

data_challenge

Zongyan Wang

March 25, 2016

Read data and packages

```
# Upload packages for data clean and analysis  
require(plyr) #data clean
```

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

```
require(dplyr) #data clean
```

```
## Loading required package: dplyr
```

```
##
```

```
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:plyr':
```

```
##
```

```
##      arrange, count, desc, failwith, id, mutate, rename, summarise,  
##      summarize
```

```
## The following objects are masked from 'package:stats':
```

```
##
```

```
##      filter, lag
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
##      intersect, setdiff, setequal, union
```

```
require(tidyr) #data clean
```

```
## Loading required package: tidyr
```

```
require(ggplot2) # visualization
```

```
## Loading required package: ggplot2
```

```
require(XML) # web scriping
```

```
## Loading required package: XML
```

```
require(testthat) #test model
```

```
## Loading required package: testthat
```

```
require(kernlab) # kernel and SVM
```

```
## Loading required package: kernlab
```

```
##
```

```
## Attaching package: 'kernlab'
```

```
## The following object is masked from 'package:ggplot2':
```

```
##
```

```
##      alpha
```

```
require(reshape2)
```

```
## Loading required package: reshape2
```

```
require(caret) # Cross - Validation k fold
```

```
## Loading required package: caret
```

```
## Loading required package: lattice
```

```
require(maps) # heatmap: map_data
```

```
## Loading required package: maps
```

```
##
```

```
## # maps v3.1: updated 'world': all lakes moved to separate new #
```

```
## # 'lakes' database. Type '?world' or 'news(package="maps")'. #
```

```
##
```

```
## Attaching package: 'maps'
```

```
## The following object is masked from 'package:plyr':
```

```
##
```

```
##      ozone
```

```
require(ggmap) # for heatmap
```

```
## Loading required package: ggmap
```

Read file

```
# Read Data and function files
LC <- read.csv(file.choose(), header = T) #Please read LC_biz_all.csv
source("clean_f2.R") # functions for cleaning the data
source("plot_f2.R") # functions for plot(heatmap)
source("analysis_f2.R") # functions for plot and data analysis
```

Data clean

```
lc <- perToN(LC) # transfer percentage to numeric
lc <- replaceBlank_all(lc) # replace the blank
lc <- dateToNum(lc) # transform earliest_cr_line from date to numeric
lc <- lc %>% select(-id, -emp_title,
                  -last_pymnt_d, -next_pymnt_d,
                  -zip_code, -issue_d, -last_credit_pull_d)
```

Data Summary

```
n.factor_all(lc)
```

```
##      term emp_length home_ownership verification_status loan_status purpose
## 1      2          12              3              3          7          1
##   addr_state
## 1          48
```

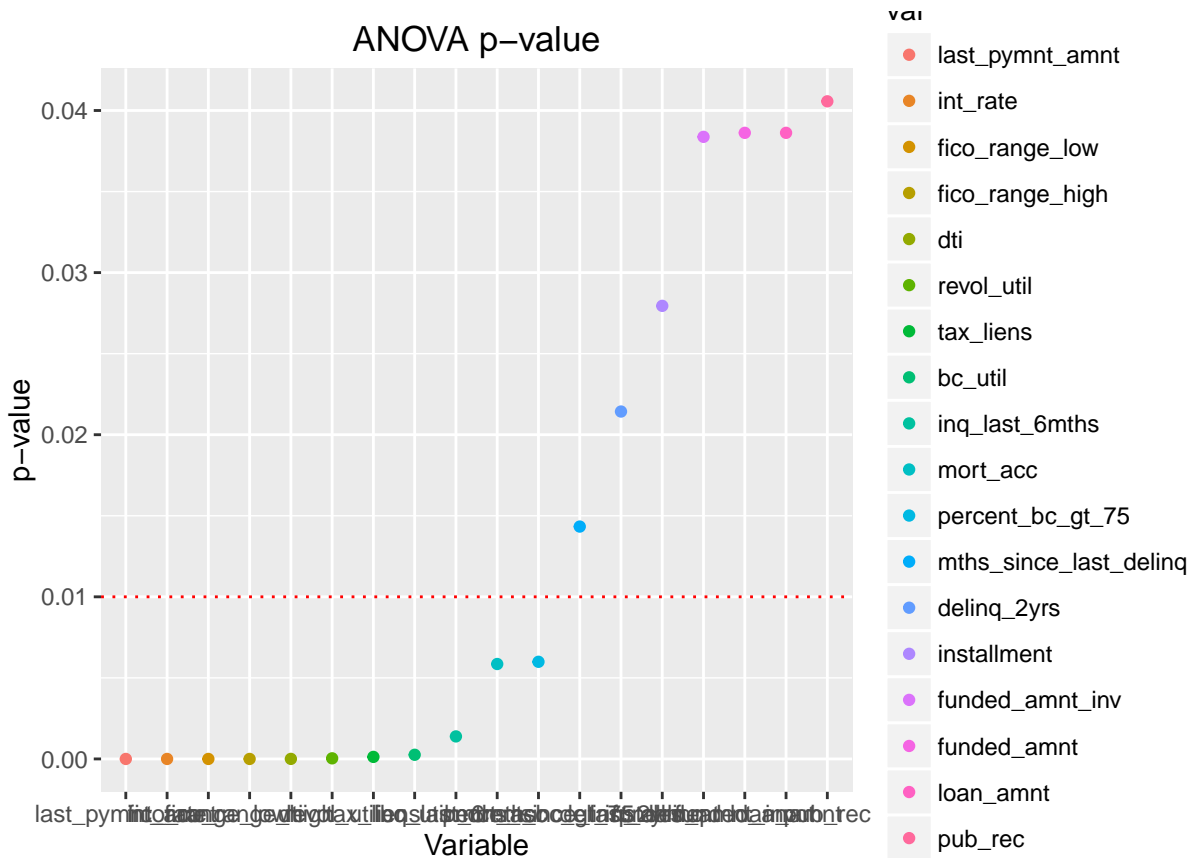
```
hasNA_all(lc)
```

```
##      mths_since_last_delinq      mths_since_last_record
##                2421                4302
##      revol_util mths_since_last_major_derog
##                3                3674
##      bc_util      mths_since_recent_bc
##                77                68
##      mths_since_recent_bc_dlq      mths_since_recent_inq
##                3932                381
##      num_tl_120dpd_2m      percent_bc_gt_75
##                275                77
```

