

# Computer Assignment 10 - Linear Regression

## Machine Learning, Spring 2020

YOUR NAME

In this assignment, we will learn how to implement three important regression techniques - linear regression, ridge regression and the LASSO - in **R**. We'll apply these three modeling techniques to the preloaded **R** dataset *mtcars*.

### 1. Regression overview and loading the data

Linear regression is used to model a continuous response variable  $y$  as a linear combination of  $p$  predictors taking values  $x_1, \dots, x_p$ . Suppose that one observes  $n$  samples of  $y$  and associated predictors (in this case  $y, x_1, \dots, x_p$  are all  $n$  dimensional vectors). Define  $X := (x_1 x_2 \dots x_p)$  as the  $n \times p$  design matrix. The stochastic linear regression model of  $y$  on  $x_1, \dots, x_p$  is given by

$$(1): y = X\beta + \epsilon$$

where  $\epsilon = (\epsilon_1, \dots, \epsilon_p)^T$  is a vector of uncorrelated errors with mean 0 and variance 1 and  $\beta = (\beta_1, \dots, \beta_p)^T$  is the coefficient vector of unknown parameters. We typically assume a stronger condition that  $\epsilon_i \stackrel{iid}{\sim} N(0, 1)$  to simplify statistical inference on  $\beta_0, \dots, \beta_p$ . One should be careful when applying model (1) as there are many conditions that should be verified. We don't discuss these conditions here, though we recommend reading more about model selection for linear regression.

Recall from class that the least squares estimates  $\hat{\beta}$  is given by the *normal equations*:

$$(2): \hat{\beta} = (X^T X)^{-1} X^T y$$

As we can see from (2), the calculation of  $\hat{\beta}$  relies upon the invertibility of  $X^T X$ . Even if we assume  $p < n$ , we still require that there is no perfectly linear dependence between the predictor vectors  $x_1, \dots, x_p$ . In other words, we require the design matrix  $X$  to have full rank  $p$ . When  $rank(X) < p$ , then  $X$  suffers from *multicollinearity* in which case additional tools are needed for estimation of  $\beta$ . For example, penalization methods like ridge regression (squared penalization) or Lasso (L1 penalization) can be used to "shrink" the estimates of  $\beta$  towards the origin.

We will use the *mtcars* dataset available in **R** as an example throughout this assignment. This dataset describes various quantitative features of 32 different automobiles. There are 11 total variables in this dataset. We will study *miles per gallon (mpg)* as a function of four other predictors:

1. disp: displacement (cu.in.)
2. hp: gross horsepower
3. drat: rear axle ratio
4. wt: weight (lb/1000)

Load and parse the data with the following code:

```
#load the data
data(mtcars)
```

```
#create design matrix
X = data.frame(displ = mtcars$displ, hp = mtcars$hp, wt = mtcars$wt, drat = mtcars$drat)

#create response
Y = mtcars$mpg
```

### Questions:

1. To get an initial idea of pairwise relationships among the predictors, plot a grid of pairwise scatterplots on  $X$  using the `pairs()` command. Comment on the grid of pairwise scatterplots. Do any pair of predictors appear to share a linear relationship with one another?

YOUR CODE HERE

2. What if a pair of predictors have a perfect linear relationship (i.e their correlation is 1 or -1)? Can we still estimate  $\beta$  using non-penalized linear regression? In this case, is  $X^T X$  invertible? Why or why not?

YOUR ANSWER HERE

3. What can you say about the least squares estimates  $\hat{\beta}$  when the correlation between a pair of predictors gets close to 1 or -1? (HINT: think about the what happens to the empirical variance of  $X$  in this case. How does this affect the estimates  $\hat{\beta}$ ?)

YOUR ANSWER HERE

## 2. Linear Regression:

The `lm(y ~ x, data)` command can be used to run a linear regression of  $y$  on  $x$ . Here,  $y$  and  $x$  are both  $n$  dimensional vectors. The `data` argument is optional and specifies the source (a data frame) which  $x$  is contained. Once a linear regression has been run, we can use the `summary()` command to obtain coefficient estimates, standard errors of estimates, and p-values measuring the significance of each coefficient in the fitted model. Please type `?lm` for more details. Fit and summarize a linear model of `mpg` on the remaining variables using the following code:

```
linear.reg = lm(Y ~ displ + hp + wt + drat, data = X)
summary(linear.reg)
```

Now to demonstrate a situation of perfect collinearity, consider including a fifth covariate which is exactly twice the value of the `hp` variable. Construct a new design matrix and try fitting a linear model using the following code:

```
#construct a new design matrix
X.new = data.frame(X, two.hp = 2*X$hp)
#attempt to fit a linear model
linear.reg.fail = lm(Y ~ displ + hp + wt + drat + two.hp, data = X.new)
#summarize the regression
summary(linear.reg.fail)
```

### Questions:

1. What are the estimated coefficients on each predictor in the `linear.reg` model?
2. Which of these coefficients are statistically significant (at a 0.05 level)? Write the p-value for each of the estimated coefficients.

3. What did **R** do when we introduced perfect collinearity in *linear.reg.fail*? Since we know that  $X^T X$  is not invertible in this case, how did **R** still fit the model? (HINT: note that the coefficients on the original 4 variables remained the same as those in *linear.reg*.)

