

Predicting Bad Lending Club Loans for Fixed Loan Grades with Multiple Different Models

Bill Anderson (william.david.anderson@gmail.com)

17 May 2016

Introduction and Executive Summary

This document presents an analysis of lending club data for loans issued between June 2007 and December 2011, with the goal of predicting which loans will go “bad” (i.e., the borrower misses a payment or defaults). This analysis is done with the loan grade held constant (e.g., analysis for all A loans, analysis for all B loans, etc.), which can be useful; for example, if we could identify all the grade D loans that would not go bad, we would have the best of both worlds: high interest rates, but no risk of loss from default. For this study, loans with grade A, B, C, and D were considered (not enough data for grade E loans). Also, this study used four different model types: logistic regression, random forest, gradient boost, and support vector machines. The results from the different models were similar, although the gradient boost results were slightly more accurate (at least for the given random number seed).

For the grade C and D loans (the ones with the most defaults), we can correctly identify over 65% of the loans that will go bad. Also for these same loan grades, all four of the models identified the same two predictors that were most important in predicting which loans will go bad: FICO score and the number of credit inquiries in the past six months.

Details on these and other results are shown below.

Data Ingest and Initialization Steps

```
# read in the lending club data
setwd("/Users/andersnb/lending-club/my-analysis")
loans <- read.csv("../data/LoanStats3a_securev1.csv")
str(loans)
```

```
## 'data.frame':    42536 obs. of  115 variables:
## $ id              : Factor w/ 42536 levels "1000007","1000030",...: 4388 4387 4386 4385 ...
## $ member_id       : int   1296599 1314167 1313524 1277178 1311748 1311441 1304742 1288...
## $ loan_amnt        : int   5000 2500 2400 10000 3000 5000 7000 3000 5600 5375 ...
## $ funded_amnt      : int   5000 2500 2400 10000 3000 5000 7000 3000 5600 5375 ...
## $ funded_amnt_inv  : num   4975 2500 2400 10000 3000 ...
## $ term             : Factor w/ 3 levels "", " 36 months",...: 2 3 2 2 3 2 3 2 3 3 ...
## $ int_rate         : Factor w/ 395 levels "", " 5.42%"," 5.79%",...: 80 223 241 162 13...
## $ installment      : num   162.9 59.8 84.3 339.3 67.8 ...
## $ grade            : Factor w/ 8 levels "", "A","B","C",...: 3 4 4 4 3 2 4 6 7 3 ...
## $ sub_grade        : Factor w/ 36 levels "", "A1","A2","A3",...: 8 15 16 12 11 5 16 22 2...
## $ emp_title        : Factor w/ 30661 levels "", " old palm inc",...: 1 22922 1 791 2823...
## $ emp_length       : Factor w/ 13 levels "", "< 1 year",...: 4 2 4 4 3 6 11 12 7 2 ...
## $ home_ownership   : Factor w/ 6 levels "", "MORTGAGE",...: 6 6 6 6 6 6 6 5 6 ...
## $ annual_inc       : num   24000 30000 12252 49200 80000 ...
## $ verification_status : Factor w/ 4 levels "", "Not Verified",...: 4 3 2 3 3 3 2 3 3 4 ...
```

```

## $ issue_d : Factor w/ 56 levels "", "Apr-2008",...: 15 15 15 15 15 15 15 15 15 15 ...
## $ loan_status : Factor w/ 10 levels "", "Charged Off",...: 7 2 7 7 3 7 3 7 2 2 ...
## $ pymnt_plan : Factor w/ 3 levels "", "n", "y": 2 2 2 2 2 2 2 2 2 2 ...
## $ url : Factor w/ 42536 levels "", "https://www.lendingclub.com/browse/loan...: 1 1 1 1 1 1 1 1 1 1 ...
## $ desc : Factor w/ 28965 levels "", "\t Member# 809768, loan description. (": 1 1 1 1 1 1 1 1 1 1 ...
## $ purpose : Factor w/ 15 levels "", "car", "credit_card",...: 3 2 13 11 11 15 4 4 4 4 ...
## $ title : Factor w/ 21267 levels "", "\tdebt consolidation",...: 3687 1869 177 177 177 177 177 177 177 177 ...
## $ zip_code : Factor w/ 838 levels "", "007xx", "010xx",...: 728 282 514 765 814 765 765 765 765 765 ...
## $ addr_state : Factor w/ 51 levels "", "AK", "AL", "AR",...: 5 12 16 6 38 5 29 6 6 4 ...
## $ dti : num 27.65 1 8.72 20 17.94 ...
## $ delinq_2yrs : int 0 0 0 0 0 0 0 0 0 0 ...
## $ earliest_cr_line : Factor w/ 531 levels "", "Apr-1964",...: 194 36 431 163 205 434 256 256 256 256 ...
## $ fico_range_low : int 735 740 735 690 695 730 690 660 675 725 ...
## $ fico_range_high : int 739 744 739 694 699 734 694 664 679 729 ...
## $ inq_last_6mths : int 1 5 2 1 0 3 1 2 2 0 ...
## $ mths_since_last_delinq : int NA NA NA 35 38 NA NA NA NA NA ...
## $ mths_since_last_record : int NA NA NA NA NA NA NA NA NA NA ...
## $ open_acc : int 3 3 2 10 15 9 7 4 11 2 ...
## $ pub_rec : int 0 0 0 0 0 0 0 0 0 0 ...
## $ revol_bal : int 13648 1687 2956 5598 27783 7963 17726 8221 5210 9279 ...
## $ revol_util : Factor w/ 1120 levels "", "0.01%", "0.03%",...: 943 1012 1105 204 594 594 594 594 594 594 ...
## $ total_acc : int 9 4 10 37 38 12 11 4 13 3 ...
## $ initial_list_status : Factor w/ 2 levels "", "f": 2 2 2 2 2 2 2 2 2 2 ...
## $ out_prncp : num 0 0 0 0 707 ...
## $ out_prncp_inv : num 0 0 0 0 707 ...
## $ total_pymnt : num 5863 1009 3006 12232 3310 ...
## $ total_pymnt_inv : num 5834 1009 3006 12232 3310 ...
## $ total_rec_prncp : num 5000 456 2400 10000 2293 ...
## $ total_rec_int : num 863 435 606 2215 1017 ...
## $ total_rec_late_fee : num 0 0 0 17 0 ...
## $ recoveries : num 0 117 0 0 0 ...
## $ collection_recovery_fee : num 0 1.11 0 0 0 0 0 0 2.09 2.52 ...
## $ last_pymnt_d : Factor w/ 100 levels "", "Apr-2008",...: 43 7 59 43 35 43 35 43 6 8 ...
## $ last_pymnt_amnt : num 171.6 119.7 649.9 357.5 67.8 ...
## $ next_pymnt_d : Factor w/ 102 levels "", "Apr-2008",...: 1 1 1 1 70 1 10 1 1 1 ...
## $ last_credit_pull_d : Factor w/ 105 levels "", "Apr-2009",...: 35 103 35 43 35 44 35 25 14 14 ...
## $ last_fico_range_high : int 719 534 679 579 674 679 644 689 499 499 ...
## $ last_fico_range_low : int 715 530 675 575 670 675 640 685 0 0 ...
## $ collections_12_mths_ex_med : int 0 0 0 0 0 0 0 0 0 0 ...
## $ mths_since_last_major_derog : logi NA NA NA NA NA NA NA ...
## $ policy_code : int 1 1 1 1 1 1 1 1 1 1 ...
## $ application_type : Factor w/ 2 levels "", "INDIVIDUAL": 2 2 2 2 2 2 2 2 2 2 ...
## $ annual_inc_joint : logi NA NA NA NA NA NA NA ...
## $ dti_joint : logi NA NA NA NA NA NA NA ...
## $ verification_status_joint : logi NA NA NA NA NA NA NA ...
## $ acc_now_delinq : int 0 0 0 0 0 0 0 0 0 0 ...
## $ tot_coll_amt : logi NA NA NA NA NA NA NA ...
## $ tot_cur_bal : logi NA NA NA NA NA NA NA ...
## $ open_acc_6m : logi NA NA NA NA NA NA NA ...
## $ open_il_6m : logi NA NA NA NA NA NA NA ...
## $ open_il_12m : logi NA NA NA NA NA NA NA ...
## $ open_il_24m : logi NA NA NA NA NA NA NA ...
## $ mths_since_rcnt_il : logi NA NA NA NA NA NA NA ...
## $ total_bal_il : logi NA NA NA NA NA NA NA ...

```

```
## $ il_util : logi NA NA NA NA NA NA ...
## $ open_rv_12m : logi NA NA NA NA NA NA ...
## $ open_rv_24m : logi NA NA NA NA NA NA ...
## $ max_bal_bc : logi NA NA NA NA NA NA ...
## $ all_util : logi NA NA NA NA NA NA ...
## $ total_rev_hi_lim : logi NA NA NA NA NA NA ...
## $ inq_fi : logi NA NA NA NA NA NA ...
## $ total_cu_tl : logi NA NA NA NA NA NA ...
## $ inq_last_12m : logi NA NA NA NA NA NA ...
## $ acc_open_past_24mths : logi NA NA NA NA NA NA ...
## $ avg_cur_bal : logi NA NA NA NA NA NA ...
## $ bc_open_to_buy : logi NA NA NA NA NA NA ...
## $ bc_util : logi NA NA NA NA NA NA ...
## $ chargeoff_within_12_mths : int 0 0 0 0 0 0 0 0 0 ...
## $ delinq_amnt : int 0 0 0 0 0 0 0 0 0 ...
## $ mo_sin_old_il_acct : logi NA NA NA NA NA NA ...
## $ mo_sin_old_rev_tl_op : logi NA NA NA NA NA NA ...
## $ mo_sin_rcnt_rev_tl_op : logi NA NA NA NA NA NA ...
## $ mo_sin_rcnt_tl : logi NA NA NA NA NA NA ...
## $ mort_acc : logi NA NA NA NA NA NA ...
## $ mths_since_recent_bc : logi NA NA NA NA NA NA ...
## $ mths_since_recent_bc_dlq : logi NA NA NA NA NA NA ...
## $ mths_since_recent_inq : logi NA NA NA NA NA NA ...
## $ mths_since_recent_revol_delinq : logi NA NA NA NA NA NA ...
## $ num_accts_ever_120_pd : logi NA NA NA NA NA NA ...
## $ num_actv_bc_tl : logi NA NA NA NA NA NA ...
## $ num_actv_rev_tl : logi NA NA NA NA NA NA ...
## $ num_bc_sats : logi NA NA NA NA NA NA ...
## $ num_bc_tl : logi NA NA NA NA NA NA ...
## $ num_il_tl : logi NA NA NA NA NA NA ...
## [list output truncated]
```

```
# initialize random number generator
set.seed(1)
```

Data Cleaning

In this section, we convert data types, get rid of unneeded data, etc.

```
#
# Loans in the dataset were issued at different times and have terms of 3 or 5 years.
# We want all loans to have the same chance to reach maturity or the results could be
# misleading. Consider an extreme case where a loan is issued the month before the end
# of when data is collected. The loan is less likely to be in default after just one
# month than if it's been outstanding for 3 (or 5) years and such loans could result in
# misleading interpretations. Thus, since this dataset ends at Feb 2016, we should only
# consider loans that were issued 5 years or more ago, or that were issued Feb 2011 or
# earlier.
#
loans <- filter(loans, issue_d != "")
loans$issue_d <- factor(loans$issue_d)
loans$issue_d <- parse_date_time(paste("01-", loans$issue_d), "%d-%b-%Y")
```

```

loans <- filter(loans, issue_d <= "2011-02-01")

#
# convert to a date type
#
loans <- filter(loans, last_pymnt_d != "")
loans$last_pymnt_d <- factor(loans$last_pymnt_d)
loans$last_pymnt_d <- parse_date_time(paste("01-", loans$last_pymnt_d), "%d-%b-%Y")

#
# convert to a date type
#
loans <- filter(loans, earliest_cr_line != "")
loans$earliest_cr_line <- factor(loans$earliest_cr_line)
loans$earliest_cr_line <- parse_date_time(paste("01-", loans$earliest_cr_line), "%d-%b-%Y")

#
# convert to a date type
#
loans <- filter(loans, last_credit_pull_d != "")
loans$last_credit_pull_d <- factor(loans$last_credit_pull_d)
loans$last_credit_pull_d <- parse_date_time(paste("01-", loans$last_credit_pull_d), "%d-%b-%Y")

# get rid of empty factor
loans <- filter(loans, term != "")
loans$term <- factor(loans$term)

# convert interest rate from string to float
loans$int_rate <- gsub("%", "", loans$int_rate)
loans$int_rate <- gsub(" ", "", loans$int_rate)
loans$int_rate <- as.numeric(loans$int_rate)

# get rid of empty factor
loans <- filter(loans, grade != "")
loans$grade <- factor(loans$grade)

# get rid of empty factor
loans <- filter(loans, sub_grade != "")
loans$sub_grade <- factor(loans$sub_grade)

# get rid of empty factor
loans <- filter(loans, emp_length != "")
loans$emp_length <- factor(loans$emp_length)

# get rid of empty factor
loans <- filter(loans, home_ownership != "")
loans$home_ownership <- factor(loans$home_ownership)

# get rid of empty factor
loans <- filter(loans, verification_status != "")
loans$verification_status <- factor(loans$verification_status)

# get rid of empty factor

```

```

loans <- filter(loans, pymnt_plan != "")
loans$pymnt_plan <- factor(loans$pymnt_plan)

# create a variable that's true if the desc is empty, else false
loans <- mutate(loans, desc_empty = as.factor(ifelse(desc == "", 1, 0)))

# get rid of empty factor
loans <- filter(loans, purpose != "")
loans$purpose <- factor(loans$purpose)

# get rid of empty factor
loans <- filter(loans, zip_code != "")
loans$zip_code <- factor(loans$zip_code)

# get rid of empty factor
loans <- filter(loans, addr_state != "")
loans$addr_state <- factor(loans$addr_state)

# convert revol_util from a factor to a numeric variable
loans$revol_util <- as.numeric(gsub("%", "", loans$revol_util))

# get rid of empty factor
loans <- filter(loans, initial_list_status != "")
loans$initial_list_status <- factor(loans$initial_list_status)

#
# the following columns are deemed not useful (for the following reasons) so we exclude them:
# mths_since_last_major_derog      (all NAs)
# annual_inc_joint                  (all NAs)
# dti_joint                         (all NAs)
# verification_status_joint        (all NAs)
# tot_coll_amt                     (all NAs)
# tot_cur_bal                      (all NAs)
# open_acc_6m                      (all NAs)
# open_il_6m                       (all NAs)
# open_il_12m                      (all NAs)
# open_il_24m                      (all NAs)
# mths_since_rcnt_il              (all NAs)
# total_bal_il                    (all NAs)
# il_util                         (all NAs)
# open_rv_12m                     (all NAs)
# open_rv_24m                     (all NAs)
# max_bal_bc                      (all NAs)
# all_util                        (all NAs)
# total_rev_hi_lim                (all NAs)
# inq_fi                          (all NAs)
# total_cu_tl                     (all NAs)
# inq_last_12m                    (all NAs)
# acc_open_past_24mths            (all NAs)
# avg_cur_bal                     (all NAs)
# bc_open_to_buy                  (all NAs)
# bc_util                         (all NAs)
# mo_sin_old_il_acct              (all NAs)

```

```

# mo_sin_old_rev_tl_op      (all NAs)
# mo_sin_rcnt_rev_tl_op    (all NAs)
# mo_sin_rcnt_tl           (all NAs)
# mort_acc                 (all NAs)
# mths_since_recent_bc     (all NAs)
# mths_since_recent_bc_dlq (all NAs)
# mths_since_recent_inq    (all NAs)
# mths_since_recent_revol_delinq (all NAs)
# num_accts_ever_120_pd    (all NAs)
# num_actv_bc_tl          (all NAs)
# num_actv_rev_tl         (all NAs)
# num_bc_sats             (all NAs)
# num_bc_tl              (all NAs)
# num_il_tl              (all NAs)
# num_op_rev_tl          (all NAs)
# num_rev_accts          (all NAs)
# num_rev_tl_bal_gt_0    (all NAs)
# num_sats               (all NAs)
# num_tl_120dpd_2m       (all NAs)
# num_tl_30dpd           (all NAs)
# num_tl_90g_dpd_24m     (all NAs)
# num_tl_op_past_12m     (all NAs)
# pct_tl_nvr_dlq         (all NAs)
# percent_bc_gt_75       (all NAs)
# tot_hi_cred_lim        (all NAs)
# total_bal_ex_mort      (all NAs)
# total_bc_limit         (all NAs)
# total_il_high_credit_limit (all NAs)
# next_pymnt_d           (doesn't seem relevant to loan status and contained a lot of missing data)
# mths_since_last_delinq  (a very large number of NAs)
# mths_since_last_record  (a very large number of NAs)
# id                     (not relevant to loan status)
# member_id              (not relevant to loan status)
# url                    (url for the loan details; not relevant to loan status)
# desc                   (it's possible the information contained in the desc. could be useful; f
# title                  (it's possible the information contained in the title could be useful; f
# emp_title              (it's possible the information contained in emp_title could be useful; f
#

```

```

loans <- subset(loans, select = -c(mths_since_last_major_derog,
  annual_inc_joint, dti_joint, verification_status_joint, tot_coll_amt,
  tot_cur_bal, open_acc_6m, open_il_6m, open_il_12m, open_il_24m,
  mths_since_rcnt_il, total_bal_il, il_util, open_rv_12m, open_rv_24m,
  max_bal_bc, all_util, total_rev_hi_lim, inq_fi, total_cu_tl,
  inq_last_12m, acc_open_past_24mths, avg_cur_bal, bc_open_to_buy,
  bc_util, mo_sin_old_il_acct, mo_sin_old_rev_tl_op, mo_sin_rcnt_rev_tl_op,
  mo_sin_rcnt_tl, mort_acc, mths_since_recent_bc, mths_since_recent_bc_dlq,
  mths_since_recent_inq, mths_since_recent_revol_delinq, num_accts_ever_120_pd,
  num_actv_bc_tl, num_actv_rev_tl, num_bc_sats, num_bc_tl,
  num_il_tl, num_op_rev_tl, num_rev_accts, num_rev_tl_bal_gt_0,
  num_sats, num_tl_120dpd_2m, num_tl_30dpd, num_tl_90g_dpd_24m,
  num_tl_op_past_12m, pct_tl_nvr_dlq, percent_bc_gt_75, tot_hi_cred_lim,
  total_bal_ex_mort, total_bc_limit, total_il_high_credit_limit,