ANOVA

```
lc1 <- lc[, !(colnames(lc) %in% c(names(n.factor_all(lc))))]
lc1$loan_status <- lc$loan_status
lc1 <- lc1 %>% subset(!loan_status == "Current")
anova.p <- data.frame(var = names(anova(lc1)), p.value = anova(lc1))
anova.p_plot <- anova.p %>% subset(p.value < 0.05) %>%
  mutate(var = reorder(x = var, X = p.value, min)) %>%
  ggplot(aes(x = var, y = p.value, color = var)) +
```

```
geom_point() +
geom_hline(yintercept = 0.01, colour = "red", linetype = 3) +
  labs(x = "Variable", y = "p-value", title = "ANOVA p-value")
anova.p_plot
```



```
sig.names <- as.character((anova.p %>% subset(p.value < 0.01))$var)
var.names <- c(sig.names, names(n.factor_all(lc)))
```

subset the data

```
lc2 <- lc[,c(var.names)]
lc2 <- lc2 %>% select(-purpose)
lc.categorical <- lc2[,names(n.factor_all(lc2))]
lc.categorical <- lc.categorical %>% select(-loan_status)
lc.numeric <- lc2[,!names(lc2) %in% names(n.factor_all(lc2))]
# Transform the categorical column to multiple numeric columns
categorical.list <- apply(lc.categorical, 2, function(x) model.matrix(~ x + 0))
lc3 <- as.data.frame(cbind(categorical.list, lc.numeric))
lc3$loan_status <- lc$loan_status
# Group the loan_status into 3 groups
lc3 <- lc3 %>% mutate(y = !(loan_status == "Current"))
lc3$y[lc3$loan_status == "Fully Paid"] <- 2
lc3 <- lc3 %>% select(-loan_status)
```

Logistic Regression

```
lc4 <- lc3 %>% subset(y != 0)
lc4$y[lc4$y == 2] <- 0
log.f <- logistic(lc4)
```

```
## Warning: glm.fit: algorithm did not converge
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
pred <- predict(log.f, lc4, type = "response")
```

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

```
pred1 <- as.numeric(pred >= 0.5)
# in-sample prediction accuracy
1 - sum(pred1 == lc4$y, na.rm = T)/length(pred1)
```

```
## [1] 0.06681191
```

```
# out-of-sample prediction accuracy
cv_error <- cv_k(lc4)
```

```
## Warning: glm.fit: algorithm did not converge
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

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

```
## Warning: glm.fit: algorithm did not converge
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

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

```
## Warning: glm.fit: algorithm did not converge
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

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

```
## Warning: glm.fit: algorithm did not converge
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```

