# Introduction to Machine Learning - Lab 2

*Gustav Sternelöv*

*Thursday, October 29, 2015*

## Assignment 1

**1.1**

The function *ridgereg_nfoldCV* which performs ridge regression by using n-fold cross validation is implemented at the first step of assignment 1. The R-code used to create the function can be seen in the appendix at the end of the report.
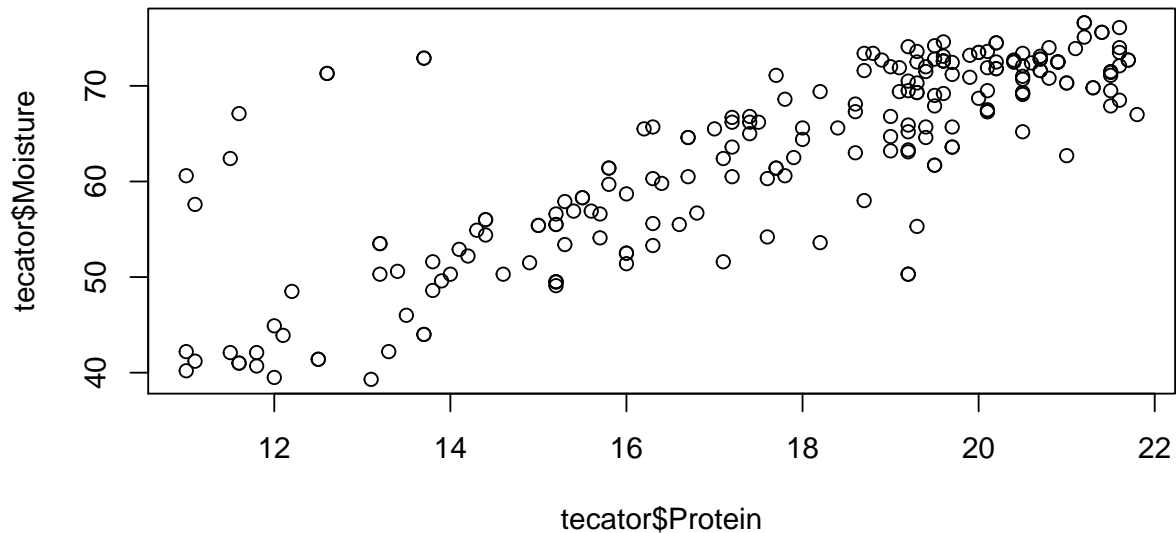
**1.2**

The function implemented in 1.2 is tested with the longley data set. The predictor variables has been centered, but not scaled. Number of folds is set to 10 and the value for lambda goes from 1 to 7, by 1. The given values for the CV score is shown in the following table below and it can be seen that the CV score increases for higher values of lambda. The conclusion from this result is that the best model is given with lambda set to 1.

```
## [1]  1220.437
## [1]  2921.322
## [1]  4750.991
## [1]  6569.879
## [1]  8316.805
## [1]  9966.241
## [1] 11510.01
```

## Assignment 2

In the second assignment the data material called *tecator* is used. This file contains information from a study where the fat content in samples of meat was investigated. Fat is the response variable and the predicitor variables in the data, observed for every meat in the sample, are 100 channel spectrum of absorbance records and the levels of moisture and protein.