```

```

next_pymnt_d, mths_since_last_delinq, mths_since_last_record,
id, member_id, url, desc, title, emp_title))

# create binary status variable; note: I define as 'bad' any
# loan that is not current or not fully paid
loans <- mutate(loans, status = factor(ifelse(loan_status ==
  "Current" | loan_status == "Fully Paid", "good", "bad"),
  levels = c("good", "bad")))

```

Exploratory Plots

In this section, we create exploratory plots and/or tables for each variable to help determine which variables are likely to have an effect on the loan status and, thus, should be used in the subsequent models. Note: to generate the various plots, set the `explPlots` and/or the `collScatterPlots` variables at the beginning of the R markdown document to `TRUE`.

```

# create exploratory plots
createExplPlots <- function(dft) {
  for (i in 1:ncol(dft)) {
    varname = names(dft)[i]
    print(paste(varname, ":"))

    if (varname == "annual_inc") {
      # annual income requires a limit of 200000 since there are
      # some outliers that make the plots hard to understand or
      # visualize
      p <- ggplot(aes_string(x = varname, group = "status",
        colour = "status"), data = dft)
      p <- p + geom_density() + xlab(varname)
      print(p)

      p <- ggplot(dft, aes_string(x = "status", y = varname)) +
        geom_boxplot() + ylab(varname) + ylim(0, 2e+05)
      print(p)
    } else if (varname == "delinq_2yrs") {
      # delinq_2yrs requires a limit of 5 since there are some
      # outliers that make the plots hard to understand
      p <- ggplot(aes_string(x = varname, group = "status",
        colour = "status"), data = dft)
      p <- p + geom_density() + xlab(varname)
      print(p)

      p <- ggplot(dft, aes_string(x = "status", y = varname)) +
        geom_boxplot() + ylab(varname) + ylim(0, 5)
      print(p)
    } else {
      # create plots that don't require special limits
      p <- ggplot(aes_string(x = varname, group = "status",
        colour = "status"), data = dft)
      p <- p + geom_density() + xlab(varname)
    }
  }
}

```

```

    print(p)

    if (class(dft[[i]]) == "numeric" || class(dft[[i]]) ==
        "integer") {
      p <- ggplot(dft, aes_string(x = "status", y = varname)) +
        geom_boxplot() + ylab(names(dft)[i])
      print(p)
    } else {
      print(table(dft[[i]], dft$status))
      print(prop.table(table(dft[[i]], dft$status),
        1))
    }
  }
  cat("\n")
}

# subset data by loan grade
a_loans <- loans[loans$grade == "A", ]
b_loans <- loans[loans$grade == "B", ]
c_loans <- loans[loans$grade == "C", ]
d_loans <- loans[loans$grade == "D", ]

# create exploratory plots by loan grade
if (explPlots == TRUE) {
  createExplPlots(a_loans)
  createExplPlots(b_loans)
  createExplPlots(c_loans)
  createExplPlots(d_loans)
}

# select predictors that have an effect on response and get
# rid of rows with NAs
a_loans <- select(a_loans, c(status, term, verification_status,
  purpose, fico_range_low, fico_range_high, inq_last_6mths,
  revol_util, last_fico_range_low, last_fico_range_high, desc_empty,
  dti))
b_loans <- select(b_loans, c(status, term, verification_status,
  purpose, fico_range_low, fico_range_high, inq_last_6mths,
  revol_util, last_fico_range_low, last_fico_range_high, desc_empty,
  dti))
c_loans <- select(c_loans, c(status, term, verification_status,
  purpose, fico_range_low, fico_range_high, inq_last_6mths,
  revol_util, last_fico_range_low, last_fico_range_high, desc_empty,
  dti))
d_loans <- select(d_loans, c(status, term, verification_status,
  purpose, fico_range_low, fico_range_high, inq_last_6mths,
  revol_util, last_fico_range_low, last_fico_range_high, desc_empty,
  dti))

a_loans <- na.omit(a_loans)

```



```

b_loans <- na.omit(b_loans)
c_loans <- na.omit(c_loans)
d_loans <- na.omit(d_loans)

# now check for collinearity
checkForColl <- function(l) {
  pairs(~term + verification_status + purpose + fico_range_low +
        fico_range_high + inq_last_6mths + revol_util + last_fico_range_low +
        last_fico_range_high + desc_empty + dti, data = l)
}

if (collScatterPlots == TRUE) {
  checkForColl(a_loans)
  checkForColl(b_loans)
  checkForColl(c_loans)
  checkForColl(d_loans)
}

# the collinearity scatterplots suggest that there's is a
# correlation between fico_range_high/fico_range_low and
# between last_fico_range_low/last_fico_range_high;
# therefore, I won't use fico_range_low or
# last_fico_range_low in the models to avoid collinearity
a_loans <- select(a_loans, c(status, term, verification_status,
  purpose, fico_range_high, inq_last_6mths, revol_util, last_fico_range_high,
  desc_empty, dti))
b_loans <- select(b_loans, c(status, term, verification_status,
  purpose, fico_range_high, inq_last_6mths, revol_util, last_fico_range_high,
  desc_empty, dti))
c_loans <- select(c_loans, c(status, term, verification_status,
  purpose, fico_range_high, inq_last_6mths, revol_util, last_fico_range_high,
  desc_empty, dti))
d_loans <- select(d_loans, c(status, term, verification_status,
  purpose, fico_range_high, inq_last_6mths, revol_util, last_fico_range_high,
  desc_empty, dti))

```

Model Construction and Execution

The next section builds several model types (logistic, random forest, gradient boost, and support vector machine (SVM)), makes predictions and identifies the important variables in each model.

```

createDataForInput <- function(dft) {
  # partition the data into a training portion and test portion
  inTraining <- createDataPartition(dft$status, p = 0.75, list = FALSE)
  dft_orig <- dft
  dft_train <- dft_orig[inTraining, ]
  dft_test <- dft_orig[-inTraining, ]

  return(list(dft_train = dft_train, dft_test = dft_test))
}

```

```

# function to create logistic regression model
logRegModel <- function(dft_train, dft_test) {
  modLogReg <- train(status ~ ., data = dft_train, method = "glm")
  print(modLogReg)
  print(summary(modLogReg))
  print(varImp(modLogReg))

  testPred <- predict(modLogReg, dft_test)
  print(confusionMatrix(testPred, dft_test$status, positive = "bad"))

  testProbs <- predict(modLogReg, dft_test, type = "prob")
  rocObj <- roc(dft_test$status, testProbs[, "bad"])
  plot(rocObj, type = "S", print.thres = 0.5)
}

# function to create random forest model
rfModel <- function(dft_train, dft_test) {
  modRandFor <- train(status ~ ., data = dft_train, method = "rf")
  print(modRandFor)
  print(summary(modRandFor))
  print(varImp(modRandFor))

  testPred <- predict(modRandFor, dft_test)
  print(confusionMatrix(testPred, dft_test$status, positive = "bad"))

  testProbs <- predict(modRandFor, dft_test, type = "prob")
  rocObj <- roc(dft_test$status, testProbs[, "bad"])
  plot(rocObj, type = "S", print.thres = 0.5)
}

# function to create a gradient boost model
gbModel <- function(dft_train, dft_test) {
  modGradBoost <- train(status ~ ., data = dft_train, method = "gbm", verbose = FALSE)
  print(modGradBoost)
  print(summary(modGradBoost))
  print(varImp(modGradBoost))

  testPred <- predict(modGradBoost, dft_test)
  print(confusionMatrix(testPred, dft_test$status, positive = "bad"))

  testProbs <- predict(modGradBoost, dft_test, type = "prob")
  rocObj <- roc(dft_test$status, testProbs[, "bad"])
  plot(rocObj, type = "S", print.thres = 0.5)
}

# function to create SVM Gaussian kernel model note: I use the 'cv' method
# for resampling because the default boot method results in a lot of warning
# messages about duplicate row names and the 'cv' method yields results that
# are as accurate as the 'boot' method
svmModel <- function(dft_train, dft_test) {

```

```

modSvm <- train(status ~ ., data = dft_train, method = "svmRadial", preProc = c("center",
  "scale"), trControl = trainControl(classProbs = TRUE, method = "cv"))

print(modSvm)
print(summary(modSvm))
print(varImp(modSvm))

testPred <- predict(modSvm, dft_test)
print(confusionMatrix(testPred, dft_test$status, positive = "bad"))

testProbs <- predict(modSvm, dft_test, type = "prob")
rocObj <- roc(dft_test$status, testProbs[, "bad"])
plot(rocObj, type = "S", print.thres = 0.5)
}

```

Results for Grade A Loans

Only a small percentage (~7%) of the Grade A loans go bad, making it somewhat challenging to identify those loans, but, since there are so few, it's also less important. The results show that the four models had sensitivities (i.e., ability to correctly predict the bad loans) ranging from 9% to 26%. This predictive ability is based on a 50% probability classification cutoff. As the ROC curves show, it's possible to predict the bad loans with a higher probability, but, of course, with a higher false positive rate. The FICO range and the number of inquiries in the past 6 months were important predictors.

Logistic Regression Model

```

## Generalized Linear Model
##
## 3879 samples
##    9 predictor
##    2 classes: 'good', 'bad'
##
## No pre-processing
## Resampling: Bootstrapped (25 reps)
##
## Summary of sample sizes: 3879, 3879, 3879, 3879, 3879, 3879, ...
##
## Resampling results
##
##   Accuracy   Kappa     Accuracy SD   Kappa SD
## 0.9335321 0.2797235 0.003836129 0.03338435
##
##
## Call:
## NULL
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -4.1505  -0.3149  -0.2099  -0.1442   3.2303
##

```

```

## Coefficients:
##
##              Estimate Std. Error z value
## (Intercept)      13.7986678   2.9075901   4.746
## `term 60 months`    0.6508329   0.2844911   2.288
## `verification_statusSource Verified` -0.3583547   0.2000624  -1.791
## verification_statusVerified      -0.1064885   0.1794718  -0.593
## purposecredit_card      -0.5959759   0.3486052  -1.710
## purposedebt_consolidation -0.3283520   0.2963099  -1.108
## purposeeducational       0.4707276   0.4926210   0.956
## purposehome_improvement  -0.4675126   0.3742106  -1.249
## purposehouse      -0.3180473   0.7039476  -0.452
## purposemajor_purchase -0.7079184   0.3973222  -1.782
## purposemedical      -0.1293572   0.5060129  -0.256
## purposemoving      -0.9334389   0.5546842  -1.683
## purposeother      -0.5040392   0.3337273  -1.510
## purposerenewable_energy  0.4920301   1.3039656   0.377
## purposesmall_business -0.0894251   0.4298221  -0.208
## purposevacation      0.1535385   0.5752085   0.267
## purposewedding     -1.0387206   0.6924558  -1.500
## fico_range_high     -0.0067497   0.0038009  -1.776
## inq_last_6mths       0.2858123   0.0565287   5.056
## revol_util      -0.0011562   0.0037267  -0.310
## last_fico_range_high -0.0162113   0.0008758 -18.511
## desc_empty1      -0.1684563   0.1697134  -0.993
## dti              0.0054574   0.0114910   0.475
##
## Pr(>|z|)
## (Intercept)      2.08e-06 ***
## `term 60 months`    0.0222 *
## `verification_statusSource Verified` 0.0733 .
## verification_statusVerified      0.5530
## purposecredit_card      0.0873 .
## purposedebt_consolidation 0.2678
## purposeeducational      0.3393
## purposehome_improvement  0.2115
## purposehouse      0.6514
## purposemajor_purchase  0.0748 .
## purposemedical      0.7982
## purposemoving      0.0924 .
## purposeother      0.1310
## purposerenewable_energy 0.7059
## purposesmall_business  0.8352
## purposevacation      0.7895
## purposewedding      0.1336
## fico_range_high      0.0758 .
## inq_last_6mths      4.28e-07 ***
## revol_util      0.7564
## last_fico_range_high < 2e-16 ***
## desc_empty1      0.3209
## dti              0.6348
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##

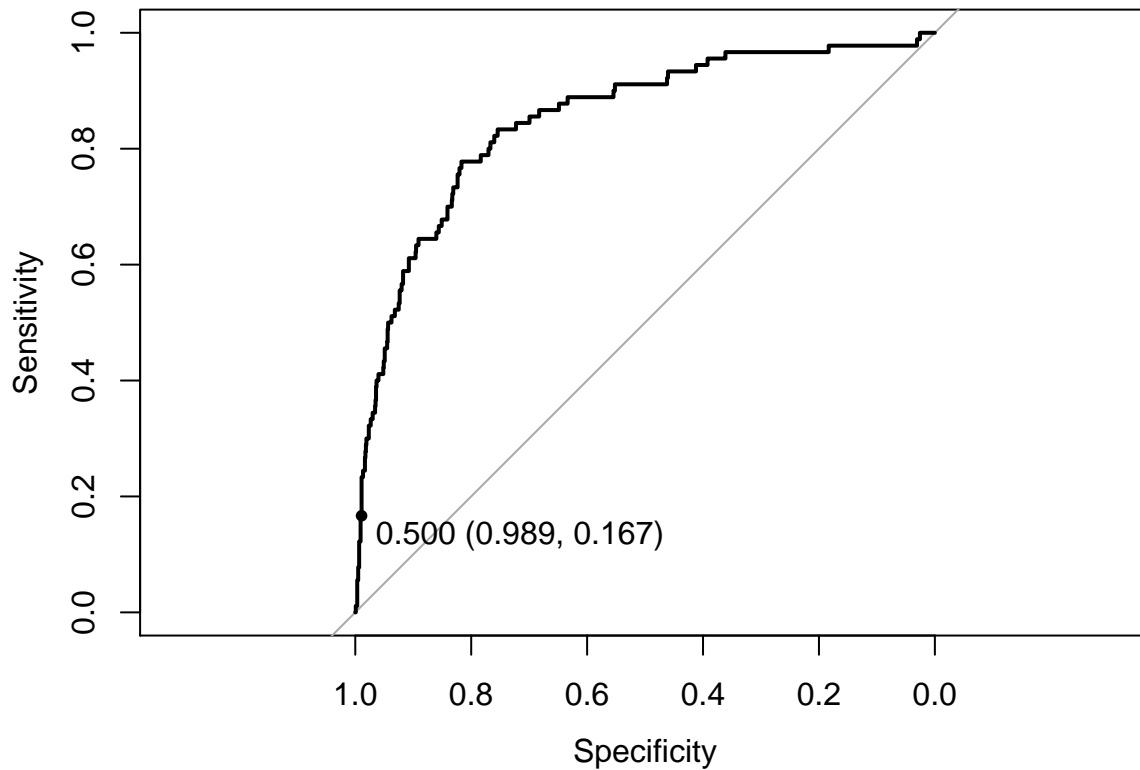
```

```

##      Null deviance: 1975.3  on 3878  degrees of freedom
## Residual deviance: 1429.0  on 3856  degrees of freedom
## AIC: 1475
##
## Number of Fisher Scoring iterations: 6
##
## glm variable importance
##
##      only 20 most important variables shown (out of 22)
##
##                                     Overall
## last_fico_range_high                100.0000
## inq_last_6mths                     26.4873
## `term 60 months`                   11.3623
## `verification_statusSource Verified` 8.6497
## purposemajor_purchase               8.5978
## fico_range_high                    8.5655
## purposecredit_card                 8.2038
## purposemoving                      8.0575
## purposeother                       7.1151
## purposewedding                     7.0589
## purposehome_improvement            5.6891
## purposedebt_consolidation          4.9177
## desc_empty1                        4.2864
## purposeeducational                 4.0840
## verification_statusVerified        2.1051
## dti                                1.4581
## purposehouse                       1.3318
## purposerenewable_energy            0.9249
## revol_util                         0.5584
## purposevacation                    0.3217
## Confusion Matrix and Statistics
##
##           Reference
## Prediction good  bad
##      good 1188   75
##      bad   13   15
##
##           Accuracy : 0.9318
##           95% CI : (0.9167, 0.945)
##      No Information Rate : 0.9303
##      P-Value [Acc > NIR] : 0.4409
##
##           Kappa : 0.2287
##      McNemar's Test P-Value : 7.893e-11
##
##           Sensitivity : 0.16667
##           Specificity : 0.98918
##      Pos Pred Value : 0.53571
##      Neg Pred Value : 0.94062
##           Prevalence : 0.06971
##      Detection Rate : 0.01162
##      Detection Prevalence : 0.02169
##      Balanced Accuracy : 0.57792

```

```
##
##      'Positive' Class : bad
##
```



```
##
## Call:
## roc.default(response = dft_test$status, predictor = testProbs[,      "bad"])
##
## Data: testProbs[, "bad"] in 1201 controls (dft_test$status good) < 90 cases (dft_test$status bad).
## Area under the curve: 0.8509
```

