymnt\_d**. While time series variables can also be useful in an exploratory analysis, none are correlated with the target.

### Location Variables

There are two variables that have data on the location of the loan: **zip\_code** and **addr\_state**, both of which are factors that have 50+ levels. The high number of levels will greatly slowdown the ML algorithm, so they will be excluded.

## What Remaining Features (Probably) Don't Matter?

The remaining features were checked against a [data dictionary](#) to decide a priori whether or not they are predictive of the target. While remaining open to surprises, the application of reasonable doubt to the features in question is a way to improve model performance with a low hypothesized Type 1 error. When in doubt, correlations were checked.

**acc\_open\_past\_24mths**: Number of trades opened in past 24 months **Keep**  
**annual\_inc**: Self-reported income **Keep**  
**avg\_cur\_bal**: How much is owed. Debt. **Keep**  
**bc\_open\_to\_buy**: Seems important to retail, not sure why it's here. **Drop**  
**bc\_util**: Current balance/credit limit. May be indicative of debt. **Keep**  
**collection\_recovery\_fee**: Record of a small collection fee. Moderate correlation. **Keep**  
**debt\_settlement\_flag**: Working with a debt-settlement company. Moderate correlation. **Keep**  
**dti**: Debt-to-income **Keep**  
**initial\_list\_status**: Whole or fractional loan. Weak correlation. **Drop**  
**installment**: Monthly amount owed. Weak correlation. **Drop**  
**last\_pymnt\_amnt**: Moderate correlation. **Keep**  
**mo\_sin\_old\_xxx**: Indicative of financial status at time of loan. NA correlation. **Drop**  
**mort\_acc**: Number of mortgage accounts. **Keep**  
**id** and **member\_id** are not predictors. **Drop** **emp\_length** (employment length) has been tidied. **\*\*Keep\***

## Training and Test Data

I will divide the data up using the *createDataPartition* function in the **caret** package, “a series of test/training partitions are created using createDataPartition while createResample creates one or more bootstrap samples.” (CRAN) The training data comprises 75% of the original data set.

```
set.seed(1234)
dindex <- createDataPartition(d$loan_status, p=.75, list=FALSE)
train <- d[dindex,]; test <- d[-dindex,]
```

