

# homework1

Na Yun Cho

```
library(glmnet)
library(caret)
library(corrplot)
library(plotmo)
library(tidyverse)
library(pls)

train_data = read.csv("./data/solubility_train.csv")
train_data <- na.omit(train_data)

test_data = read.csv("./data/solubility_test.csv")
test_data <- na.omit(test_data)

train2 <- model.matrix(Solubility ~ ., train_data)[ , -1]
test2 <- model.matrix(Solubility ~ ., test_data)[ , -1]

#matrix of predictors
x <- train2
y1 <- test2
#vector of response
y <- train_data$Solubility
y2 <- test_data$Solubility
```

## Part(a)

```
#fit a linear model using least squares
ctrl1 <- trainControl(method = "cv", number = 10)
set.seed(1)
lm.fit <- train(x, y,
               method = "lm",
               trControl = ctrl1)

#calculate mean squared error using test data
linear.pred <- predict(lm.fit, newdata = test2)
mean((linear.pred - test_data$Solubility)^2)
```

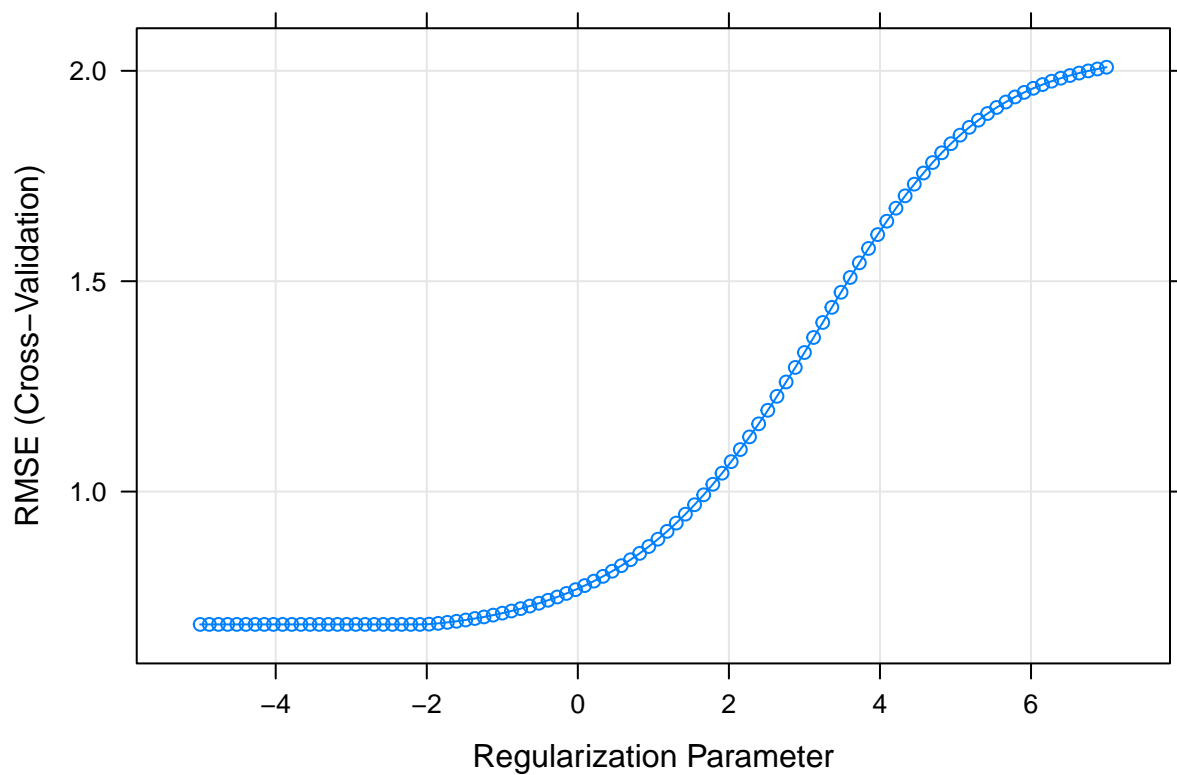
```
## [1] 0.5558898
```

The mean squared error using the test data is 0.5558898

## Part(b)

```
#fit a ridge regression model
ctrl1 <- trainControl(method = "cv", number = 10)
set.seed(1)
ridge.fit <- train(x, y,
  method = "glmnet",
  tuneGrid = expand.grid(alpha = 0,
    lambda = exp(seq(7, -5, length=100))),
  preProc = c("center", "scale"),
  trControl = ctrl1)

#choose lambda
plot(ridge.fit, xTrans = log)
```



```
best_ridge = ridge.fit$bestTune

#calculate test error
ridge.pred <- predict(ridge.fit, newdata = test2)
mean((ridge.pred - test_data$Solubility)^2)
```