Random Forest Model

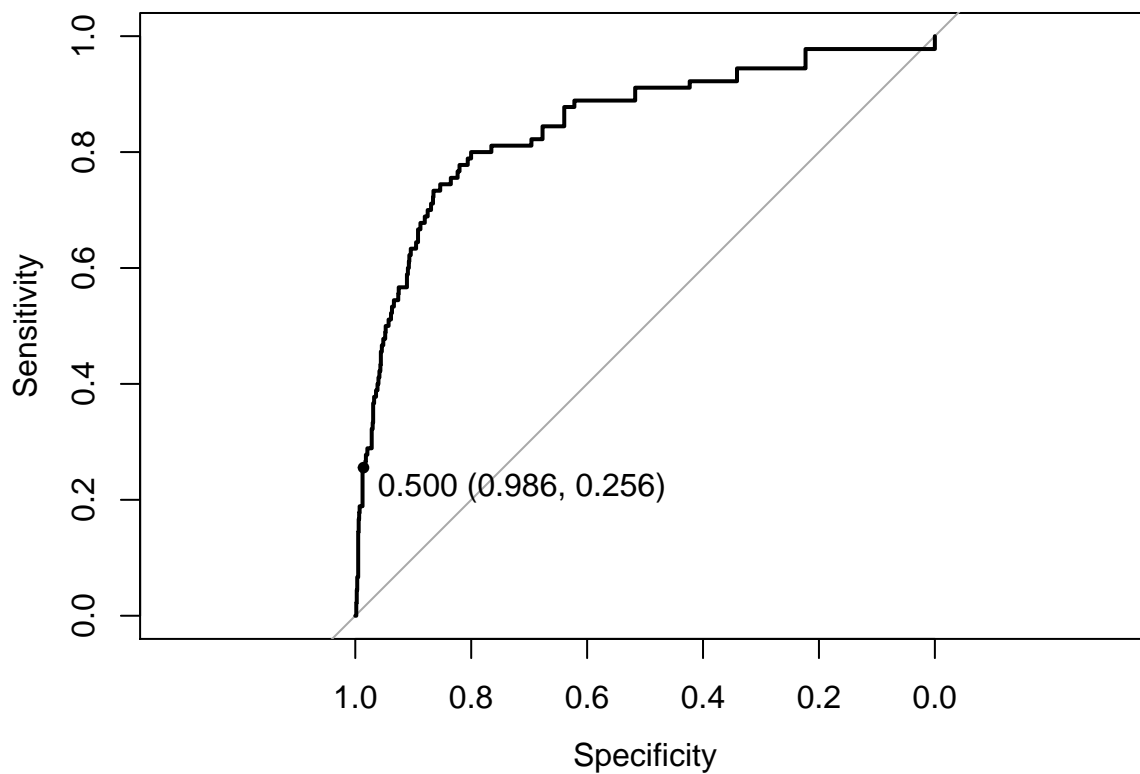
```
## Random Forest
##
## 3879 samples
##    9 predictor
##    2 classes: 'good', 'bad'
##
## No pre-processing
## Resampling: Bootstrapped (25 reps)
##
## Summary of sample sizes: 3879, 3879, 3879, 3879, 3879, 3879, ...
##
## Resampling results across tuning parameters:
##
##      mtry  Accuracy  Kappa      Accuracy SD  Kappa SD
##      2     0.9286334  0.002294058  0.005948366  0.006716373
```

```

## 12 0.9294704 0.320081829 0.005578916 0.036346900
## 22 0.9257695 0.315731616 0.005629033 0.025405788
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was mtry = 12.
##      Length Class      Mode
## call          4 -none-    call
## type          1 -none-    character
## predicted     3879 factor   numeric
## err.rate      1500 -none-    numeric
## confusion      6 -none-    numeric
## votes        7758 matrix   numeric
## oob.times     3879 -none-    numeric
## classes       2 -none-    character
## importance    22 -none-    numeric
## importanceSD   0 -none-    NULL
## localImportance 0 -none-    NULL
## proximity      0 -none-    NULL
## ntree         1 -none-    numeric
## mtry          1 -none-    numeric
## forest       14 -none-    list
## y            3879 factor   numeric
## test          0 -none-    NULL
## inbag          0 -none-    NULL
## xNames        22 -none-    character
## problemType    1 -none-    character
## tuneValue      1 data.frame list
## obsLevels      2 -none-    character
## rf variable importance
##
## only 20 most important variables shown (out of 22)
##
##      Overall
## last_fico_range_high 100.0000
## dti                  48.6456
## revol_util           45.6345
## fico_range_high      33.4729
## inq_last_6mths       19.1502
## purposedebt_consolidation 5.4778
## verification_statusVerified 4.9122
## purposeother         4.2978
## verification_statusSource Verified 3.9842
## desc_empty1          3.9830
## purposecredit_card   3.5645
## term 60 months       3.2203
## purposehome_improvement 2.5997
## purposemajor_purchase 2.3445
## purposesmall_business 1.5694
## purposemedical       1.4680
## purposeeducational   1.2634
## purposevacation      1.0744
## purposemoving        0.8869
## purposehouse         0.6622
## Confusion Matrix and Statistics

```

```
##
##           Reference
## Prediction good  bad
##           good 1185  67
##           bad   16  23
##
##           Accuracy : 0.9357
##           95% CI : (0.9209, 0.9485)
##           No Information Rate : 0.9303
##           P-Value [Acc > NIR] : 0.2412
##
##           Kappa : 0.3283
## Mcnemar's Test P-Value : 4.06e-08
##
##           Sensitivity : 0.25556
##           Specificity : 0.98668
##           Pos Pred Value : 0.58974
##           Neg Pred Value : 0.94649
##           Prevalence : 0.06971
##           Detection Rate : 0.01782
##           Detection Prevalence : 0.03021
##           Balanced Accuracy : 0.62112
##
##           'Positive' Class : bad
##
```



```
##
## Call:
## roc.default(response = dft_test$status, predictor = testProbs[,      "bad"])
##
```



```
##
## Data: testProbs[, "bad"] in 1201 controls (dft_test$status good) < 90 cases (dft_test$status bad).
## Area under the curve: 0.8529
```

Gradient Boost Model

```
## Loading required package: plyr
```

```
## Warning: package 'plyr' was built under R version 3.1.3
```

```
## -----
## You have loaded plyr after dplyr - this is likely to cause problems.
## If you need functions from both plyr and dplyr, please load plyr first, then dplyr:
## library(plyr); library(dplyr)
## -----
```

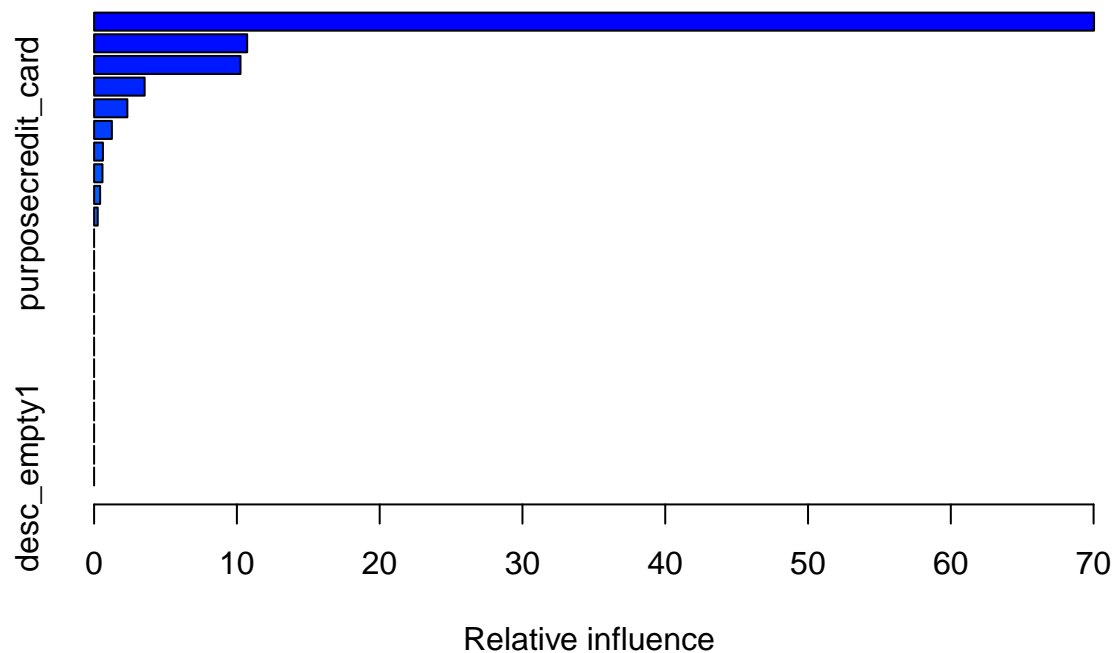
```
##
## Attaching package: 'plyr'
##
## The following objects are masked from 'package:reshape':
##
##   rename, round_any
##
## The following object is masked from 'package:lubridate':
##
##   here
##
## The following objects are masked from 'package:dplyr':
##
##   arrange, count, desc, failwith, id, mutate, rename, summarise,
##   summarize
```

```
## Stochastic Gradient Boosting
##
## 3879 samples
##   9 predictor
##   2 classes: 'good', 'bad'
##
## No pre-processing
## Resampling: Bootstrapped (25 reps)
##
## Summary of sample sizes: 3879, 3879, 3879, 3879, 3879, 3879, ...
##
```

```
## Resampling results across tuning parameters:
```

```
##
##   interaction.depth  n.trees  Accuracy  Kappa    Accuracy SD
##   1                 50      0.9315846  0.2746793  0.004460033
##   1                 100      0.9305221  0.2731812  0.004811056
##   1                 150      0.9302392  0.2781188  0.005544541
##   2                  50      0.9332636  0.3065577  0.004945600
##   2                 100      0.9317234  0.3080410  0.005155197
##   2                 150      0.9308576  0.3097073  0.005087139
##   3                  50      0.9325115  0.3126589  0.005051484
##   3                 100      0.9317023  0.3194381  0.004382830
```

```
##      3              150      0.9295464  0.3060641  0.004869173
##      Kappa SD
##      0.06294179
##      0.06290319
##      0.05216209
##      0.04617534
##      0.04941934
##      0.04880348
##      0.04759473
##      0.04865336
##      0.04405054
##
## Tuning parameter 'shrinkage' was held constant at a value of 0.1
##
## Tuning parameter 'n.minobsinnode' was held constant at a value of 10
## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were n.trees = 50, interaction.depth
## = 2, shrinkage = 0.1 and n.minobsinnode = 10.
```



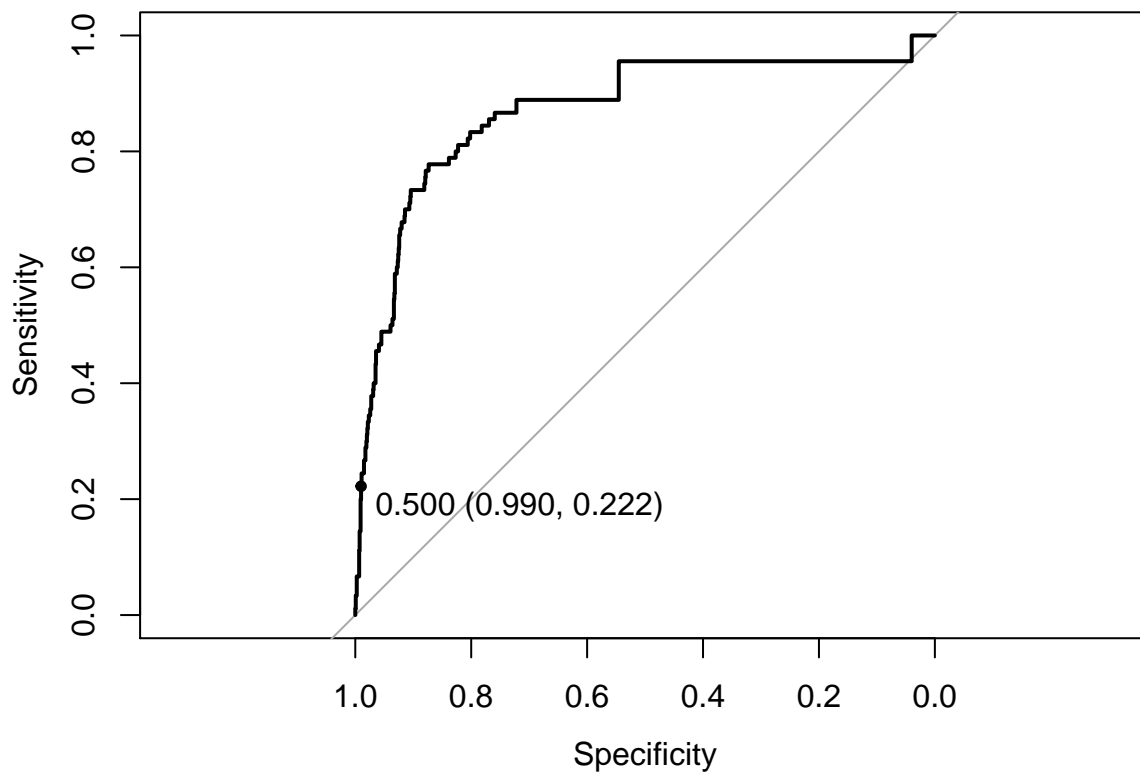
```
##
## last_fico_range_high      last_fico_range_high
## fico_range_high          fico_range_high
## inq_last_6mths           inq_last_6mths
## revol_util               revol_util
## dti                      dti
## purposeother              purposeother
## term 60 months            term 60 months
## purposecredit_card        purposecredit_card
## purposedebt_consolidation  purposedebt_consolidation
## purposeeducational        purposeeducational
## verification_statusSource Verified verification_statusSource Verified
## verification_statusVerified verification_statusVerified
```

## purposehome_improvement		purposehome_improvement
## purposehouse		purposehouse
## purposemajor_purchase		purposemajor_purchase
## purposemedical		purposemedical
## purposemoving		purposemoving
## purposerenewable_energy		purposerenewable_energy
## purposesmall_business		purposesmall_business
## purposevacation		purposevacation
## purposewedding		purposewedding
## desc_empty1		desc_empty1
##	rel.inf	
## last_fico_range_high	70.0311765	
## fico_range_high	10.7205994	
## inq_last_6mths	10.2527697	
## revol_util	3.5395880	
## dti	2.3295672	
## purposeother	1.2455554	
## term 60 months	0.6185364	
## purposecredit_card	0.5857127	
## purposedebt_consolidation	0.4216319	
## purposeeducational	0.2548628	
## verification_statusSource Verified	0.0000000	
## verification_statusVerified	0.0000000	
## purposehome_improvement	0.0000000	
## purposehouse	0.0000000	
## purposemajor_purchase	0.0000000	
## purposemedical	0.0000000	
## purposemoving	0.0000000	
## purposerenewable_energy	0.0000000	
## purposesmall_business	0.0000000	
## purposevacation	0.0000000	
## purposewedding	0.0000000	
## desc_empty1	0.0000000	
## gbm variable importance		
##		
##	only 20 most important variables shown (out of 22)	
##		
##	Overall	
## last_fico_range_high	100.0000	
## fico_range_high	15.3083	
## inq_last_6mths	14.6403	
## revol_util	5.0543	
## dti	3.3265	
## purposeother	1.7786	
## term 60 months	0.8832	
## purposecredit_card	0.8364	
## purposedebt_consolidation	0.6021	
## purposeeducational	0.3639	
## purposevacation	0.0000	
## desc_empty1	0.0000	
## purposemoving	0.0000	
## purposewedding	0.0000	
## purposehome_improvement	0.0000	
## purposemedical	0.0000	

```

## verification_statusVerified      0.0000
## purposesmall_business            0.0000
## purposehouse                      0.0000
## verification_statusSource Verified 0.0000
## Confusion Matrix and Statistics
##
##           Reference
## Prediction good  bad
##      good 1189   70
##      bad   12   20
##
##           Accuracy : 0.9365
##           95% CI : (0.9218, 0.9492)
##      No Information Rate : 0.9303
##      P-Value [Acc > NIR] : 0.2076
##
##           Kappa : 0.3024
##  McNemar's Test P-Value : 3.082e-10
##
##           Sensitivity : 0.22222
##           Specificity : 0.99001
##      Pos Pred Value : 0.62500
##      Neg Pred Value : 0.94440
##           Prevalence : 0.06971
##      Detection Rate : 0.01549
##      Detection Prevalence : 0.02479
##      Balanced Accuracy : 0.60612
##
##      'Positive' Class : bad
##

```



```
##
## Call:
## roc.default(response = dft_test$status, predictor = testProbs[,      "bad"])
##
## Data: testProbs[, "bad"] in 1201 controls (dft_test$status good) < 90 cases (dft_test$status bad).
## Area under the curve: 0.888
```