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

## Warning: glm.fit: algorithm did not converge

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

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

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

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

## Warning: glm.fit: algorithm did not converge

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

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

## Warning: glm.fit: algorithm did not converge

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

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

## Warning: glm.fit: algorithm did not converge

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

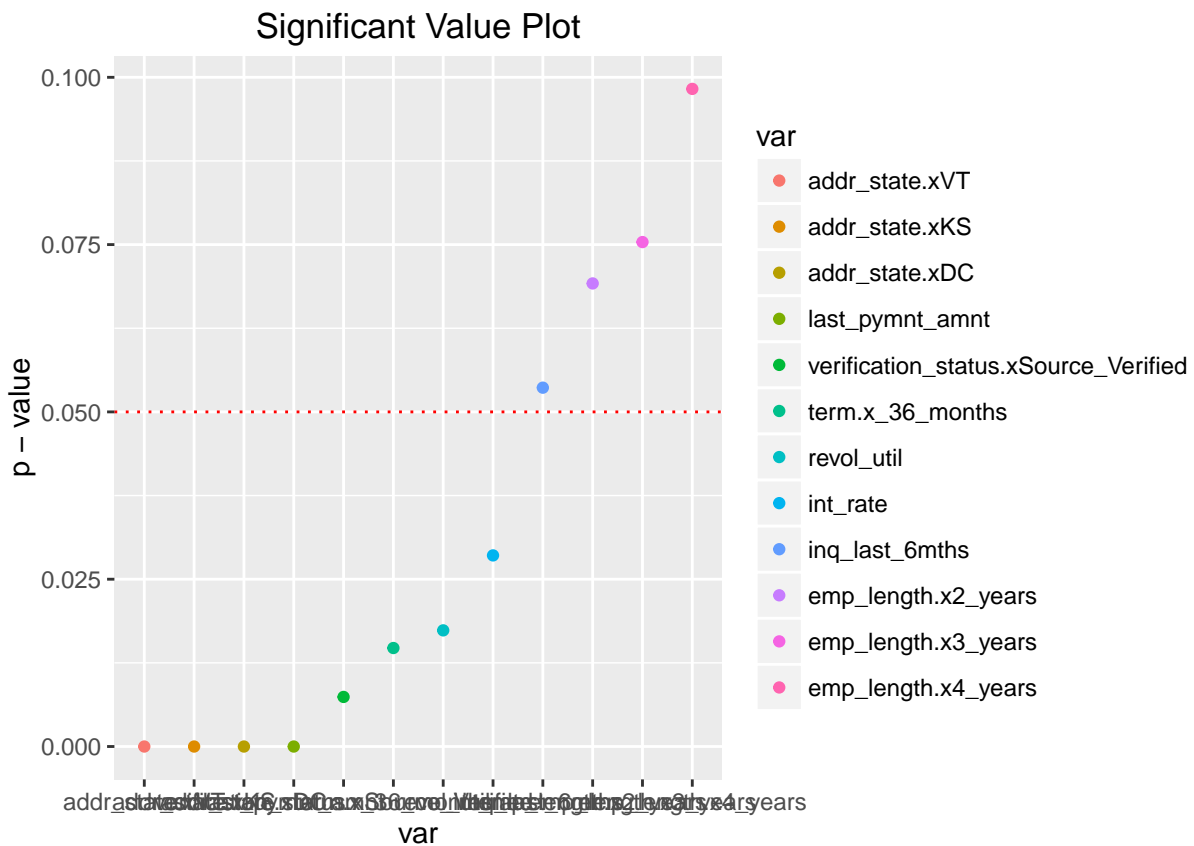
## Warning in predict.lm(object, newdata, se.fit, scale = 1, type =
## ifelse(type == : prediction from a rank-deficient fit may be misleading