## Step 3: Training a model on the data

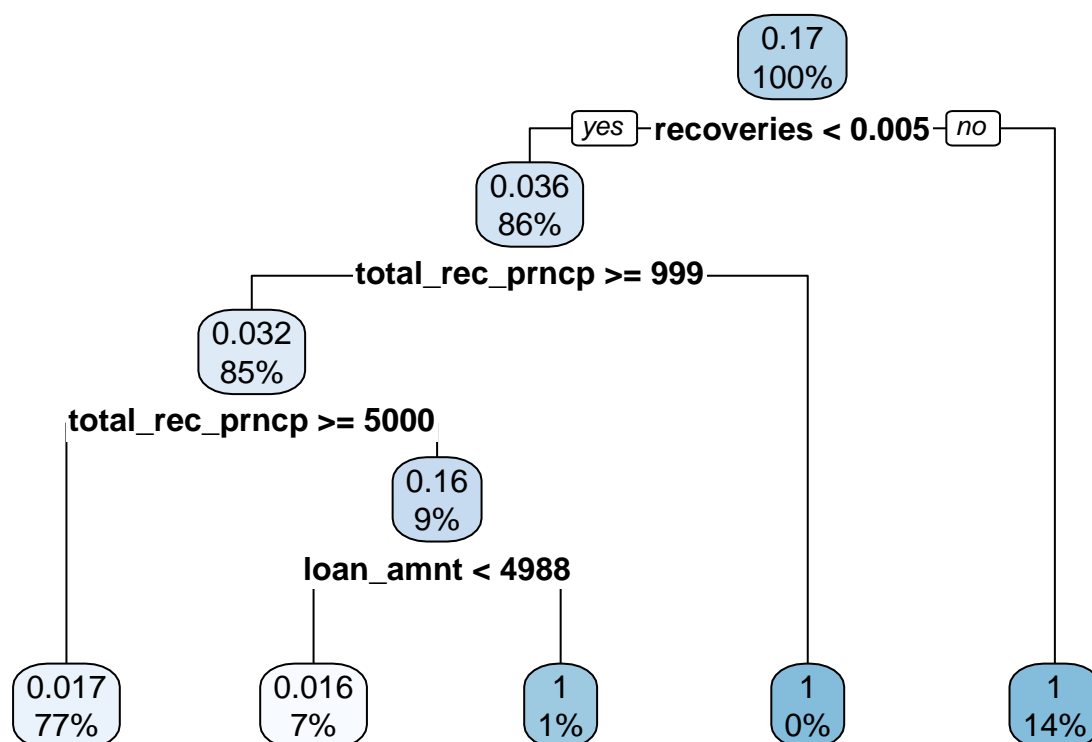
### Model 0: Null Model

First, I'll impose the null model by calling 5/7 classifiers *Current*. Subsequent models are conditional on this subset of the data. The remaining subset is called below:

```
traincat <- train %>% filter(loan_status == "Charged Off" | loan_status == "Fully Paid")
traincat$loan_status <- as.numeric(traincat$loan_status) #recoding
traincat$loan_status <- ifelse(traincat$loan_status==1, 1, 0) #0 is paid, 1 is charged
```

### Model 1: Classification Tree

The classification tree model using recursive partitioning executed properly with all of the features.



The model mistakenly reports *Fully Paid* predictions as slightly greater than 0. Other than that, **recoveries** (which has a correlation of 0.5 with the outcome) is the sole predictor of whether someone has charged off, that is, has not paid their loan.

**recoveries** is defined as “post charge off gross recovery,” suggesting that it represents a certain amount of written-off debt. This makes a lot of sense. No loans that have been fully paid would need a gross recovery, and loans that have been charged off would have some partial forgiveness.

The receipt of payments on the principal **total\_rec\_prncp** indicates that a loan is paid off. In short, paying principal is the best predictor that the loan will become *Fully Paid*.

## Model 2: Logistic Regression

A logistic regression works as a method of classification. I tried to fit as many remaining predictors as possible, while still allowing the model to converge.

```
logit.mod <- glm(loan_status ~  
  ###low correlation  
  verification_status + acc_open_past_24mths +  
  avg_cur_bal + mo_sin_rcnt_rev_tl_op + int_rate + grade + emp_length +  
  num_accts_ever_120_pd + num_actv_bc_tl + num_actv_rev_tl + num_bc_sats +  
  num_bc_tl + pct_tl_nvr_dlq + loan_amnt +  
  ###moderate correlation  
  debt_settlement_flag #(r= 0.27)  
  , data=traincat, family = binomial)
```

Four moderately correlated predictors lead the model to overfit: **recoveries** (0.53), **last\_pymnt\_amnt** (-.29), **collection\_recovery\_fee** (0.46), **total\_rec\_prncp**(-0.39). R gives the warning: “glm.fit: fitted probabilities numerically 0 or 1 occurred.”

Attempts to add additional weaker predictors to the model caused the model to fail to converge and exclude observations: **home\_ownership** + **annual\_inc** + **purpose** + **dti** + **inq\_last\_6mths** + **open\_acc** + **revol\_bal** + **revol\_util** + **total\_acc** + **total\_pymnt** + **total\_pymnt\_inv** + **total\_rec\_int** + **total\_rec\_late\_fee** + **tot\_coll\_amt** + **tot\_cur\_bal** + **total\_rev\_hi\_lim** + **total\_bc\_limit** + **total\_il\_high\_credit\_limit** + **num\_il\_tl** + **num\_op\_rev\_tl** + **num\_rev\_accts** + **num\_rev\_tl\_bal\_gt\_0** + **mo\_sin\_rcnt\_tl** + **mths\_since\_recent\_bc** + **mths\_since\_recent\_inq** + **percent\_bc\_gt\_75** + **tot\_hi\_cred\_lim** + **num\_sats** + **num\_tl\_120dpd\_2m** + **num\_tl\_30dpd** + **num\_tl\_90g\_dpd\_24m** + **num\_tl\_op\_past\_12m** + **total\_bal\_ex\_mort**.