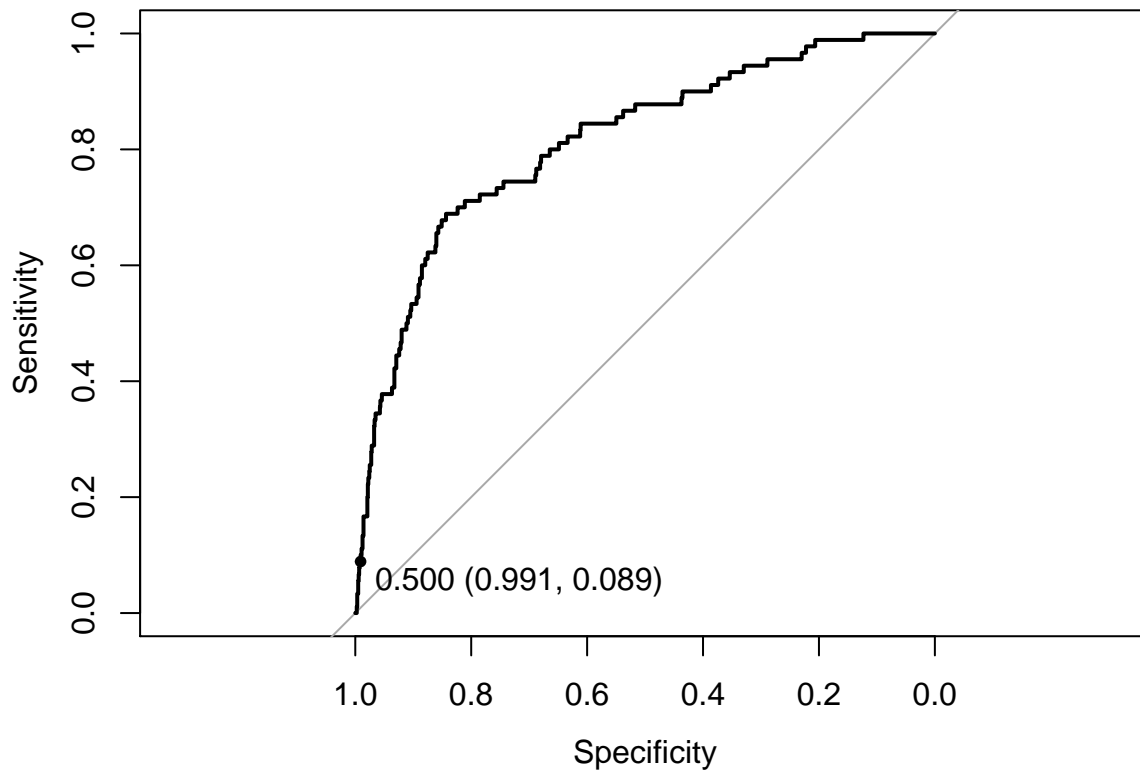
SVM Model

```
## Support Vector Machines with Radial Basis Function Kernel
##
## 3879 samples
##    9 predictor
##    2 classes: 'good', 'bad'
##
## Pre-processing: centered, scaled
## Resampling: Cross-Validated (10 fold)
##
## Summary of sample sizes: 3491, 3491, 3491, 3491, 3492, 3491, ...
##
## Resampling results across tuning parameters:
##
##    C      Accuracy  Kappa      Accuracy SD  Kappa SD
##    0.25  0.9283411  0.1374781  0.008142182  0.10385356
##    0.50  0.9283398  0.1378456  0.007593441  0.09922288
##    1.00  0.9293707  0.1402285  0.007239328  0.10303976
##
## Tuning parameter 'sigma' was held constant at a value of 0.04150303
```

```

## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were sigma = 0.04150303 and C = 1.
## Length Class Mode
##      1      ksvm      S4
## ROC curve variable importance
##
##                               Importance
## last_fico_range_high      100.000
## fico_range_high           38.782
## inq_last_6mths            28.886
## revol_util                18.958
## dti                       13.235
## purpose                   11.062
## term                      6.388
## verification_status       2.016
## desc_empty                 0.000
## Confusion Matrix and Statistics
##
##              Reference
## Prediction good  bad
##      good 1190   82
##      bad   11    8
##
##              Accuracy : 0.928
##              95% CI : (0.9125, 0.9415)
##      No Information Rate : 0.9303
##      P-Value [Acc > NIR] : 0.6539
##
##              Kappa : 0.1255
##      McNemar's Test P-Value : 3.909e-13
##
##              Sensitivity : 0.088889
##              Specificity : 0.990841
##              Pos Pred Value : 0.421053
##              Neg Pred Value : 0.935535
##              Prevalence : 0.069713
##              Detection Rate : 0.006197
##      Detection Prevalence : 0.014717
##              Balanced Accuracy : 0.539865
##
##      'Positive' Class : bad
##

```



```
##
## Call:
## roc.default(response = dft_test$status, predictor = testProbs[,      "bad"])
##
## Data: testProbs[, "bad"] in 1201 controls (dft_test$status good) < 90 cases (dft_test$status bad).
## Area under the curve: 0.8141
```

Results for Grade B Loans

Approximately 16% of the Grade B loans in this dataset went bad. With the four models, we were able to predict between 40% and 52% of the bad loans. This predictive ability is based on a 50% probability classification cutoff. As the ROC curves show, it's possible to predict the bad loans with a higher probability, of course, with a higher false positive rate, though. The FICO range and the number of inquiries in the past 6 months were also important predictors for this loan grade.

Logistic Regression Model

```
## Generalized Linear Model
##
## 4929 samples
##    9 predictor
##    2 classes: 'good', 'bad'
##
## No pre-processing
## Resampling: Bootstrapped (25 reps)
##
## Summary of sample sizes: 4929, 4929, 4929, 4929, 4929, 4929, ...
```

```

##
## Resampling results
##
##   Accuracy   Kappa   Accuracy SD   Kappa SD
##   0.8790792  0.4665746  0.005767996  0.02478362
##
##
## Call:
## NULL
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -4.6278  -0.4408  -0.2694  -0.1564   3.1895
##
## Coefficients:
##                                     Estimate Std. Error z value
## (Intercept)                       7.1089236  1.7249985   4.121
## `term 60 months`                   0.6568411  0.1369624   4.796
## `verification_statusSource Verified` -0.1767463  0.1353809  -1.306
## verification_statusVerified        -0.0766389  0.1173121  -0.653
## purposecredit_card                  0.3174541  0.2942137   1.079
## purposedebt_consolidation           0.1159055  0.2697443   0.430
## purposeeducational                  0.4195458  0.4299855   0.976
## purposehome_improvement            -0.0639845  0.3134996  -0.204
## purposehouse                       -0.2207673  0.5420417  -0.407
## purposemajor_purchase               0.0280041  0.3352412   0.084
## purposemedical                     0.5750696  0.4044557   1.422
## purposemoving                      0.4152181  0.4316807   0.962
## purposeother                       0.3365332  0.2878080   1.169
## purposerenewable_energy             0.8386232  0.7044937   1.190
## purposesmall_business              0.5745717  0.3295441   1.744
## purposevacation                    -0.0621005  0.5337315  -0.116
## purposewedding                     -0.1365327  0.4914831  -0.278
## fico_range_high                    0.0023927  0.0023237   1.030
## inq_last_6mths                     0.6073043  0.0346215  17.541
## revol_util                         0.0039184  0.0021216   1.847
## last_fico_range_high               -0.0176608  0.0006585 -26.818
## desc_empty1                       -0.1391766  0.1181128  -1.178
## dti                               0.0001462  0.0077470   0.019
##                                     Pr(>|z|)
## (Intercept)                       3.77e-05 ***
## `term 60 months`                   1.62e-06 ***
## `verification_statusSource Verified` 0.1917
## verification_statusVerified        0.5136
## purposecredit_card                  0.2806
## purposedebt_consolidation           0.6674
## purposeeducational                  0.3292
## purposehome_improvement            0.8383
## purposehouse                       0.6838
## purposemajor_purchase               0.9334
## purposemedical                     0.1551
## purposemoving                      0.3361
## purposeother                       0.2423

```

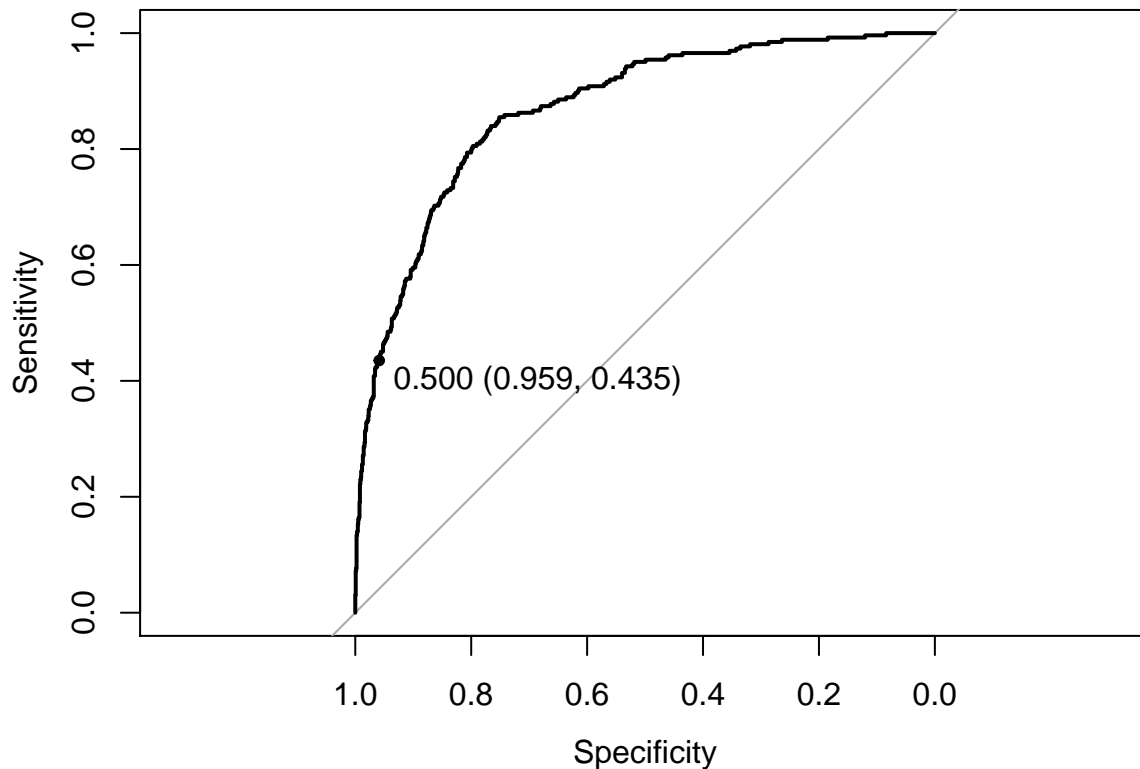


```

## purposerenewable_energy          0.2339
## purposesmall_business            0.0812 .
## purposevacation                   0.9074
## purposewedding                    0.7812
## fico_range_high                   0.3032
## inq_last_6mths                    < 2e-16 ***
## revol_util                        0.0648 .
## last_fico_range_high              < 2e-16 ***
## desc_empty1                       0.2387
## dti                              0.9849
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 4325.5  on 4928  degrees of freedom
## Residual deviance: 2900.0  on 4906  degrees of freedom
## AIC: 2946
##
## Number of Fisher Scoring iterations: 6
##
## glm variable importance
##
##    only 20 most important variables shown (out of 22)
##
##
##                                Overall
## last_fico_range_high          100.0000
## inq_last_6mths                 65.3832
## `term 60 months`              17.8246
## revol_util                     6.8210
## purposesmall_business         6.4355
## purposemedical                 5.2351
## `verification_statusSource Verified` 4.8011
## purposerenewable_energy       4.3714
## desc_empty1                   4.3265
## purposeother                  4.2927
## purposecredit_card            3.9558
## fico_range_high               3.7717
## purposeeducational            3.5704
## purposemoving                 3.5187
## verification_statusVerified   2.3673
## purposedebt_consolidation     1.5329
## purposehouse                  1.4494
## purposewedding                0.9662
## purposehome_improvement       0.6912
## purposevacation               0.3638
## Confusion Matrix and Statistics
##
##           Reference
## Prediction good  bad
##      good 1323  148
##      bad   57   114
##
##           Accuracy : 0.8752

```

```
##          95% CI : (0.8582, 0.8908)
##    No Information Rate : 0.8404
##    P-Value [Acc > NIR] : 4.379e-05
##
##          Kappa : 0.4583
## Mcnemar's Test P-Value : 3.260e-10
##
##      Sensitivity : 0.43511
##      Specificity : 0.95870
##      Pos Pred Value : 0.66667
##      Neg Pred Value : 0.89939
##      Prevalence : 0.15956
##      Detection Rate : 0.06943
##      Detection Prevalence : 0.10414
##      Balanced Accuracy : 0.69691
##
##      'Positive' Class : bad
##
```



```
##
## Call:
## roc.default(response = dft_test$status, predictor = testProbs[,      "bad"])
##
## Data: testProbs[, "bad"] in 1380 controls (dft_test$status good) < 262 cases (dft_test$status bad).
## Area under the curve: 0.8701
```

Random Forest Model

```
## Random Forest
```

```

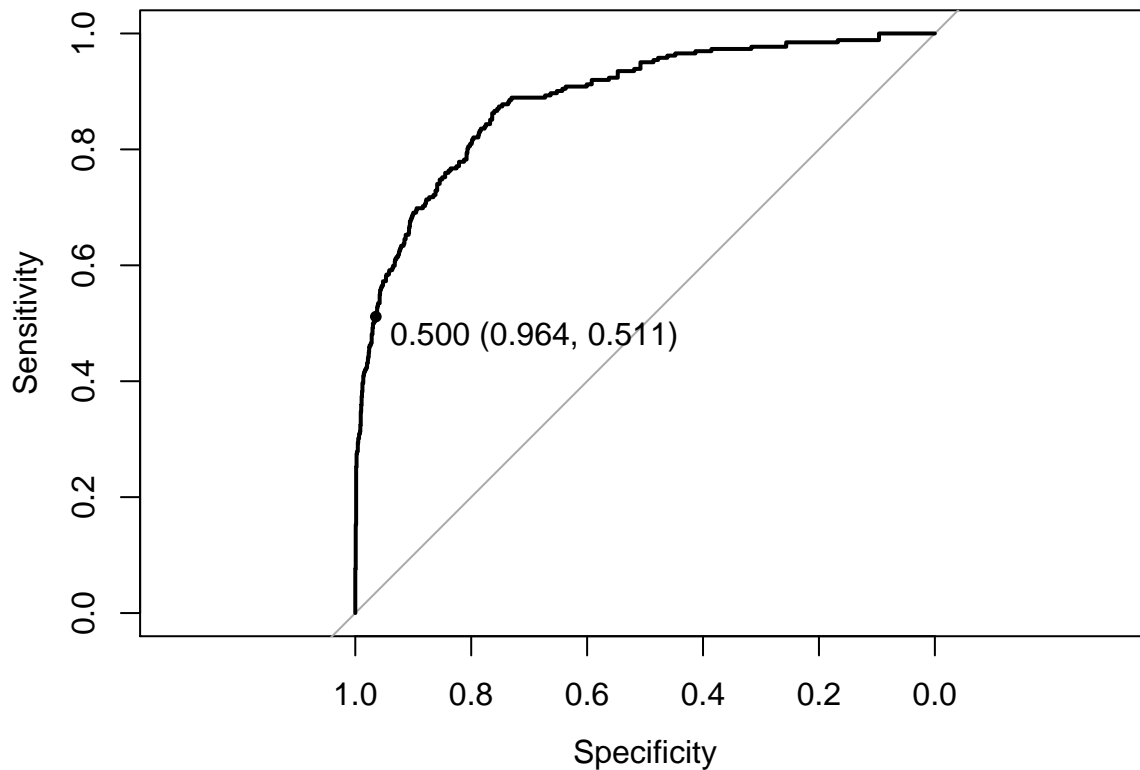
##
## 4929 samples
##    9 predictor
##    2 classes: 'good', 'bad'
##
## No pre-processing
## Resampling: Bootstrapped (25 reps)
##
## Summary of sample sizes: 4929, 4929, 4929, 4929, 4929, 4929, ...
##
## Resampling results across tuning parameters:
##
##  mtry  Accuracy   Kappa     Accuracy SD  Kappa SD
##    2    0.8702462 0.2897623  0.008941716  0.04889492
##   12    0.8964821 0.5694612  0.006941764  0.02595673
##   22    0.8923925 0.5568188  0.007532415  0.02818076
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was mtry = 12.
##
##      Length Class      Mode
## call           4  -none-    call
## type           1  -none-   character
## predicted      4929 factor    numeric
## err.rate       1500 -none-    numeric
## confusion        6  -none-    numeric
## votes          9858 matrix    numeric
## oob.times       4929 -none-    numeric
## classes         2  -none-   character
## importance      22  -none-    numeric
## importanceSD      0  -none-     NULL
## localImportance  0  -none-     NULL
## proximity        0  -none-     NULL
## ntree           1  -none-    numeric
## mtry            1  -none-    numeric
## forest          14  -none-     list
## y               4929 factor    numeric
## test            0  -none-     NULL
## inbag            0  -none-     NULL
## xNames          22  -none-   character
## problemType      1  -none-   character
## tuneValue        1 data.frame list
## obsLevels        2  -none-   character
## rf variable importance
##
##    only 20 most important variables shown (out of 22)
##
##
##      Overall
## last_fico_range_high 100.0000
## inq_last_6mths       47.8090
## revol_util           36.6473
## dti                   36.2931
## fico_range_high      35.1854
## term 60 months       4.4510
## desc_empty1           3.8201

```

```

## verification_statusVerified      3.4280
## purposedebt_consolidation        3.3190
## verification_statusSource Verified 2.9665
## purposeother                     2.5541
## purposesmall_business            2.1953
## purposehome_improvement          1.8584
## purposecredit_card               1.7846
## purposemajor_purchase             1.5046
## purposemedical                   0.8645
## purposeeducational               0.6455
## purposemoving                    0.5585
## purposehouse                     0.5020
## purposerenewable_energy          0.2257
## Confusion Matrix and Statistics
##
##           Reference
## Prediction good  bad
##      good 1331  128
##      bad   49  134
##
##           Accuracy : 0.8922
##           95% CI : (0.8762, 0.9068)
##      No Information Rate : 0.8404
##      P-Value [Acc > NIR] : 1.125e-09
##
##           Kappa : 0.5422
##  McNemar's Test P-Value : 4.550e-09
##
##           Sensitivity : 0.51145
##           Specificity : 0.96449
##           Pos Pred Value : 0.73224
##           Neg Pred Value : 0.91227
##           Prevalence : 0.15956
##           Detection Rate : 0.08161
##      Detection Prevalence : 0.11145
##           Balanced Accuracy : 0.73797
##
##           'Positive' Class : bad
##

```

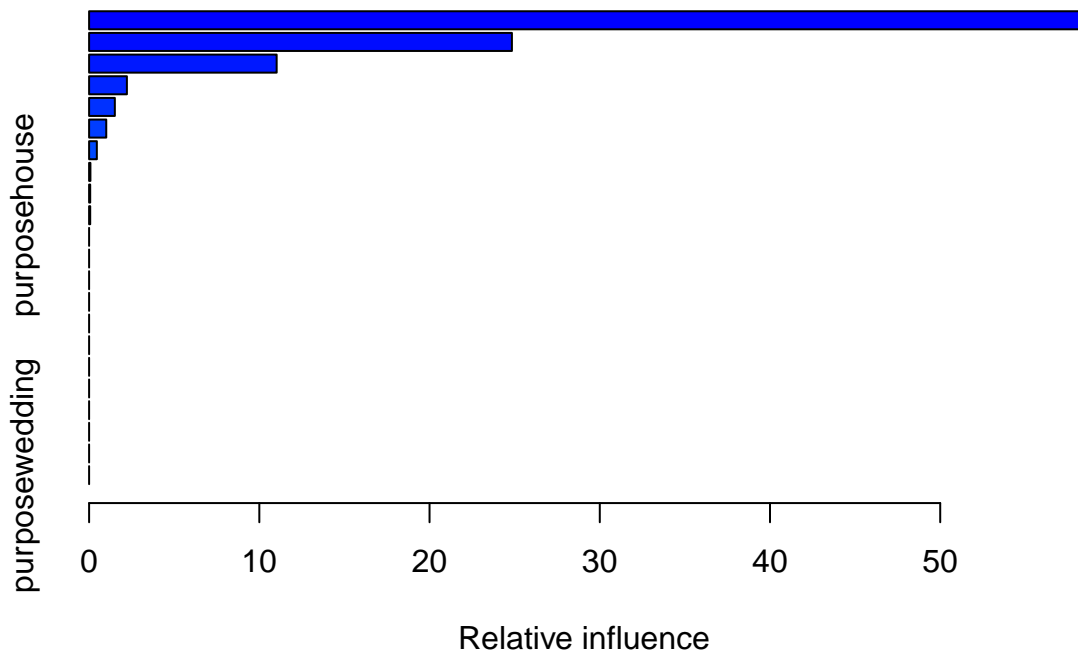


```
##
## Call:
## roc.default(response = dft_test$status, predictor = testProbs[,      "bad"])
##
## Data: testProbs[, "bad"] in 1380 controls (dft_test$status good) < 262 cases (dft_test$status bad).
## Area under the curve: 0.8885
```