cv_error

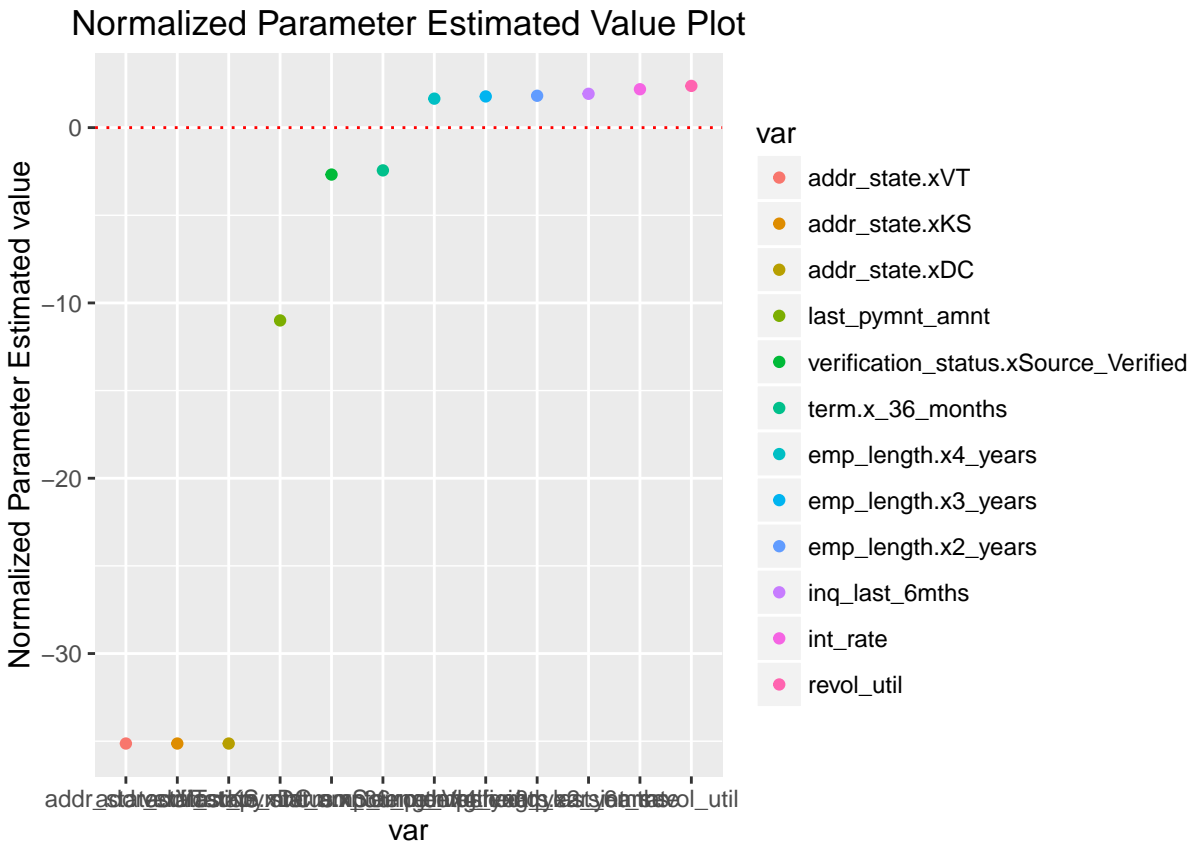
## [1] 0.0879505

```

```
#Plot the significant value
significant_plot <- coeff_plot1(log.f, sig = T, 0.1)
print(significant_plot)
```



```
#Plot the coefficient value
coefficient_plot <- coeff_plot1(log.f, sig = F, 0.1)
print(coefficient_plot)
```



```
## Kernel SVM
```

```
# Prepare for kernel SVM
lc4 <- lc4[-which(is.na(lc4$percent_bc_gt_75) | is.na(lc4$revol_util) | is.na(lc4$bc_util)),]

l.kern <- svm(lc4, "vanilladot")
```

```
## Setting default kernel parameters
```

```
## Warning in .local(x, ...): Variable(s) `` constant. Cannot scale data.
```

```
l.kern
```

```
## Support Vector Machine object of class "ksvm"
##
## SV type: C-svc (classification)
## parameter : cost C = 1
##
## Linear (vanilla) kernel function.
##
## Number of Support Vectors : 193
##
## Objective Function Value : -207.8558
## Training error : 0.051546
```



```
# New data
newdata <- as.matrix(lc4$y %>% select(-y))
pred.kern <- predict(l.kern, newdata)
# in-sample prediction accuracy
1 - sum(pred.kern == lc4$y, na.rm = T)/length(lc4$y)
```

```
## [1] 0.05154639
```

```
# cross validation error
cv.kernel.error <- cv_ksvm(lc4, "vanilladot", k = 10)
```

```
## Setting default kernel parameters
```

```
## Warning in .local(x, ...): Variable(s) ` ` constant. Cannot scale data.
```

```
## Setting default kernel parameters
```

```
## Warning in .local(x, ...): Variable(s) `` constant. Cannot scale data.
```

```
## Setting default kernel parameters
```

```
## Warning in .local(x, ...): Variable(s) ` ` constant. Cannot scale data.
```

```
## Setting default kernel parameters
```

```
## Warning in .local(x, ...): Variable(s) `` constant. Cannot scale data.
```

```
## Setting default kernel parameters
```

```
## Warning in .local(x, ...): Variable(s) ` ` constant. Cannot scale data.
```

```
## Setting default kernel parameters
```

```
## Warning in .local(x, ...): Variable(s) ` ` constant. Cannot scale data.
```

```
## Setting default kernel parameters
```

```
## Warning in .local(x, ...): Variable(s) ` ` constant. Cannot scale data.
```

```
## Setting default kernel parameters
```

```
## Warning in .local(x, ...): Variable(s) ``' constant. Cannot scale data.
```

```
## Setting default kernel parameters
```

```
## Warning in .local(x, ...): Variable(s) ``' constant. Cannot scale data.
```

```
## Setting default kernel parameters
```

```
## Warning in .local(x, ...): Variable(s) `` constant. Cannot scale data.
```

```
cv.kernel.error
```

```
## [1] 0.05745098
```