### Model 3: Logistic (Bayesian)

Given the failure of the logit model to include all moderately correlated predictors, a Bayesian logistic model with adaptive priors will be introduced.

There were two executions of the model. The first execution ran 1000 iterations on the full training data set and the second ran 500 iterations on 25% of the training data (in order to decrease execution time).

```
#wrangle
traincatsmall$debt_settlement_flag <- as.numeric(traincatsmall$debt_settlement_flag)
set.seed(1234)
options(mc.cores = parallel::detectCores())
Sys.setenv(LOCAL_CPPFLAGS = '-march=native')
logit.bayes.mod <- map2stan(
  alist(
    loan_status ~ dbinom( 1 , p ) ,
    logit(p) <- a + b * debt_settlement_flag + c * recoveries +
              d * last_pymnt_amnt + e * collection_recovery_fee +
              f * total_rec_prncp,

    #adaptive priors
    a ~ dnorm( 0 , 1.5 ),
    b ~ dnorm(b_mu, b_sigma),
    b_mu ~ dnorm(0,1),
    b_sigma ~ dcauchy(0,2),
    c ~ dnorm(c_mu, c_sigma),
    c_mu ~ dnorm(0,1),
    c_sigma ~ dcauchy(0,2),
    d ~ dnorm(d_mu, d_sigma),
    d_mu ~ dnorm(0,1),
    d_sigma ~ dcauchy(0,2),
    e ~ dnorm(e_mu, e_sigma),
    e_mu ~ dnorm(0,1),
    e_sigma ~ dcauchy(0,2),
    f ~ dnorm(f_mu, f_sigma),
    f_mu ~ dnorm(0,1),
    f_sigma ~ dcauchy(0,2)
  ) , data=traincat,
  iter=1000, warmup=250, chains=1, cores=4)
```

Model convergence was poor with a large number of divergent samples. However, coefficients were calculated for all predictors. The regression model is below:

$$\log \frac{p}{1-p} = 0.4756 - 0.762 * debtsettlementflag + 2.434 * recoveries - 0.0008 * lastpaymentamount - 1.641 * collectionrecoveryfee - 0.0002 * totalrecprncp$$



## Step 4: Evaluating Model Performance

Firstly, I have to subset the data in order to apply the null model **loan\_status** = *Current* to 5/7 levels. For each model evaluation, I will add this subset back into the final analysis after running the remaining data through the algorithm.

Given that there are more than two levels of the outcome, accuracy is the preferred method of evaluation, along with confusion matrices. Note that for confusion matrices, the level 3 *Default (n=1)* was excluded. As a result of the programming, the level *Fully Paid* may be represented by either a 0 or a 4.

### Null Model Evaluation

After the data set is subsetting, a null model would assume that all loans are *Fully Paid*. This leads to an overall accuracy of 83.1%.

### CART Model Evaluation

The CART model has an overall accuracy of 98%.

Table 3: Confusion Matrix- CART Model

	1	2	4	5	6	7
1	16179	0	0	0	0	0
2	0	3646	0	90	31	104
4	1474	0	84425	0	0	0

### Logit Evaluation

The logit model has an overall accuracy of 79%.

Table 4: Confusion Matrix- Logit Model

	0	1	2	5	6	7
0	78532	14987	0	0	0	0
1	43	1528	0	0	0	0
2	0	0	3646	90	31	104

## Logit Bayes Evaluation

The logit Bayes model had an accuracy of 84% and 96%. The model that executed on a larger proportion of the training data with more iterations was more accurate.

Table 5: Confusion Matrix- Bayes Logit

	0	1	2	5	6	7
0	84426	3544	0	0	0	0
1	0	2917	0	0	0	0
2	0	0	3646	90	31	104

## Conclusion

On the entryway to Bruce Lee’s dojo, there was a sign that read: “Reject what is useless. Accept what is useful.” Bruce Lee is remembered as a martial artist who judged his ability as a fighter against all possibilities, ie. randomness. He synthesized what worked and what didn’t into his own system of classification.

There are many models that I could have used in this project. But, because I have achieved 95%+ accuracy with both a statistical learning and machine learning algorithms, I will consider the assignment a success. I also put a lot of credit on the model success on a thorough exploratory analysis.