Gradient Boost Model

```
## Stochastic Gradient Boosting
##
## 4929 samples
##    9 predictor
##    2 classes: 'good', 'bad'
##
## No pre-processing
## Resampling: Bootstrapped (25 reps)
##
## Summary of sample sizes: 4929, 4929, 4929, 4929, 4929, 4929, ...
##
## Resampling results across tuning parameters:
##
##  interaction.depth  n.trees  Accuracy  Kappa    Accuracy SD
##    1                50      0.8844194  0.4871892  0.004767854
##    1                100      0.8885364  0.5228538  0.004690834
##    1                150      0.8902577  0.5355176  0.004658145
##    2                 50      0.9013938  0.5826713  0.004021311
##    2                100      0.9003550  0.5804478  0.003644686
```

```
##      2          150      0.8997678  0.5780676  0.003672264
##      3           50      0.9015587  0.5839425  0.005022749
##      3          100      0.9001326  0.5787466  0.004784543
##      3          150      0.8994194  0.5770137  0.004370545
## Kappa SD
## 0.02906046
## 0.02712741
## 0.02003669
## 0.01583832
## 0.01826659
## 0.01778459
## 0.02052217
## 0.02157765
## 0.02015747
##
## Tuning parameter 'shrinkage' was held constant at a value of 0.1
##
## Tuning parameter 'n.minobsinnode' was held constant at a value of 10
## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were n.trees = 50, interaction.depth
## = 3, shrinkage = 0.1 and n.minobsinnode = 10.
```



```
##                                     var
## last_fico_range_high               last_fico_range_high
## inq_last_6mths                     inq_last_6mths
## fico_range_high                    fico_range_high
## term 60 months                     term 60 months
## revol_util                         revol_util
## dti                               dti
## purposesmall_business              purposesmall_business
## desc_empty1                       desc_empty1
## purposemoving                      purposemoving
```

```

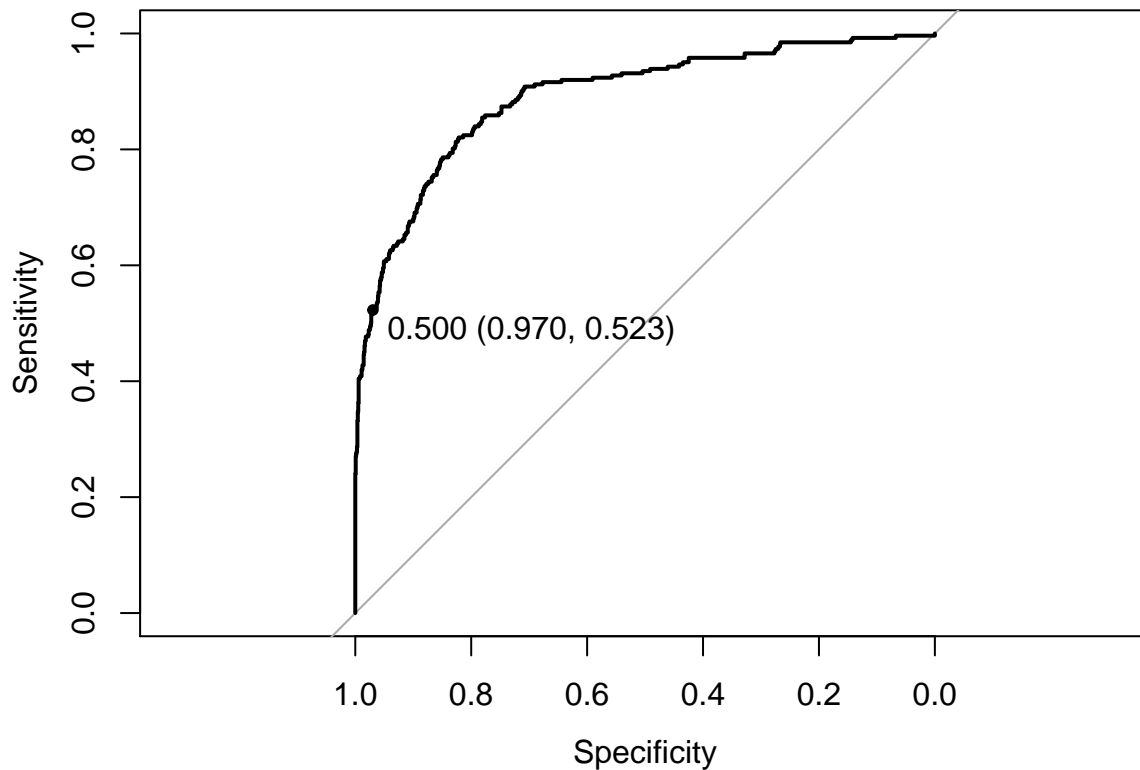
## purposehouse                                purposehouse
## verification_statusSource Verified verification_statusSource Verified
## verification_statusVerified                verification_statusVerified
## purposecredit_card                        purposecredit_card
## purposedebt_consolidation                purposedebt_consolidation
## purposeeducational                      purposeeducational
## purposehome_improvement                  purposehome_improvement
## purposemajor_purchase                    purposemajor_purchase
## purposemedical                          purposemedical
## purposeother                            purposeother
## purposerenewable_energy                  purposerenewable_energy
## purposevacation                          purposevacation
## purposewedding                          purposewedding
##
## rel.inf
## last_fico_range_high                    58.74332638
## inq_last_6mths                        24.83920077
## fico_range_high                      11.01830964
## term 60 months                      2.22012115
## revol_util                          1.51221465
## dti                                  1.00733671
## purposesmall_business                0.45485157
## desc_empty1                          0.07678682
## purposemoving                        0.06399285
## purposehouse                        0.06385946
## verification_statusSource Verified  0.00000000
## verification_statusVerified          0.00000000
## purposecredit_card                  0.00000000
## purposedebt_consolidation            0.00000000
## purposeeducational                  0.00000000
## purposehome_improvement              0.00000000
## purposemajor_purchase                0.00000000
## purposemedical                      0.00000000
## purposeother                        0.00000000
## purposerenewable_energy              0.00000000
## purposevacation                      0.00000000
## purposewedding                      0.00000000
## gbm variable importance
##
## only 20 most important variables shown (out of 22)
##
## Overall
## last_fico_range_high                100.0000
## inq_last_6mths                      42.2843
## fico_range_high                    18.7567
## term 60 months                      3.7794
## revol_util                          2.5743
## dti                                  1.7148
## purposesmall_business                0.7743
## desc_empty1                          0.1307
## purposemoving                        0.1089
## purposehouse                        0.1087
## purposemajor_purchase                0.0000
## purposemedical                      0.0000
## purposeeducational                  0.0000

```

```

## purposedebt_consolidation      0.0000
## purposehome_improvement        0.0000
## purposeother                    0.0000
## purposerenewable_energy        0.0000
## purposevacation                 0.0000
## verification_statusSource Verified 0.0000
## verification_statusVerified      0.0000
## Confusion Matrix and Statistics
##
##           Reference
## Prediction good  bad
##      good 1338  125
##      bad   42  137
##
##           Accuracy : 0.8983
##           95% CI : (0.8827, 0.9125)
##      No Information Rate : 0.8404
##      P-Value [Acc > NIR] : 8.247e-12
##
##           Kappa : 0.565
##  McNemar's Test P-Value : 2.219e-10
##
##           Sensitivity : 0.52290
##           Specificity : 0.96957
##      Pos Pred Value : 0.76536
##      Neg Pred Value : 0.91456
##           Prevalence : 0.15956
##      Detection Rate : 0.08343
##      Detection Prevalence : 0.10901
##      Balanced Accuracy : 0.74623
##
##           'Positive' Class : bad
##

```

```
##
## Call:
## roc.default(response = dft_test$status, predictor = testProbs[,      "bad"])
##
## Data: testProbs[, "bad"] in 1380 controls (dft_test$status good) < 262 cases (dft_test$status bad).
## Area under the curve: 0.8917
```

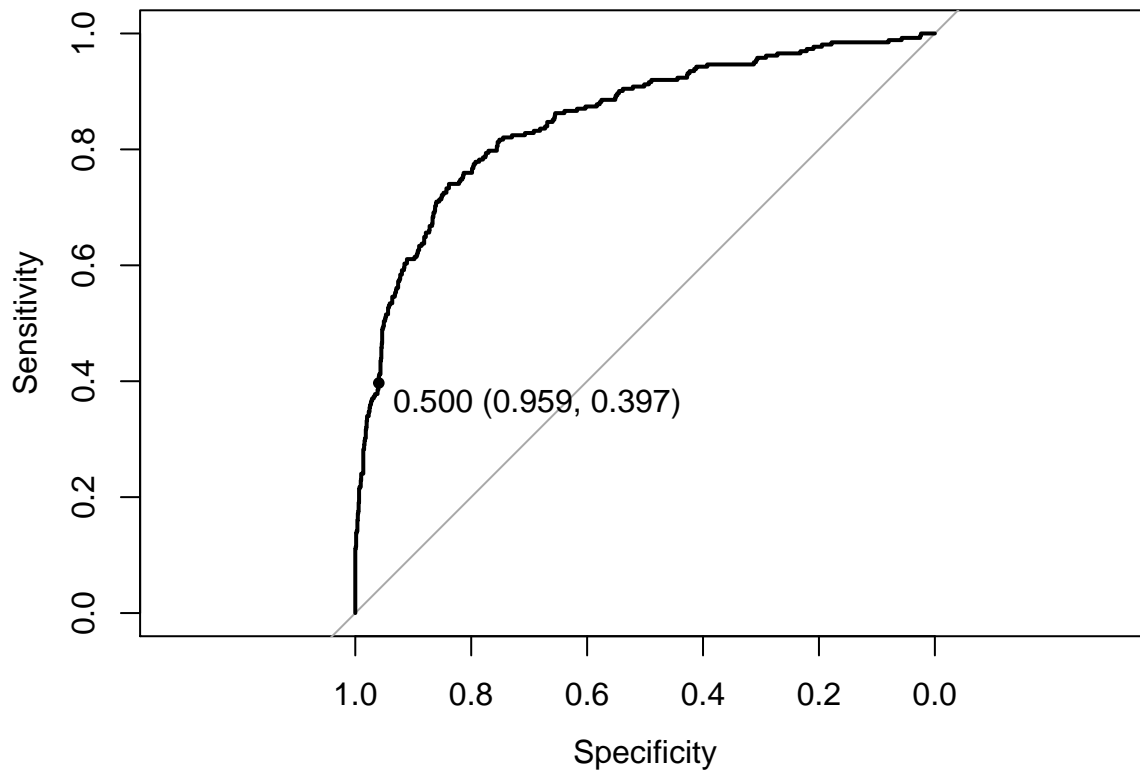
SVM Model

```
## Support Vector Machines with Radial Basis Function Kernel
##
## 4929 samples
##    9 predictor
##    2 classes: 'good', 'bad'
##
## Pre-processing: centered, scaled
## Resampling: Cross-Validated (10 fold)
##
## Summary of sample sizes: 4437, 4436, 4437, 4435, 4436, 4435, ...
##
## Resampling results across tuning parameters:
##
##    C      Accuracy  Kappa      Accuracy SD  Kappa SD
##    0.25  0.8792842  0.4893452  0.01214939  0.06537045
##    0.50  0.8800939  0.4703411  0.01349183  0.07394253
##    1.00  0.8841528  0.4883953  0.01042157  0.05627201
##
## Tuning parameter 'sigma' was held constant at a value of 0.04361204
```

```

## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were sigma = 0.04361204 and C = 1.
## Length Class Mode
##      1      ksvm      S4
## ROC curve variable importance
##
##                               Importance
## last_fico_range_high      100.0000
## inq_last_6mths            49.1108
## purpose                    21.3103
## term                       11.7052
## revol_util                 5.2943
## dti                        4.5958
## verification_status        1.7731
## fico_range_high            0.4807
## desc_empty                  0.0000
## Confusion Matrix and Statistics
##
##              Reference
## Prediction good  bad
##      good 1324  158
##      bad   56  104
##
##              Accuracy : 0.8697
##              95% CI : (0.8524, 0.8856)
##      No Information Rate : 0.8404
##      P-Value [Acc > NIR] : 0.0005217
##
##              Kappa : 0.4231
##      McNemar's Test P-Value : 5.048e-12
##
##              Sensitivity : 0.39695
##              Specificity : 0.95942
##              Pos Pred Value : 0.65000
##              Neg Pred Value : 0.89339
##              Prevalence : 0.15956
##              Detection Rate : 0.06334
##      Detection Prevalence : 0.09744
##              Balanced Accuracy : 0.67818
##
##      'Positive' Class : bad
##

```



```
##
## Call:
## roc.default(response = dft_test$status, predictor = testProbs[,      "bad"])
##
## Data: testProbs[, "bad"] in 1380 controls (dft_test$status good) < 262 cases (dft_test$status bad).
## Area under the curve: 0.8517
```

Results for Grade C Loans

Approximately 25% of the Grade C loans in this dataset went bad. With the four models, we were able to correctly predict between 58% and 66% of the bad loans. This predictive ability is based on a 50% probability classification cutoff. As the ROC curves show, it's possible to predict the bad loans with a higher probability, of course, with a higher false positive rate, though. The FICO range and the number of inquiries in the past 6 months were also important predictors for this loan grade.

Logistic Regression Model

```
## Generalized Linear Model
##
## 3919 samples
##    9 predictor
##    2 classes: 'good', 'bad'
##
## No pre-processing
## Resampling: Bootstrapped (25 reps)
##
## Summary of sample sizes: 3919, 3919, 3919, 3919, 3919, 3919, ...
```

```

##
## Resampling results
##
##   Accuracy   Kappa   Accuracy SD   Kappa SD
##   0.8344897  0.526276  0.008668034  0.02280993
##
##
## Call:
## NULL
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -4.1503  -0.5465  -0.3162   0.0414   3.0392
##
## Coefficients:
##                                Estimate Std. Error z value
## (Intercept)                   2.6795739  1.7759518   1.509
## `term 60 months`              0.4120035  0.1430904   2.879
## `verification_statusSource Verified` -0.6608699  0.1401579  -4.715
## verification_statusVerified      -0.2445471  0.1142248  -2.141
## purposecredit_card              0.1001298  0.3285056   0.305
## purposedebt_consolidation        -0.0156633  0.3063685  -0.051
## purposeeducational             -0.3521902  0.4358913  -0.808
## purposehome_improvement         -0.0570577  0.3461229  -0.165
## purposehouse                   0.0780706  0.5634877   0.139
## purposemajor_purchase           0.0413611  0.3673723   0.113
## purposemedical                 -0.1473572  0.4633574  -0.318
## purposemoving                  -0.4542210  0.4821819  -0.942
## purposeother                   -0.0002234  0.3234839  -0.001
## purposerenewable_energy         0.7425936  1.0916630   0.680
## purposesmall_business          0.3258348  0.3675664   0.886
## purposevacation                0.2648882  0.6105528   0.434
## purposewedding                 -0.0390929  0.4377212  -0.089
## fico_range_high                0.0080235  0.0025187   3.186
## inq_last_6mths                 0.7164624  0.0329162  21.766
## revol_util                     -0.0026653  0.0019165  -1.391
## last_fico_range_high           -0.0164993  0.0006725 -24.533
## desc_empty1                    -0.1117800  0.1212968  -0.922
## dti                            0.0316298  0.0078020   4.054
##                                Pr(>|z|)
## (Intercept)                   0.13135
## `term 60 months`              0.00399 **
## `verification_statusSource Verified` 2.41e-06 ***
## verification_statusVerified      0.03228 *
## purposecredit_card              0.76052
## purposedebt_consolidation        0.95923
## purposeeducational             0.41910
## purposehome_improvement         0.86906
## purposehouse                   0.88981
## purposemajor_purchase           0.91036
## purposemedical                 0.75047
## purposemoving                  0.34619
## purposeother                   0.99945

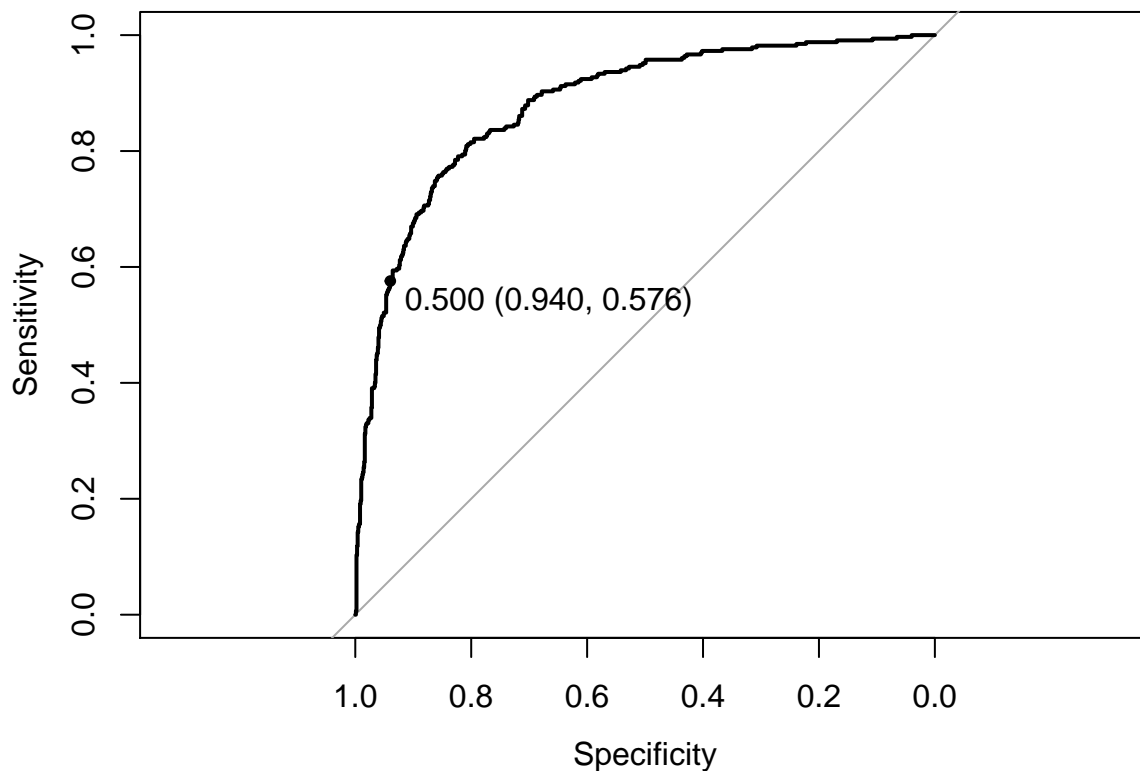
```

```

## purposerenewable_energy          0.49635
## purposesmall_business            0.37537
## purposevacation                  0.66440
## purposewedding                   0.92884
## fico_range_high                   0.00144 **
## inq_last_6mths                   < 2e-16 ***
## revol_util                       0.16431
## last_fico_range_high              < 2e-16 ***
## desc_empty1                      0.35677
## dti                              5.03e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 4430.0  on 3918  degrees of freedom
## Residual deviance: 2882.3  on 3896  degrees of freedom
## AIC: 2928.3
##
## Number of Fisher Scoring iterations: 5
##
## glm variable importance
##
##    only 20 most important variables shown (out of 22)
##
##
##                                Overall
## last_fico_range_high          100.0000
## inq_last_6mths                88.7219
## `verification_statusSource Verified` 19.2174
## dti                          16.5225
## fico_range_high              12.9824
## `term 60 months`            11.7340
## verification_statusVerified   8.7241
## revol_util                   5.6661
## purposemoving                 3.8371
## desc_empty1                  3.7536
## purposesmall_business        3.6106
## purposeeducational           3.2907
## purposerenewable_energy      2.7700
## purposevacation              1.7657
## purposemedical               1.2935
## purposecredit_card           1.2396
## purposehome_improvement      0.6691
## purposehouse                 0.5619
## purposemajor_purchase        0.4561
## purposewedding               0.3612
## Confusion Matrix and Statistics
##
##           Reference
## Prediction good bad
##      good  917 140
##      bad   59 190
##
##           Accuracy : 0.8476