**2.1**



To use a linear model to describe the levels of moisture with the levels of protein may be an appopriate approach by the look of the plot.
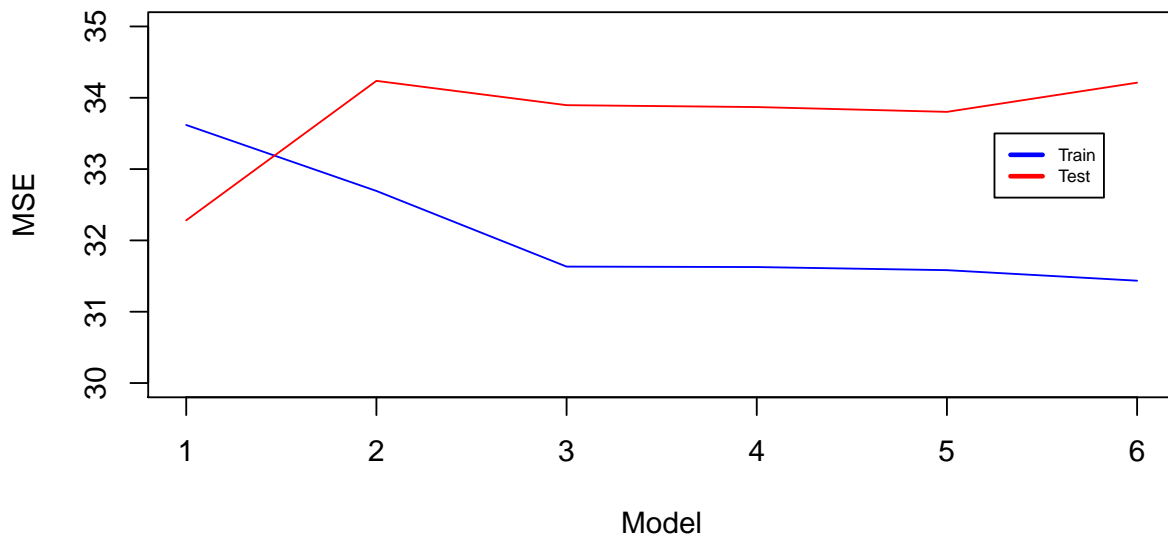
**2.2**

A Probabilistic model. . .

**2.3**

For exercise 2.3 the data set is divided into a training set and a validation set where each set contains 50 percent of the observations.
For models with polynomials of higher power the MSE for test data in general increases, and for training data it decreases. The best model for training data in terms of MSE is the, probably overfitted, model with polynomial terms up to a power of 6. This can be compared to the the test data which has the lowest MSE for a simple linear model. In terms of the bias-variance trade-off an interpretation of the plot is that the linear model gives the best fit. A conclusion that suits well with the comments made about the relationship between moisture and protein in 2.1.

**2.4**

Now the whole data set is used to construct the same models. AIC values is computed to evaluate the models and decide which model that is the best one. The criterion used for comparing the models is that the AIC value should be as low as possible. With that criterion in mind is it concluded that the model with polynomial terms up to a power of 3 is the preferred model.

```
## [1] 754.7710 756.4501 753.6657 755.3104 757.1727 759.1538
```
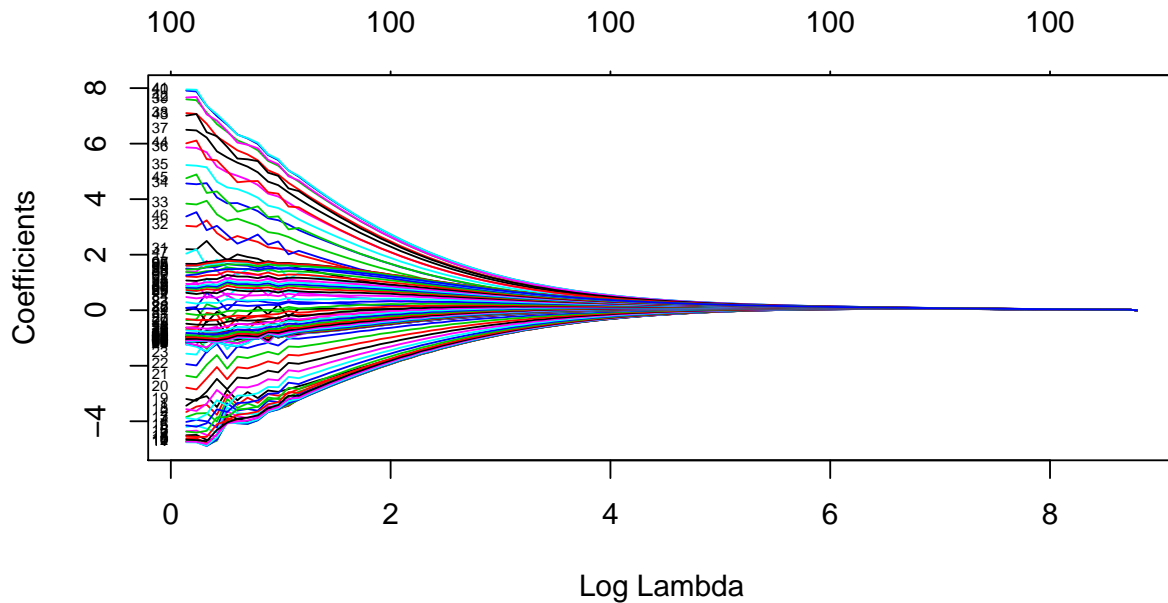
**2.5**