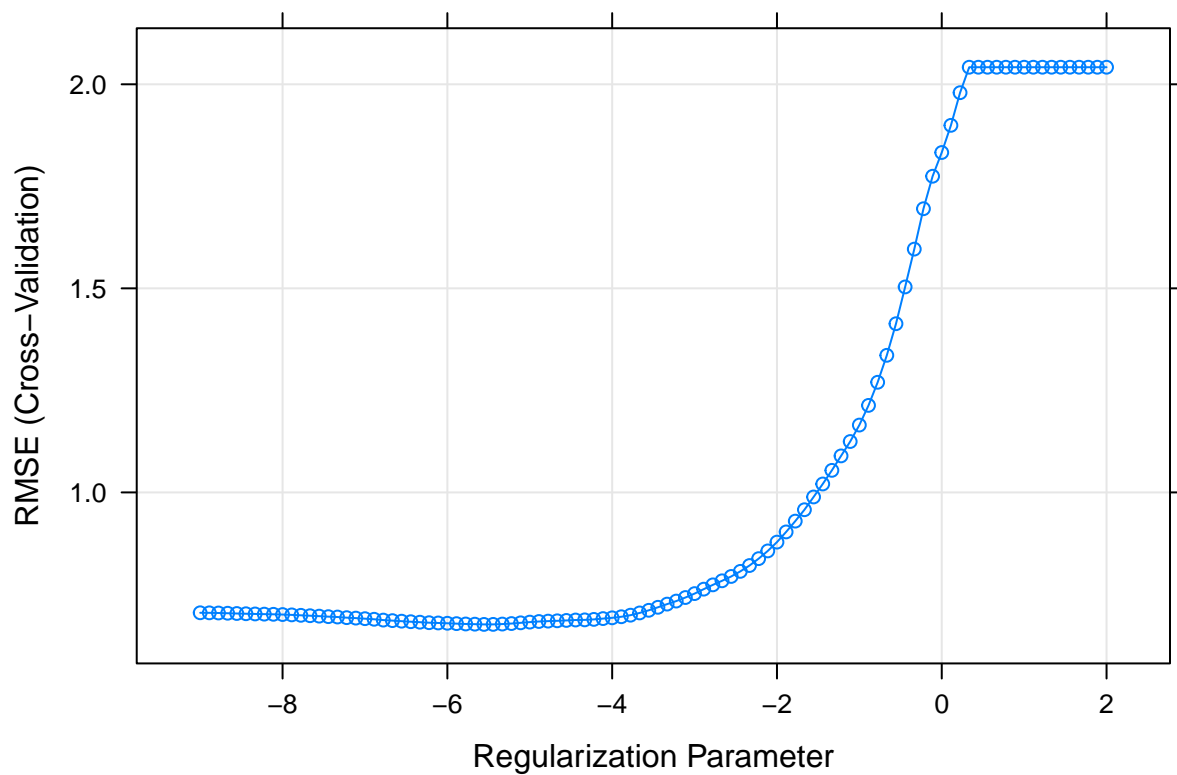
```
## [1] 0.5134603
```

The chosen lambda is 0.1235747 and the test error is 0.5134603

## Part(c)

```
#fit a lasso model
set.seed(1)
lasso.fit <- train(x, y,
  method = "glmnet",
  tuneGrid = expand.grid(alpha = 1,
    lambda = exp(seq(2, -9, length=100))),
  preProc = c("center", "scale"),
  trControl = ctrl1)

#choose lambda
plot(lasso.fit, xTrans = log)
```



```
best_lasso = lasso.fit$bestTune

#calculate test error
lasso.pred <- predict(lasso.fit, newdata = test2)
mean((lasso.pred - test_data$Solubility)^2)

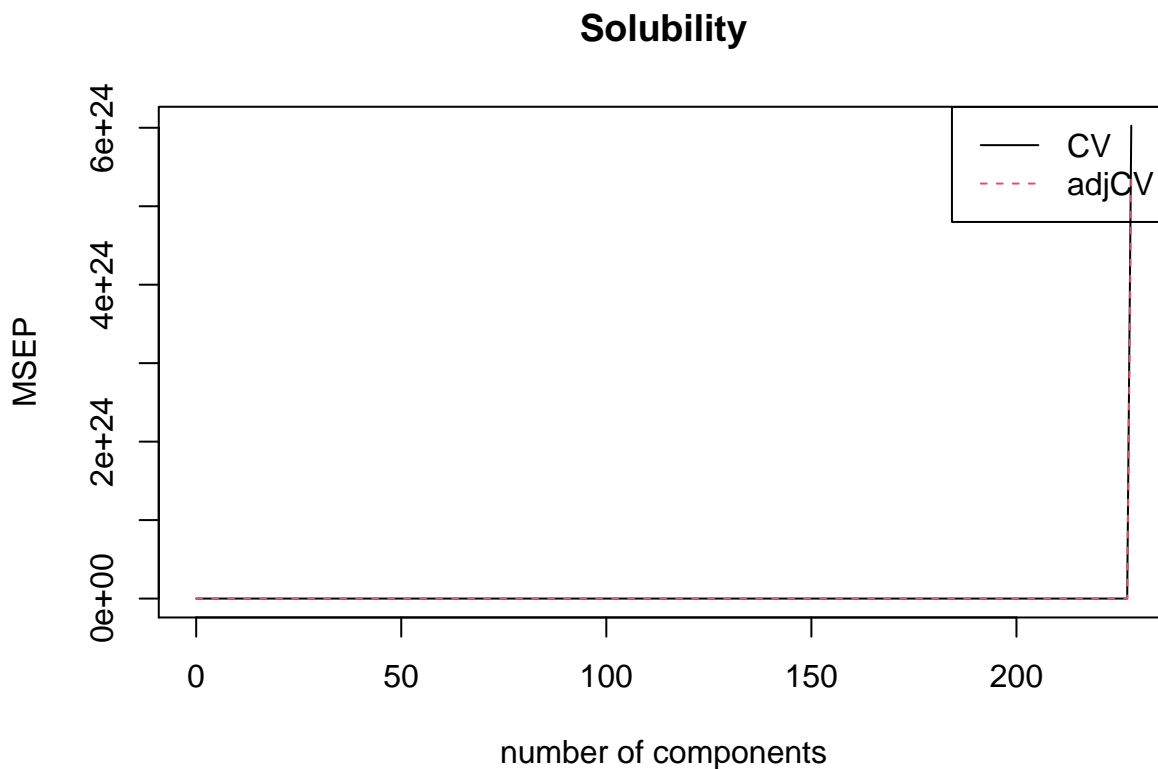
#find number of coefficient estimates
coef(lasso.fit$finalModel, lasso.fit$bestTune$lambda)
```

The chosen lambda is 0.00386592, the test error is 0.50333, and there are 149 non-zero coefficient estimates in the model.

## Part(d)

```
#fit model using function pcr()
set.seed(1)
pcr.mod <- pcr(Solubility ~ .,
data = train_data,
scale = TRUE,
validation= "CV")
summary(pcr.mod)

validationplot(pcr.mod, val.type="MSEP", legendpos = "topright")
```



```
cv.mse <- RMSEP(pcr.mod)
ncomp.cv <- which.min(cv.mse$val[1,,])-1
ncomp.cv

predy2.pcr2 <- predict(pcr.mod, newdata = y1, ncomp = ncomp.cv)
mean((y2 - predy2.pcr2)^2)

#fit model using caret
ctrl1 <- trainControl(method = "cv", selectionFunction = "best")
modelLookup("pcr")
set.seed(1)
pcr.fit <- train(x, y,
method = "pcr",
tuneGrid = data.frame(ncomp = 1:190),
```

```
trControl = ctrl1,  
preProcess = c("center", "scale"))
```

The chosen M is 152 and the test error is 0.5477905