```

```
##          95% CI : (0.827, 0.8667)
##    No Information Rate : 0.7473
##    P-Value [Acc > NIR] : < 2.2e-16
##
##          Kappa : 0.5609
## Mcnemar's Test P-Value : 1.419e-08
##
##      Sensitivity : 0.5758
##      Specificity : 0.9395
##      Pos Pred Value : 0.7631
##      Neg Pred Value : 0.8675
##      Prevalence : 0.2527
##      Detection Rate : 0.1455
##      Detection Prevalence : 0.1907
##      Balanced Accuracy : 0.7577
##
##      'Positive' Class : bad
##
```



```
##
## Call:
## roc.default(response = dft_test$status, predictor = testProbs[,      "bad"])
##
## Data: testProbs[, "bad"] in 976 controls (dft_test$status good) < 330 cases (dft_test$status bad).
## Area under the curve: 0.8813
```

Random Forest Model

```
## Random Forest
```

```

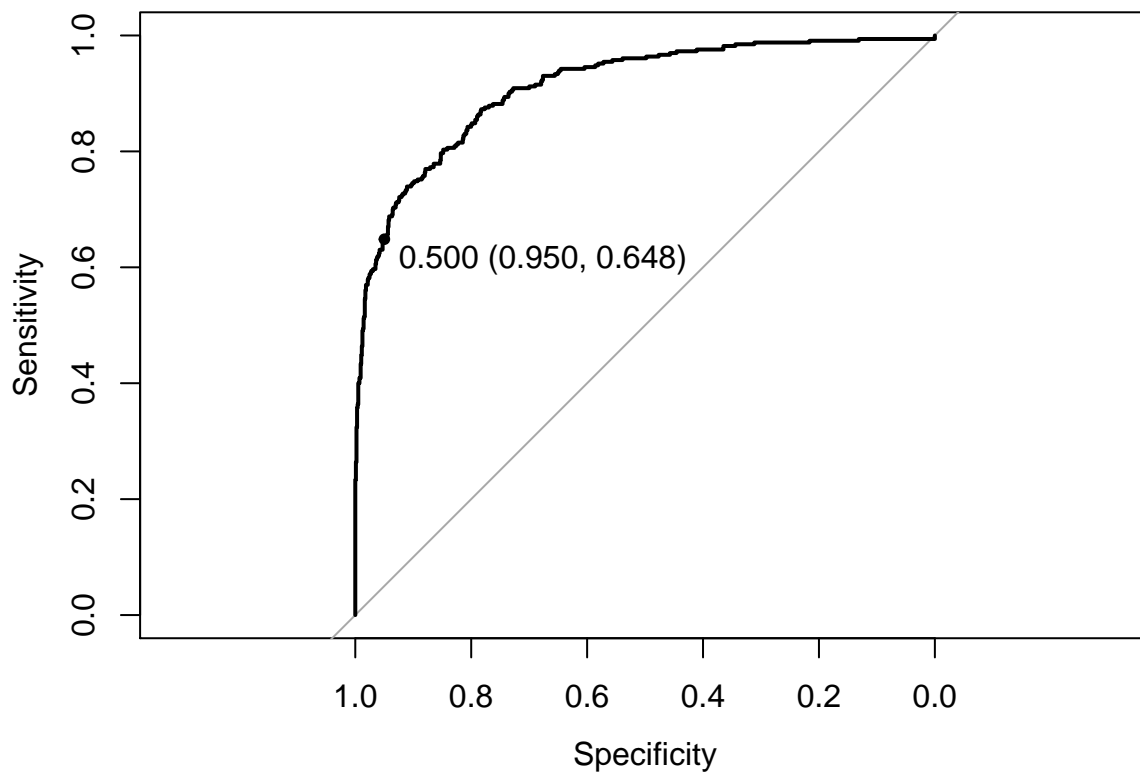
##
## 3919 samples
##    9 predictor
##    2 classes: 'good', 'bad'
##
## No pre-processing
## Resampling: Bootstrapped (25 reps)
##
## Summary of sample sizes: 3919, 3919, 3919, 3919, 3919, 3919, ...
##
## Resampling results across tuning parameters:
##
##  mtry  Accuracy   Kappa     Accuracy SD  Kappa SD
##    2    0.8536866  0.5282783  0.009429435  0.03045501
##   12    0.8743988  0.6469866  0.008030578  0.01851626
##   22    0.8705122  0.6378193  0.009296249  0.02179501
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was mtry = 12.
##
##      Length Class      Mode
## call           4  -none-    call
## type           1  -none-    character
## predicted      3919 factor    numeric
## err.rate       1500 -none-    numeric
## confusion       6  -none-    numeric
## votes          7838 matrix    numeric
## oob.times       3919 -none-    numeric
## classes         2  -none-    character
## importance      22  -none-    numeric
## importanceSD     0  -none-    NULL
## localImportance  0  -none-    NULL
## proximity        0  -none-    NULL
## ntree           1  -none-    numeric
## mtry            1  -none-    numeric
## forest          14  -none-    list
## y              3919 factor    numeric
## test            0  -none-    NULL
## inbag            0  -none-    NULL
## xNames          22  -none-    character
## problemType     1  -none-    character
## tuneValue        1 data.frame list
## obsLevels        2  -none-    character
## rf variable importance
##
##    only 20 most important variables shown (out of 22)
##
##
##                                Overall
## inq_last_6mths                 100.00000
## last_fico_range_high            90.48403
## dti                             37.57445
## revol_util                      36.77269
## fico_range_high                 28.22595
## term 60 months                  4.89063
## verification_statusVerified     3.74457

```

```

## purposedebt_consolidation      3.40080
## purposecredit_card             3.13706
## desc_empty1                   2.99225
## verification_statusSource Verified 2.67772
## purposeother                   2.61639
## purposehome_improvement        2.16608
## purposesmall_business          1.53603
## purposemajor_purchase          1.07948
## purposeeducational             0.73490
## purposemoving                  0.50541
## purposemedical                 0.46270
## purposewedding                 0.38302
## purposehouse                   0.06192
## Confusion Matrix and Statistics
##
##           Reference
## Prediction good bad
##      good  928 116
##      bad   48 214
##
##           Accuracy : 0.8744
##           95% CI : (0.8552, 0.8919)
##      No Information Rate : 0.7473
##      P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.6432
##      McNemar's Test P-Value : 1.678e-07
##
##           Sensitivity : 0.6485
##           Specificity : 0.9508
##           Pos Pred Value : 0.8168
##           Neg Pred Value : 0.8889
##           Prevalence : 0.2527
##           Detection Rate : 0.1639
##      Detection Prevalence : 0.2006
##           Balanced Accuracy : 0.7997
##
##           'Positive' Class : bad
##

```

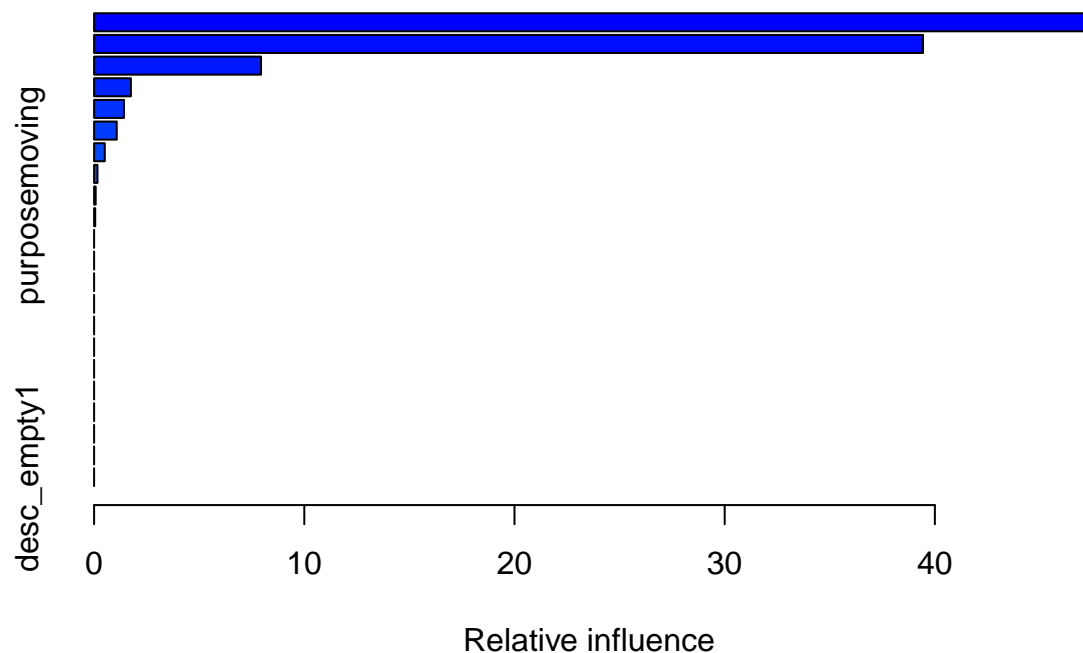



```
##
## Call:
## roc.default(response = dft_test$status, predictor = testProbs[,      "bad"])
##
## Data: testProbs[, "bad"] in 976 controls (dft_test$status good) < 330 cases (dft_test$status bad).
## Area under the curve: 0.9123
```

Gradient Boost Model

```
## Stochastic Gradient Boosting
##
## 3919 samples
##   9 predictor
##   2 classes: 'good', 'bad'
##
## No pre-processing
## Resampling: Bootstrapped (25 reps)
##
## Summary of sample sizes: 3919, 3919, 3919, 3919, 3919, 3919, ...
##
## Resampling results across tuning parameters:
##
##   interaction.depth  n.trees  Accuracy  Kappa    Accuracy SD
##   1                   50      0.8706639  0.6247338  0.005731525
##   1                   100     0.8708167  0.6284428  0.007067068
##   1                   150     0.8712036  0.6310435  0.007057299
##   2                    50      0.8773309  0.6465619  0.007626409
##   2                   100     0.8766634  0.6485005  0.006919372
```

```
##      2      150      0.8755198  0.6461899  0.006555678
##      3       50      0.8774845  0.6489184  0.007435278
##      3      100      0.8762049  0.6480520  0.005813588
##      3      150      0.8756033  0.6467955  0.006548232
## Kappa SD
## 0.01725275
## 0.02168284
## 0.02143313
## 0.02235166
## 0.02071361
## 0.01897217
## 0.02100417
## 0.01721180
## 0.02020654
##
## Tuning parameter 'shrinkage' was held constant at a value of 0.1
##
## Tuning parameter 'n.minobsinnode' was held constant at a value of 10
## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were n.trees = 50, interaction.depth
## = 3, shrinkage = 0.1 and n.minobsinnode = 10.
```



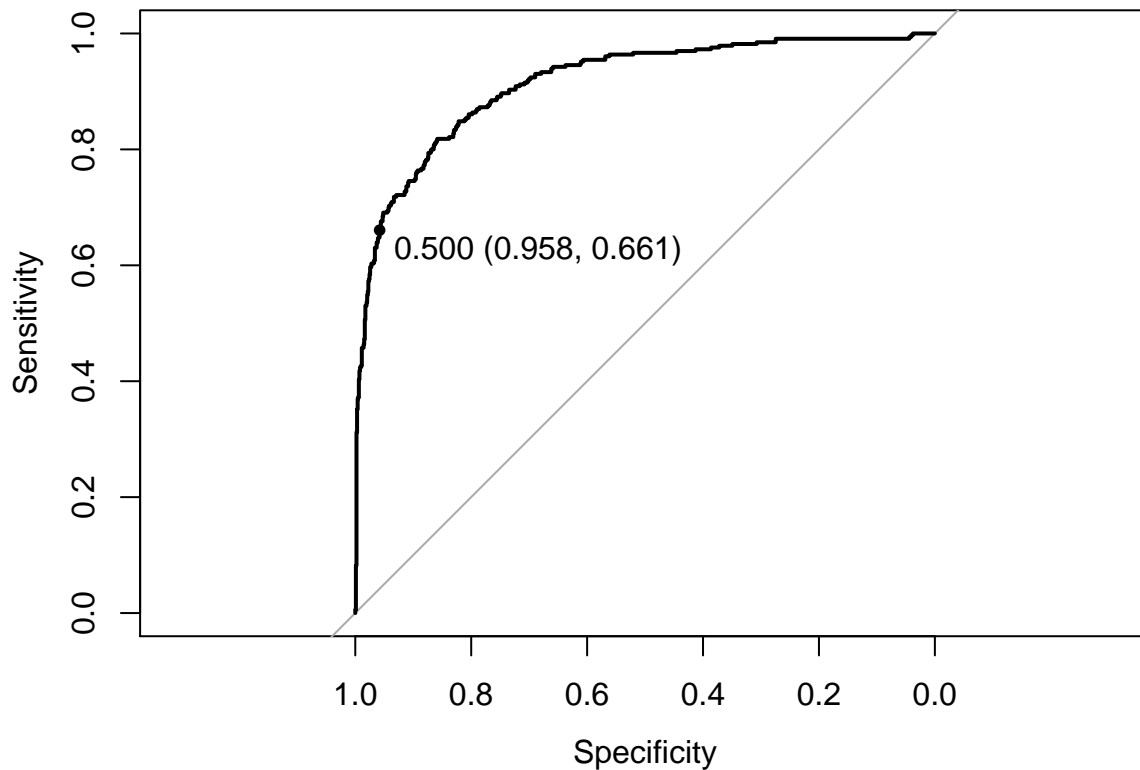
```
##                               var
## inq_last_6mths                inq_last_6mths
## last_fico_range_high          last_fico_range_high
## fico_range_high               fico_range_high
## dti                           dti
## revol_util                    revol_util
## term 60 months                term 60 months
## verification_statusSource Verified verification_statusSource Verified
## purposesmall_business         purposesmall_business
## purposemoving                 purposemoving
```