### 3.Ridge Regression:

The `lm.ridge(y ~ x, data, lambda)` command is used to run a ridge regression of  $y$  on  $x$  with ridge parameter  $lambda$ . Here,  $y$ ,  $x$  and  $data$  have the same interpretation as in the `lm()` function described in question 2. The `lm.ridge()` command is available in the *MASS* package in **R**, so be sure to load this package before use. Importantly, one must specify the ridge parameter  $\lambda$  for his/her/their choice of model. One principled way to choose  $\lambda$  is through cross validation.

For the moment, let's try a few fixed values of  $\lambda$ . First, fit a ridge regression for a fixed value of  $\lambda = 1$ . Also, calculate the distance  $\hat{\beta}$  is from the origin (a.k.a the magnitude of  $\hat{\beta}$ ) using the following code:

```
#run ridge regression with lambda = 1
ridge.reg.1 = lm.ridge(Y ~ disp + hp + wt + drat, data = X, lambda = 1)

#get a summary of the fit
coef.ridge.1 = coef(ridge.reg.1)

#calculate the magnitude of the estimated coefficient vector
dist.reg.1 = sum(abs(coef.ridge.1))
```

#### Questions:

1. Write the estimated coefficients for the fitted model with  $\lambda = 1$ . What is magnitude of  $\hat{\beta}$  in this model?
2. Fit a ridge regression with  $\lambda = 0$ . Write down the estimated coefficients and the magnitude of  $\hat{\beta}$ . How do the estimated coefficients here compare to those found with non-penalized linear regression? Explain why your observation makes sense.

YOUR CODE HERE

3. Fit a ridge regression with  $\lambda = 10, 50, 100$ , and  $1000$ . For each value, calculate the magnitude of the estimated coefficient vector. What happens to the magnitude of  $\hat{\beta}$  as you increase  $\lambda$ ? This is an example of the "shrinkage" effect of ridge regression.

YOUR CODE HERE

### 4.LASSO:

The `glmnet(X, y, lambda)` command can be used to fit the LASSO model of  $y$  on  $X$ . The *glmnet* package contains the functions required to conduct LASSO. Download this package before proceeding using the `install.packages()` and `library()` commands. Like ridge regression, the LASSO relies on the choice of a penalty parameter, (call it  $\lambda_1$ ). The `glmnet(X, y, lambda)` command is used to fit a LASSO regression with specified parameter  $lambda$ . In contrast to the `lm()` and `ridge.lm()` commands, the `glmnet(X, y, lambda)` command requires  $X$  to be an  $n \times p$  design matrix. As usual,  $y$  is the  $n$  dimensional vector of responses. Once again, we can choose the "best"  $\lambda_1$  using cross validation but we'll come back to this later. The `coef()` command is used to summarize the estimated coefficients of the model. Fit a LASSO model with  $\lambda_1 = 1$  to the *mtcars* data and calculate the magnitude of the estimated coefficients using the following commands:

```
#conduct a cross validation study to fit minimum MSE model
lasso.fit.1 = glmnet(as.matrix(X),Y,lambda = 1)
```

```
#summarize the estimated coefficients
coef.lasso.1 = coef(lasso.fit.1)

#calculate the magnitude of estimated coefficients
dist.lasso.1 = sum(abs(coef.lasso.1))
```

### Questions:

1. Write the estimated coefficients for the fitted model with  $\lambda_1 = 1$ . What is the magnitude of  $\hat{\beta}$  in this model?
2. Fit the LASSO with  $\lambda_1 = 0$ . Write down the estimated coefficients and the magnitude of  $\hat{\beta}$ . Do these estimates match (or are within a precision error close to) those found in the non-penalized linear regression?

YOUR CODE HERE

3. Fit the LASSO with  $\lambda_1 = 10, 50, 100$ , and 1000. For each value, calculate the magnitude of the estimated coefficient vector. What happens to the magnitude of  $\hat{\beta}$  as you increase  $\lambda$ ? Why do you think this is?

YOUR CODE HERE

4. In this assignment, we considered a dataset where  $n > p$ . If we had a situation where  $p > n$ , what modeling framework would you consider using to fit a linear regression? If you do not remove any of the variables, can we use non-penalized linear regression when  $p > n$ ?

## 5. High Dimensional LASSO:

We will repeat the analysis that was done in class. Start by uncommenting and running the following code to install the `BiocManager` and `bcellViper` packages. If you are having any difficulties installing these packages, you can seek further instruction [here](#).

```
# if (!requireNamespace("BiocManager", quietly = TRUE))
#   install.packages("BiocManager")
#
# BiocManager::install("bcellViper")
# install.packages("HDCI")
library(bcellViper)
data(bcellViper)
gene_expressions = data.frame(t(assayDataElement(dset, 'exprs')))
```

The installation of the above packages needs only to be done once, and can be re-commented after this initial run.

### Questions

1. Run an OLS model on the `gene_expressions` data, using `ADA` as your response variables, and the remaining variables as your predictors. Comment on any irregularities in the model output. Explain why you are not getting a reasonable number for your degrees of freedom in the model summary.

YOUR CODE, AND ANALYSIS, HERE

2. Now run 5 different LASSO models on this data, using any  $\lambda$  values of your choosing. Report how many non-zero  $\beta$  coefficients each model has. Plot the number of non-zero coefficients as a function of  $\lambda$ . Comment on any trend you observe.

YOUR CODE, AND ANALYSIS, HERE