Another area where the AIC criterion can be useful is in the variable selection. Here the function *stepAIC* starts with a linear model that has *fat* as response variable and *channel1-channel100* as predictors. After running the variable selection procedure it is concluded that 63 of the predictors are selected.
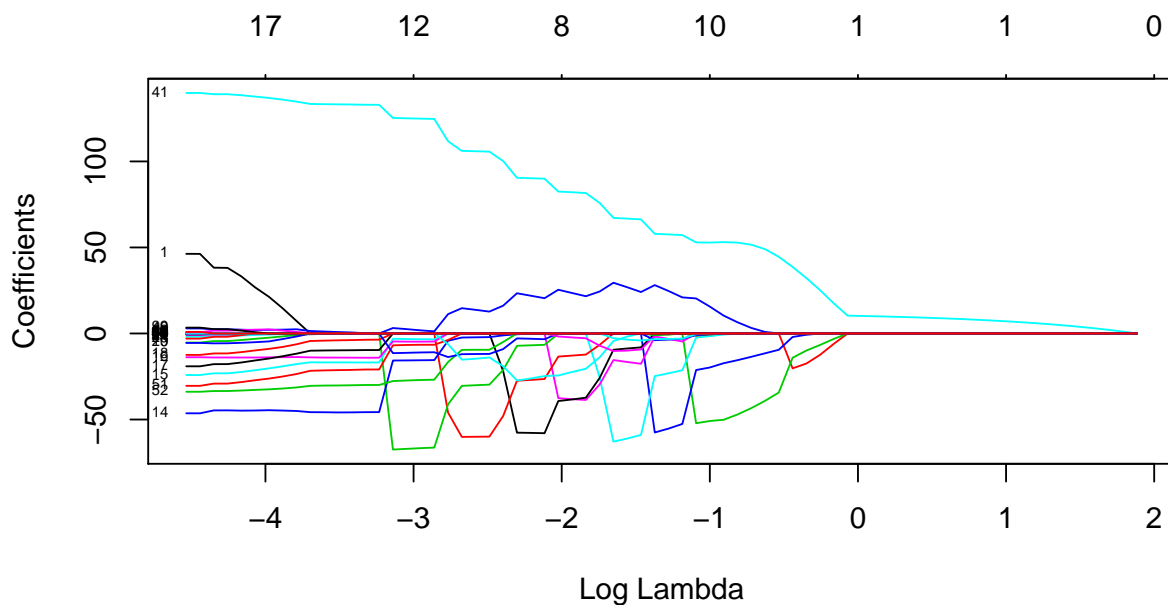
**2.6 - 2.7**

With the same response and predictor variables as in 2.5 a ridge regression model is fitted. The difference between a ridge regression model and a linear model is the use of a penalty factor $\lambda$ in the former model. The value of $\lambda$ affects the coefficients in such a way that they shrinks.
How the values of the coefficients in the ridge regression model depends on the log of $\lambda$ is illustrated by plotting these values.