## purposehome_improvement		purposehome_improvement
## verification_statusVerified		verification_statusVerified
## purposecredit_card		purposecredit_card
## purposedebt_consolidation		purposedebt_consolidation
## purposeeducational		purposeeducational
## purposehouse		purposehouse
## purposemajor_purchase		purposemajor_purchase
## purposemedical		purposemedical
## purposeother		purposeother
## purposerenewable_energy		purposerenewable_energy
## purposevacation		purposevacation
## purposewedding		purposewedding
## desc_empty1		desc_empty1
##	rel.inf	
## inq_last_6mths	47.57363013	
## last_fico_range_high	39.43710594	
## fico_range_high	7.94244463	
## dti	1.75123514	
## revol_util	1.42006993	
## term 60 months	1.07767425	
## verification_statusSource Verified	0.50936980	
## purposesmall_business	0.16157295	
## purposemoving	0.07453077	
## purposehome_improvement	0.05236646	
## verification_statusVerified	0.00000000	
## purposecredit_card	0.00000000	
## purposedebt_consolidation	0.00000000	
## purposeeducational	0.00000000	
## purposehouse	0.00000000	
## purposemajor_purchase	0.00000000	
## purposemedical	0.00000000	
## purposeother	0.00000000	
## purposerenewable_energy	0.00000000	
## purposevacation	0.00000000	
## purposewedding	0.00000000	
## desc_empty1	0.00000000	
## gbm variable importance		
##		
##	only 20 most important variables shown (out of 22)	
##		
##	Overall	
## inq_last_6mths	100.0000	
## last_fico_range_high	82.8970	
## fico_range_high	16.6951	
## dti	3.6811	
## revol_util	2.9850	
## term 60 months	2.2653	
## verification_statusSource Verified	1.0707	
## purposesmall_business	0.3396	
## purposemoving	0.1567	
## purposehome_improvement	0.1101	
## purposemedical	0.0000	
## purposehouse	0.0000	
## purposedebt_consolidation	0.0000	

```

## purposewedding          0.0000
## purposerenewable_energy 0.0000
## desc_empty1             0.0000
## verification_statusVerified 0.0000
## purposevacation         0.0000
## purposeeducational      0.0000
## purposeother            0.0000
## Confusion Matrix and Statistics
##
##           Reference
## Prediction good bad
##      good  935 112
##      bad   41 218
##
##           Accuracy : 0.8828
##           95% CI : (0.8642, 0.8998)
##      No Information Rate : 0.7473
##      P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.666
##  McNemar's Test P-Value : 1.521e-08
##
##           Sensitivity : 0.6606
##           Specificity : 0.9580
##      Pos Pred Value : 0.8417
##      Neg Pred Value : 0.8930
##           Prevalence : 0.2527
##      Detection Rate : 0.1669
##      Detection Prevalence : 0.1983
##      Balanced Accuracy : 0.8093
##
##      'Positive' Class : bad
##

```



```
##
## Call:
## roc.default(response = dft_test$status, predictor = testProbs[,      "bad"])
##
## Data: testProbs[, "bad"] in 976 controls (dft_test$status good) < 330 cases (dft_test$status bad).
## Area under the curve: 0.9153
```

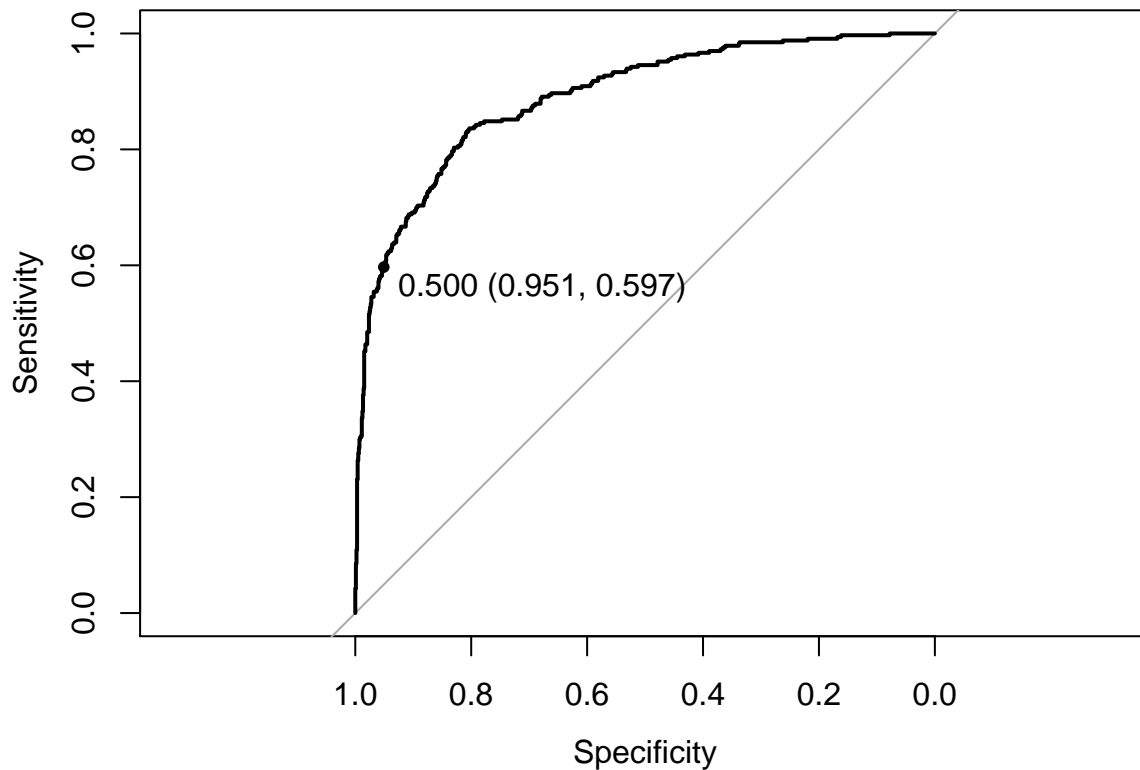
SVM Model

```
## Support Vector Machines with Radial Basis Function Kernel
##
## 3919 samples
##   9 predictor
##   2 classes: 'good', 'bad'
##
## Pre-processing: centered, scaled
## Resampling: Cross-Validated (10 fold)
##
## Summary of sample sizes: 3527, 3527, 3527, 3527, 3527, 3527, ...
##
## Resampling results across tuning parameters:
##
##   C      Accuracy  Kappa      Accuracy SD  Kappa SD
##   0.25  0.8415334  0.5484558  0.01656389  0.05347900
##   0.50  0.8481719  0.5653249  0.01244967  0.04018334
##   1.00  0.8517427  0.5740472  0.01332602  0.04561921
##
## Tuning parameter 'sigma' was held constant at a value of 0.04862674
```

```

## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were sigma = 0.04862674 and C = 1.
## Length Class Mode
##      1      ksvm      S4
## ROC curve variable importance
##
##                               Importance
## last_fico_range_high      100.000
## inq_last_6mths            79.399
## revol_util                33.918
## purpose                   25.494
## fico_range_high           22.987
## dti                       14.281
## term                      10.520
## desc_empty                2.972
## verification_status        0.000
## Confusion Matrix and Statistics
##
##           Reference
## Prediction good bad
##      good  928 133
##      bad   48 197
##
##           Accuracy : 0.8614
##           95% CI : (0.8415, 0.8797)
##      No Information Rate : 0.7473
##      P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.5988
##      McNemar's Test P-Value : 4.274e-10
##
##           Sensitivity : 0.5970
##           Specificity : 0.9508
##           Pos Pred Value : 0.8041
##           Neg Pred Value : 0.8746
##           Prevalence : 0.2527
##           Detection Rate : 0.1508
##      Detection Prevalence : 0.1876
##           Balanced Accuracy : 0.7739
##
##           'Positive' Class : bad
##

```



```
##
## Call:
## roc.default(response = dft_test$status, predictor = testProbs[,      "bad"])
##
## Data: testProbs[, "bad"] in 976 controls (dft_test$status good) < 330 cases (dft_test$status bad).
## Area under the curve: 0.8909
```

Results for Grade D Loans

Approximately 35% of the Grade D loans in this dataset went bad. With the four models, we were able to correctly predict between 65% and 77% of the bad loans. This predictive ability is based on a 50% probability classification cutoff. As the ROC curves show, it's possible to predict the bad loans with a higher probability, of course, with a higher false positive rate, though. The FICO range and the number of inquiries in the past 6 months were also important predictors for this loan grade.

Logistic Regression Model

```
## Warning in predict.lm(object, newdata, se.fit, scale = 1, type =
## ifelse(type == : prediction from a rank-deficient fit may be misleading

## Generalized Linear Model
##
## 2643 samples
##    9 predictor
##    2 classes: 'good', 'bad'
##
## No pre-processing
```

```

## Resampling: Bootstrapped (25 reps)
##
## Summary of sample sizes: 2643, 2643, 2643, 2643, 2643, 2643, ...
##
## Resampling results
##
##   Accuracy   Kappa     Accuracy SD   Kappa SD
##   0.8207444  0.5928377  0.01100261  0.02344815
##
##
## Call:
## NULL
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.5752  -0.6130  -0.3308   0.4964   3.2215
##
## Coefficients:
##                                     Estimate Std. Error z value
## (Intercept)                      5.2171095  2.0471105   2.549
## `term 60 months`                  0.0610230  0.1462789   0.417
## `verification_statusSource Verified` -0.4372620  0.1597814  -2.737
## verification_statusVerified        -0.2843958  0.1296191  -2.194
## purposecredit_card                 1.2086063  0.3926973   3.078
## purposedebt_consolidation          0.4878831  0.3643532   1.339
## purposeeducational                1.3136310  0.5702381   2.304
## purposehome_improvement           0.2729951  0.4242896   0.643
## purposehouse                      0.6811404  0.6894088   0.988
## purposemajor_purchase             0.5530122  0.4252511   1.300
## purposemedical                    0.5246664  0.5701721   0.920
## purposemoving                     0.3915786  0.5886544   0.665
## purposeother                      0.8453995  0.3914984   2.159
## purposerenewable_energy          -1.4262403  1.5926639  -0.896
## purposesmall_business             0.6601729  0.4265684   1.548
## purposevacation                   0.2137037  0.7363585   0.290
## purposewedding                    0.2009450  0.5146825   0.390
## fico_range_high                   0.0039129  0.0029939   1.307
## inq_last_6mths                    0.7840041  0.0377272  20.781
## revol_util                        -0.0033938  0.0022172  -1.531
## last_fico_range_high              -0.0156092  0.0007861 -19.856
## desc_empty1                      -0.4371316  0.1401979  -3.118
## dti                               0.0035591  0.0087704   0.406
##
## Pr(>|z|)
## (Intercept)                      0.01082 *
## `term 60 months`                  0.67656
## `verification_statusSource Verified` 0.00621 **
## verification_statusVerified        0.02823 *
## purposecredit_card                 0.00209 **
## purposedebt_consolidation          0.18056
## purposeeducational                0.02124 *
## purposehome_improvement           0.51995
## purposehouse                      0.32315
## purposemajor_purchase             0.19345

```

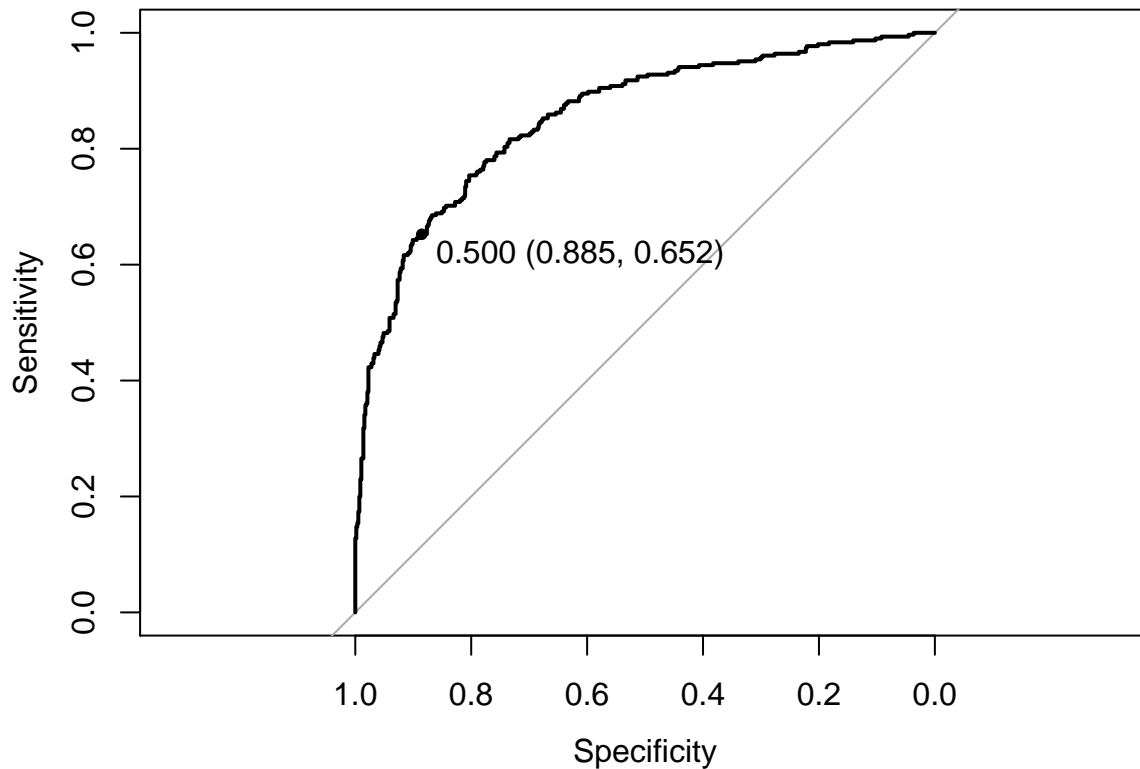


```

## purposemedical          0.35747
## purposemoving           0.50592
## purposeother            0.03082 *
## purposerenewable_energy 0.37052
## purposesmall_business   0.12171
## purposevacation         0.77165
## purposewedding          0.69622
## fico_range_high         0.19124
## inq_last_6mths          < 2e-16 ***
## revol_util              0.12585
## last_fico_range_high    < 2e-16 ***
## desc_empty1             0.00182 **
## dti                     0.68488
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 3409.8  on 2642  degrees of freedom
## Residual deviance: 2127.4  on 2620  degrees of freedom
## AIC: 2173.4
##
## Number of Fisher Scoring iterations: 5
##
## glm variable importance
##
##    only 20 most important variables shown (out of 22)
##
##
##                                Overall
## inq_last_6mths                 100.0000
## last_fico_range_high           95.4858
## desc_empty1                   13.8002
## purposecredit_card             13.6037
## `verification_statusSource Verified` 11.9391
## purposeeducational             9.8261
## verification_statusVerified     9.2914
## purposeother                   9.1221
## purposesmall_business          6.1365
## revol_util                     6.0538
## purposedebt_consolidation       5.1185
## fico_range_high                4.9618
## purposemajor_purchase          4.9301
## purposehouse                   3.4054
## purposemedical                 3.0744
## purposerenewable_energy        2.9540
## purposemoving                  1.8301
## purposehome_improvement        1.7237
## `term 60 months`              0.6196
## dti                            0.5641
## Confusion Matrix and Statistics
##
##              Reference
## Prediction good bad
##      good  509 106

```

```
##      bad      66 199
##
##      Accuracy : 0.8045
##      95% CI : (0.7768, 0.8303)
##      No Information Rate : 0.6534
##      P-Value [Acc > NIR] : < 2.2e-16
##
##      Kappa : 0.5548
##      McNemar's Test P-Value : 0.002942
##
##      Sensitivity : 0.6525
##      Specificity : 0.8852
##      Pos Pred Value : 0.7509
##      Neg Pred Value : 0.8276
##      Prevalence : 0.3466
##      Detection Rate : 0.2261
##      Detection Prevalence : 0.3011
##      Balanced Accuracy : 0.7688
##
##      'Positive' Class : bad
##
```



```
##
## Call:
## roc.default(response = dft_test$status, predictor = testProbs[,      "bad"])
##
## Data: testProbs[, "bad"] in 575 controls (dft_test$status good) < 305 cases (dft_test$status bad).
## Area under the curve: 0.8547
```

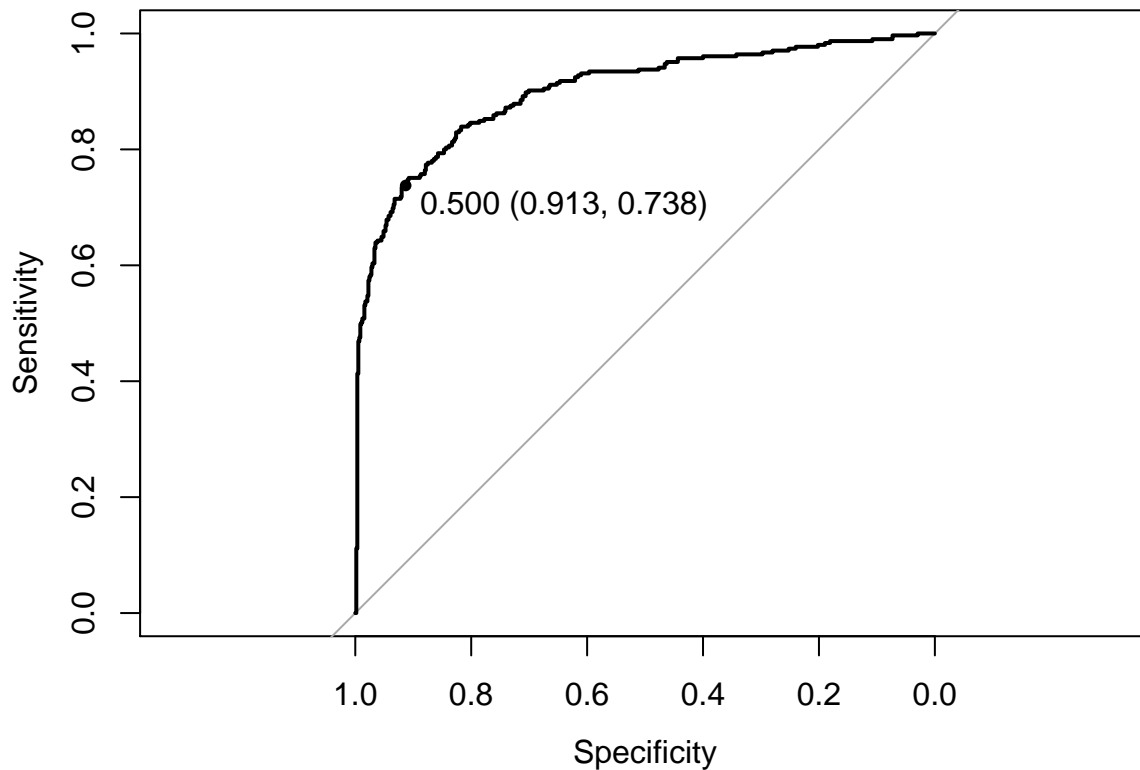
Random Forest Model

```
## Random Forest
##
## 2643 samples
##    9 predictor
##    2 classes: 'good', 'bad'
##
## No pre-processing
## Resampling: Bootstrapped (25 reps)
##
## Summary of sample sizes: 2643, 2643, 2643, 2643, 2643, 2643, ...
##
## Resampling results across tuning parameters:
##
##  mtry  Accuracy   Kappa     Accuracy SD  Kappa SD
##    2    0.8395431 0.6090853 0.014538549 0.03667554
##   12    0.8690558 0.7025225 0.008439853 0.02002008
##   22    0.8603770 0.6835663 0.010225713 0.02431546
##
## Accuracy was used to select the optimal model using  the largest value.
## The final value used for the model was mtry = 12.
##
##      Length Class      Mode
## call           4  -none-    call
## type           1  -none-    character
## predicted      2643 factor    numeric
## err.rate       1500 -none-    numeric
## confusion       6  -none-    numeric
## votes          5286 matrix    numeric
## oob.times       2643 -none-    numeric
## classes         2  -none-    character
## importance      22  -none-    numeric
## importanceSD     0  -none-    NULL
## localImportance  0  -none-    NULL
## proximity        0  -none-    NULL
## ntree           1  -none-    numeric
## mtry            1  -none-    numeric
## forest          14  -none-    list
## y               2643 factor    numeric
## test            0  -none-    NULL
## inbag            0  -none-    NULL
## xNames          22  -none-    character
## problemType     1  -none-    character
## tuneValue        1 data.frame list
## obsLevels        2  -none-    character
## rf variable importance
##
##    only 20 most important variables shown (out of 22)
##
##
##                                Overall
## inq_last_6mths                 100.0000
## last_fico_range_high            72.2777
## fico_range_high                 33.1404
## dti                             32.4094
```

```

## revol_util 27.4758
## term 60 months 4.0244
## verification_statusVerified 3.4194
## purposedebt_consolidation 2.9417
## desc_empty1 2.8152
## purposecredit_card 2.3297
## verification_statusSource Verified 2.3072
## purposeother 2.0685
## purposesmall_business 1.5509
## purposehome_improvement 1.1140
## purposemajor_purchase 1.0537
## purposemedical 0.9300
## purposeeducational 0.7923
## purposehouse 0.6226
## purposemoving 0.5915
## purposewedding 0.4943
## Confusion Matrix and Statistics
##
##           Reference
## Prediction good bad
##      good  525  80
##      bad   50 225
##
##           Accuracy : 0.8523
##           95% CI : (0.8271, 0.8751)
##      No Information Rate : 0.6534
##      P-Value [Acc > NIR] : < 2e-16
##
##           Kappa : 0.6661
##      McNemar's Test P-Value : 0.01098
##
##           Sensitivity : 0.7377
##           Specificity : 0.9130
##           Pos Pred Value : 0.8182
##           Neg Pred Value : 0.8678
##           Prevalence : 0.3466
##           Detection Rate : 0.2557
##      Detection Prevalence : 0.3125
##           Balanced Accuracy : 0.8254
##
##           'Positive' Class : bad
##