When fitting a LASSO model instead the dependence between the coefficients and log of $\lambda$ changes. The difference between a Ridge regression model and a LASSO model is that LASSO shrinks the dimensionality by setting some features to zero. In practice LASSO works a bit like a variable selection procedure.
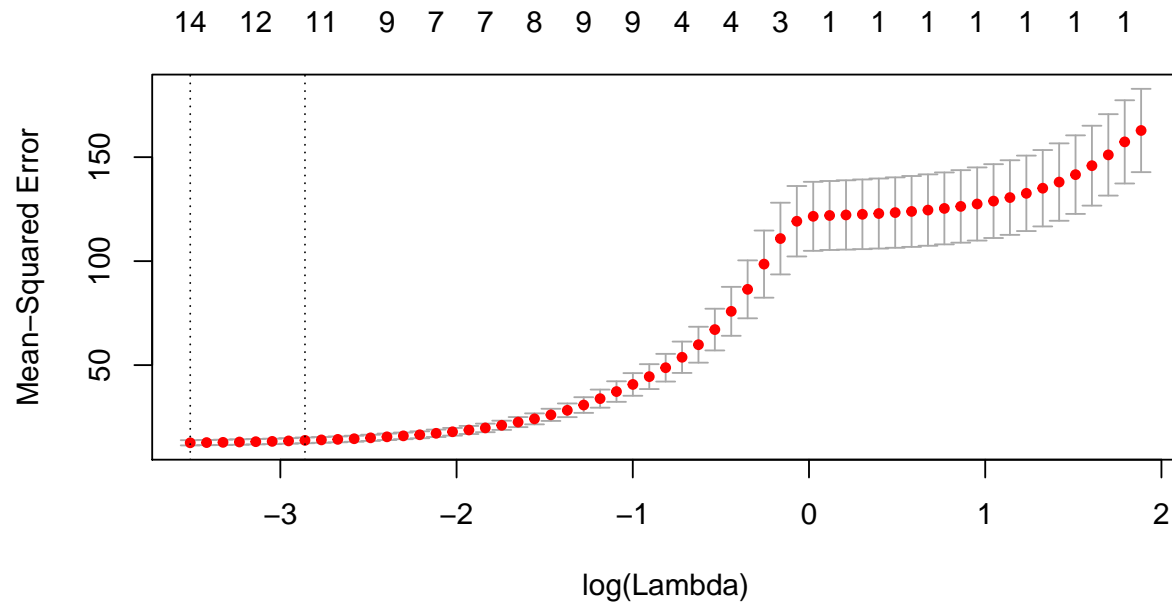


When comparing the plot for the ridge regression and for LASSO the characteristics for the respective methods becomes clear. For the ridge model the coefficients shrinks with higher values of lambda and slowly approaches zero. For the LASSO models coefficients it is harder to se a clear pattern in the dependence between the coefficients and log of $\lambda$. The majority of the coefficients quickly equals zero, some has a sort of wiggly curve and some of the coefficients decreases slowly for higher values of log $\lambda$.

Conclusions?

**2.8**

To find a optimal LASSO model a possible approach is to use cross-validation. This results in a model with 0.03 as the optimal value of lambda and 14 variables.
How the CV score depend on the log of $\lambda$ can be investigated with the following plot.



For higher values of the log $\lambda$ the CV-score increases. The optimal $\lambda$ was suggested to be approximately 0.03, but as can be seen in the graph it is not that clear which value of $\lambda$ that is optimal. The difference between the CV scores for log $\lambda$ values around -3 is very small, whilst the number of variables in the model differs. Therefore the reliability of the value chosen for the penalty factor, $\lambda$, may be questioned. A simpler model with for example 10 variables instead of 14 has its confidence interval for the CV score overlapping the confidence interval for the optimal models CV score. This might cause some uncertainty over the decision of the value for the penalty factor.

**2.9**

The variable selection done with *stepAIC* and the selection done with the cross-validated LASSO model choose a significantly different amount of variables. With *stepAIC* a model with 63 variables were obtained and from the cross-validated LASSO computations a model with 14 variables.