```

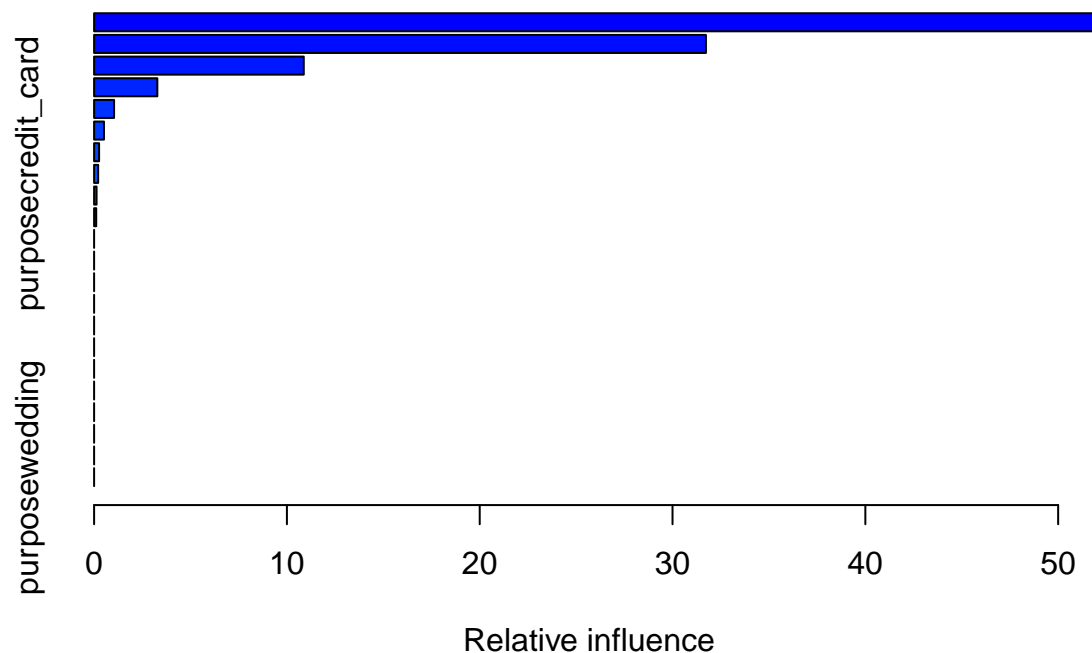


```
##
## Call:
## roc.default(response = dft_test$status, predictor = testProbs[,      "bad"])
##
## Data: testProbs[, "bad"] in 575 controls (dft_test$status good) < 305 cases (dft_test$status bad).
## Area under the curve: 0.9012
```

Gradient Boost Model

```
## Stochastic Gradient Boosting
##
## 2643 samples
##   9 predictor
##   2 classes: 'good', 'bad'
##
## No pre-processing
## Resampling: Bootstrapped (25 reps)
##
## Summary of sample sizes: 2643, 2643, 2643, 2643, 2643, 2643, ...
##
## Resampling results across tuning parameters:
##
##   interaction.depth  n.trees  Accuracy  Kappa    Accuracy SD
##   1                  50      0.8660421  0.6902063  0.009447489
##   1                  100     0.8693640  0.6997814  0.009290477
##   1                  150     0.8706833  0.7033118  0.009872404
##   2                   50      0.8712706  0.7045937  0.008514818
##   2                  100     0.8714019  0.7056883  0.009501993
```

```
##      2          150      0.8700134 0.7028419 0.009229986
##      3           50      0.8716738 0.7060678 0.009554344
##      3          100      0.8704110 0.7040435 0.009526368
##      3          150      0.8690296 0.7012735 0.009475418
##      Kappa SD
##      0.02053097
##      0.02104011
##      0.02163652
##      0.01860828
##      0.02047966
##      0.01985766
##      0.02083295
##      0.02078379
##      0.02117697
##
## Tuning parameter 'shrinkage' was held constant at a value of 0.1
##
## Tuning parameter 'n.minobsinnode' was held constant at a value of 10
## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were n.trees = 50, interaction.depth
## = 3, shrinkage = 0.1 and n.minobsinnode = 10.
```



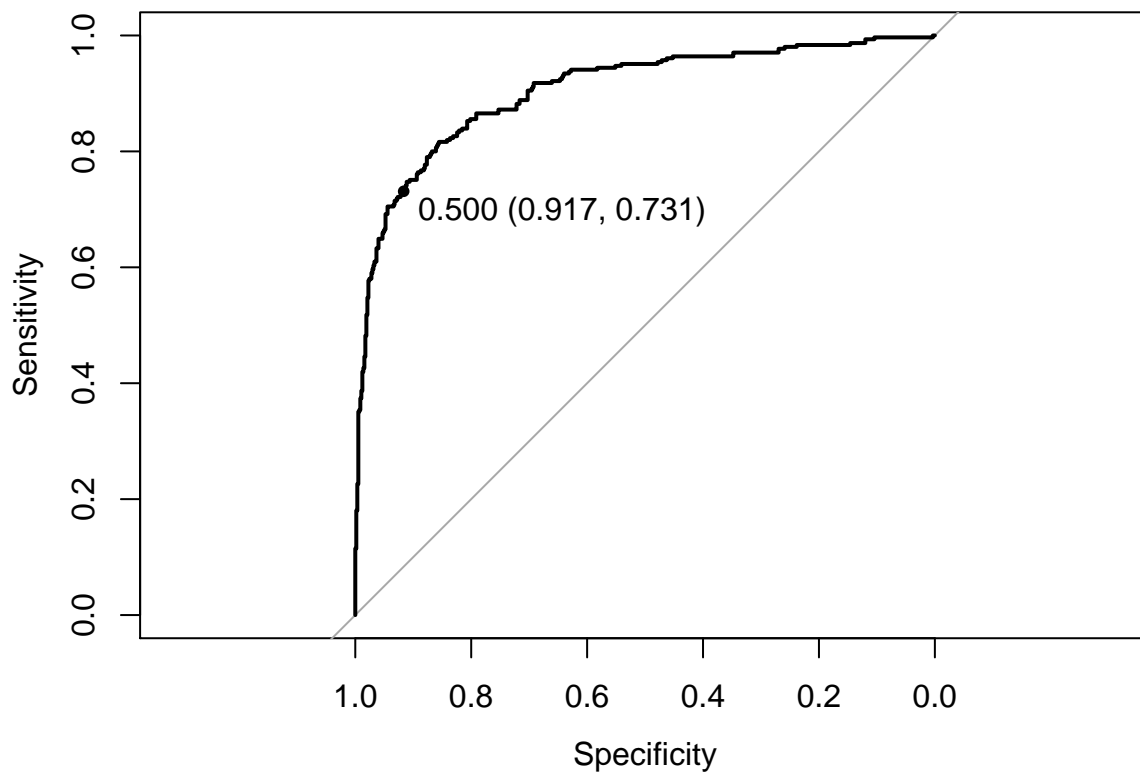
```
##      var
## inq_last_6mths      inq_last_6mths
## last_fico_range_high last_fico_range_high
## fico_range_high      fico_range_high
## dti      dti
## revol_util      revol_util
## term 60 months      term 60 months
## desc_empty1      desc_empty1
## purposecredit_card      purposecredit_card
## verification_statusSource Verified verification_statusSource Verified
```

## purposeeducational		purposeeducational
## verification_statusVerified		verification_statusVerified
## purposedebt_consolidation		purposedebt_consolidation
## purposehome_improvement		purposehome_improvement
## purposehouse		purposehouse
## purposemajor_purchase		purposemajor_purchase
## purposemedical		purposemedical
## purposemoving		purposemoving
## purposeother		purposeother
## purposerenewable_energy		purposerenewable_energy
## purposesmall_business		purposesmall_business
## purposevacation		purposevacation
## purposewedding		purposewedding
##	rel.inf	
## inq_last_6mths	51.8619695	
## last_fico_range_high	31.7323615	
## fico_range_high	10.8741519	
## dti	3.2794162	
## revol_util	1.0362737	
## term 60 months	0.5114909	
## desc_empty1	0.2581910	
## purposecredit_card	0.2100264	
## verification_statusSource Verified	0.1273153	
## purposeeducational	0.1088036	
## verification_statusVerified	0.0000000	
## purposedebt_consolidation	0.0000000	
## purposehome_improvement	0.0000000	
## purposehouse	0.0000000	
## purposemajor_purchase	0.0000000	
## purposemedical	0.0000000	
## purposemoving	0.0000000	
## purposeother	0.0000000	
## purposerenewable_energy	0.0000000	
## purposesmall_business	0.0000000	
## purposevacation	0.0000000	
## purposewedding	0.0000000	
## gbm variable importance		
##		
##	only 20 most important variables shown (out of 22)	
##		
##	Overall	
## inq_last_6mths	100.0000	
## last_fico_range_high	61.1862	
## fico_range_high	20.9675	
## dti	6.3234	
## revol_util	1.9981	
## term 60 months	0.9863	
## desc_empty1	0.4978	
## purposecredit_card	0.4050	
## verification_statusSource Verified	0.2455	
## purposeeducational	0.2098	
## purposemoving	0.0000	
## purposevacation	0.0000	
## purposewedding	0.0000	

```

## purposemedical 0.0000
## verification_statusVerified 0.0000
## purposedebt_consolidation 0.0000
## purposesmall_business 0.0000
## purposemajor_purchase 0.0000
## purposeother 0.0000
## purposerenewable_energy 0.0000
## Confusion Matrix and Statistics
##
##           Reference
## Prediction good bad
##      good  527  82
##      bad   48 223
##
##           Accuracy : 0.8523
##           95% CI : (0.8271, 0.8751)
##      No Information Rate : 0.6534
##      P-Value [Acc > NIR] : <2e-16
##
##           Kappa : 0.6651
##  McNemar's Test P-Value : 0.0038
##
##           Sensitivity : 0.7311
##           Specificity : 0.9165
##      Pos Pred Value : 0.8229
##      Neg Pred Value : 0.8654
##           Prevalence : 0.3466
##      Detection Rate : 0.2534
##      Detection Prevalence : 0.3080
##      Balanced Accuracy : 0.8238
##
##      'Positive' Class : bad
##

```

```
##
## Call:
## roc.default(response = dft_test$status, predictor = testProbs[,      "bad"])
##
## Data: testProbs[, "bad"] in 575 controls (dft_test$status good) < 305 cases (dft_test$status bad).
## Area under the curve: 0.9058
```

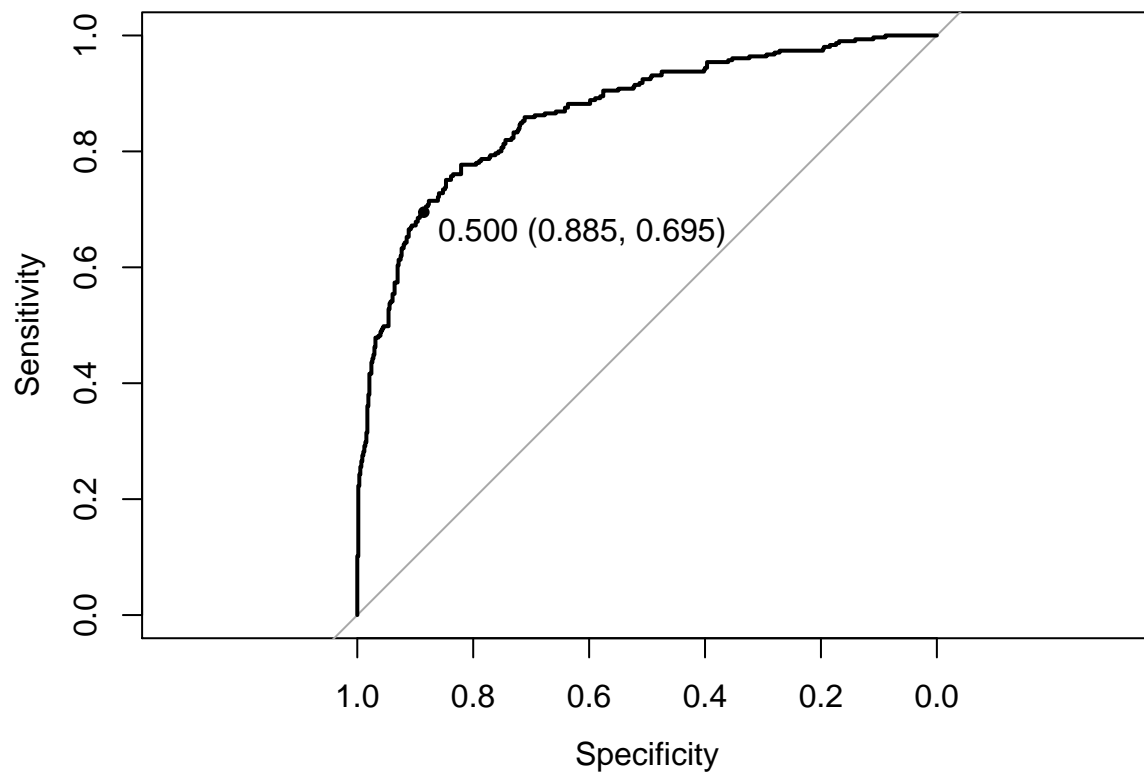
SVM Model

```
## Support Vector Machines with Radial Basis Function Kernel
##
## 2643 samples
##   9 predictor
##   2 classes: 'good', 'bad'
##
## Pre-processing: centered, scaled
## Resampling: Cross-Validated (10 fold)
##
## Summary of sample sizes: 2379, 2378, 2378, 2379, 2378, 2380, ...
##
## Resampling results across tuning parameters:
##
##   C      Accuracy  Kappa      Accuracy SD  Kappa SD
##   0.25  0.8187223  0.5890269  0.03031166  0.07075603
##   0.50  0.8285737  0.6102924  0.02797389  0.06552996
##   1.00  0.8304763  0.6155251  0.02572944  0.06043479
##
## Tuning parameter 'sigma' was held constant at a value of 0.04529224
```

```

## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were sigma = 0.04529224 and C = 1.
## Length Class Mode
##      1      ksvm      S4
## ROC curve variable importance
##
##                               Importance
## last_fico_range_high      100.000
## inq_last_6mths           95.726
## verification_status       34.398
## revol_util                31.174
## fico_range_high           24.776
## purpose                    22.957
## dti                       14.213
## term                       3.347
## desc_empty                 0.000
## Confusion Matrix and Statistics
##
##              Reference
## Prediction good bad
##      good  509  93
##      bad   66 212
##
##              Accuracy : 0.8193
##              95% CI : (0.7923, 0.8442)
##      No Information Rate : 0.6534
##      P-Value [Acc > NIR] : < 2e-16
##
##              Kappa : 0.5926
##      McNemar's Test P-Value : 0.03921
##
##              Sensitivity : 0.6951
##              Specificity : 0.8852
##              Pos Pred Value : 0.7626
##              Neg Pred Value : 0.8455
##              Prevalence : 0.3466
##              Detection Rate : 0.2409
##      Detection Prevalence : 0.3159
##              Balanced Accuracy : 0.7901
##
##      'Positive' Class : bad
##

```



```
##
## Call:
## roc.default(response = dft_test$status, predictor = testProbs[,      "bad"])
##
## Data: testProbs[, "bad"] in 575 controls (dft_test$status good) < 305 cases (dft_test$status bad).
## Area under the curve: 0.8677
```