## Appendix - R-code

```
# Assignment 1
library(MASS)

# Reads in data set Longley
data(longley)
```

```r
data <- longley

longley.x <- data.matrix(longley[, 1:6])
longley.y <- longley[, "Employed"]

longley.x <- scale(longley.x, center=TRUE, scale=FALSE)

# Implement ridgereg function
ridgereg_nfoldCV <- function(x, y, lambda, nfolds){
  # Create the folds and initialize CV vector
  n <- length(y)
  seques <- floor(n/nfolds)
  reps <- n%%nfolds
  groups <- rep(seq(1,nfolds, 1), seques)
  end_values <- rep(nfolds, reps)
  folds <- c(groups, end_values)
  folds <- sort(folds)

  x <- cbind(rep(1, nrow(x)), x)
  CV <- 0
  for (i in 1:nfolds){
    testIndexes <- which(folds==i,arr.ind=TRUE)
    testData_x <- x[testIndexes, ]
    testData_y <- y[testIndexes]
    trainData_x <- x[-testIndexes, ]
    trainData_y <- y[-testIndexes]

    # Perform ridge regression on train data
    x_t <- t(trainData_x)
    I <- diag(ncol(trainData_x))
    BetaRidge <- solve(x_t %*% trainData_x + lambda * I) %*% x_t %*% trainData_y
    #y_hat <- trainData_x %*% BetaRidge
    # Test regression on test data and compare with true values for test data
    y_hat_test <- testData_x %*% BetaRidge
    # Calculates CV
    CV[i] <- sum((testData_y - y_hat_test)^2)
  }
  CV_score <- (1/nfolds) * sum(CV)
  return(CV_score)
}
for (i in 1:7){
  print(ridgereg_nfoldCV(longley.x, longley.y, lambda=i, nfolds=10))
}
# Assignment 2
# 2.1
tecator <- read.csv("Lab 2/tecator.csv", sep=";", header = TRUE)

plot(tecator$Moisture, tecator$Protein)

# By looking at the plot it does not seem completely of the charts to, at least to start with,
# consider a linear model.

# 2.2
```

```r
# How the Mi model looks up to the power n.
# lm(Moisture ~ poly(Protein, n, raw=TRUE))


# 2.3
# Divides data into train and test
n=dim(tecator)[1]
set.seed(12345)
id=sample(1:n, floor(n*0.5))
train=tecator[id,]
test=tecator[-id,]

# Extract data sets for moisture and protein data from train and test
Protein_train <- train$Protein
Moisture_train <- train$Moisture
Protein_test <- test$Protein
Moisture_test <- test$Moisture

poly_modelTrain <- list()
MSE_train <- 0
MSE_test <- 0

for (i in 1:6){
poly_modelTrain[[i]] <-  lm(Moisture_train ~ poly(Protein_train, i, raw=TRUE))
MSE_train[i] <- 1/length(Moisture_train) * sum(poly_modelTrain[[i]]$residuals^2)
y_hat_test <- predict(lm(Moisture_train ~ poly(Protein_train, i)), data.frame(Protein_train=Protein_test
MSE_test[i] <- 1/length(Moisture_test)  * sum((y_hat_test-Moisture_test)^2)
}


plot(1:6, MSE_train, type="l", ylim=c(30,35), col="blue", ylab="MSE", xlab="Model")
lines(1:6, MSE_test, type="l", col="red")
legend(5.25,33.5,c("Train","Test"), lty=c(1,1),
  lwd=c(2.5,2.5),col=c("blue","red"),  cex=0.6)

# 2.4
AIC <- 0
AIC2 <- 0
AIC3 <- 0
for (i in 1:6){
  poly_modelTrain[[i]] <-  lm(tecator$Moisture ~ poly(tecator$Protein, i, raw=TRUE))
  n <- length(tecator$Moisture)
  RSS <- sum(poly_modelTrain[[i]]$residuals^2)
  AIC[i] <- 2*(i+1) + n * log(RSS/n)
  AIC2[i] <- extractAIC(poly_modelTrain[[i]])[2]
  AIC3[i] <- AIC(poly_modelTrain[[i]])
}

# 2.5
vars <- data.frame(tecator[,2:102])

FullModel <- lm(Fat ~ ., data=vars)
Step_selec <- stepAIC(FullModel, direction = "both")
```

```r
# 2.6 - 2.7
library(glmnet)
# Create matrix with x variables
mat_vars <- as.matrix(vars[,1:100])
# y variable
mat_y <- as.matrix(vars[,101])

ridge_mod <- glmnet(mat_vars, mat_y, alpha=0, family = "gaussian")
plot(ridge_mod, xvar="lambda",label=TRUE)

lasso_mod <- glmnet(mat_vars, mat_y, alpha=1, family = "gaussian")
plot(lasso_mod, xvar="lambda",label=TRUE)

# 2.8
set.seed(12345)
lasso_cv <- cv.glmnet(mat_vars, mat_y, alpha=1, family = "gaussian")

lasso_cv$lambda.min
plot(lasso_cv)
coef(lasso_cv, s = "lambda.min")
```