

# Gov 2018: Lab 4 Cross Validation

Adeline Lo

Tuesday February 15, 2022

This lab on Ridge Regression and the Lasso in R is based off of p. 251-255 of “Introduction to Statistical Learning with Applications in R” by Gareth James, Daniela Witten, Trevor Hastie and Robert Tibshirani.

## CV Ridge Regression and the Lasso

```
rm(list=ls())
library(ISLR)
library(glmnet)
```

```
## Loading required package: Matrix
```

```
## Loaded glmnet 4.1-3
```

```
library(dplyr)
```

```
##
```

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

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

```
##
```

```
##      filter, lag
```

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

```
##
```

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

```
library(tidyr)
```

```
##
```

```
## Attaching package: 'tidyr'
```

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

```
##
```

```
##      expand, pack, unpack
```

```
#set seed
```

```
lab.seed<-202202
```

Use the `glmnet` package in order to perform ridge regression and the lasso. The main function in this package is `glmnet()`, which can be used to fit ridge regression models, lasso models, and more.

Load and remove NAs.

```
Hitters = na.omit(Hitters)
```

Execute a ridge regression and the lasso in order to predict **Salary** on the **Hitters** data.

Set up data:

```
x = model.matrix(Salary~., Hitters)[-1] # trim off the first column
                                         # leaving only the predictors
y = Hitters$Salary
```

The `model.matrix()` function is particularly useful for creating  $x$ ; not only does it produce a matrix corresponding to the 19 predictors but it also automatically transforms any qualitative variables into dummy variables. The latter property is important because `glmnet()` can only take numerical, quantitative inputs.

## Question 1. Ridge Regression

The `glmnet()` function has an `alpha` argument that determines what type of model is fit. If `alpha = 0` then a ridge regression model is fit, and if `alpha = 1` then a lasso model is fit. Fit a ridge regression model on  $x$  and  $y$  using the grid of lambda values from below.

```
grid = 10^seq(10, -2, length = 100)

ridge_mod <- glmnet(x, y, alpha = grid)
```

```
## Warning in if (alpha > 1) {: the condition has length > 1 and only the first
## element will be used
```

```
## Warning in glmnet(x, y, alpha = grid): alpha >1; set to 1
```

By default the `glmnet()` function performs ridge regression for an automatically selected range of  $\lambda$  values. However, here we have chosen to implement the function over a grid of values ranging from  $\lambda = 10^{10}$  to  $\lambda = 10^{-2}$ , essentially covering the full range of scenarios from the null model containing only the intercept, to the least squares fit.

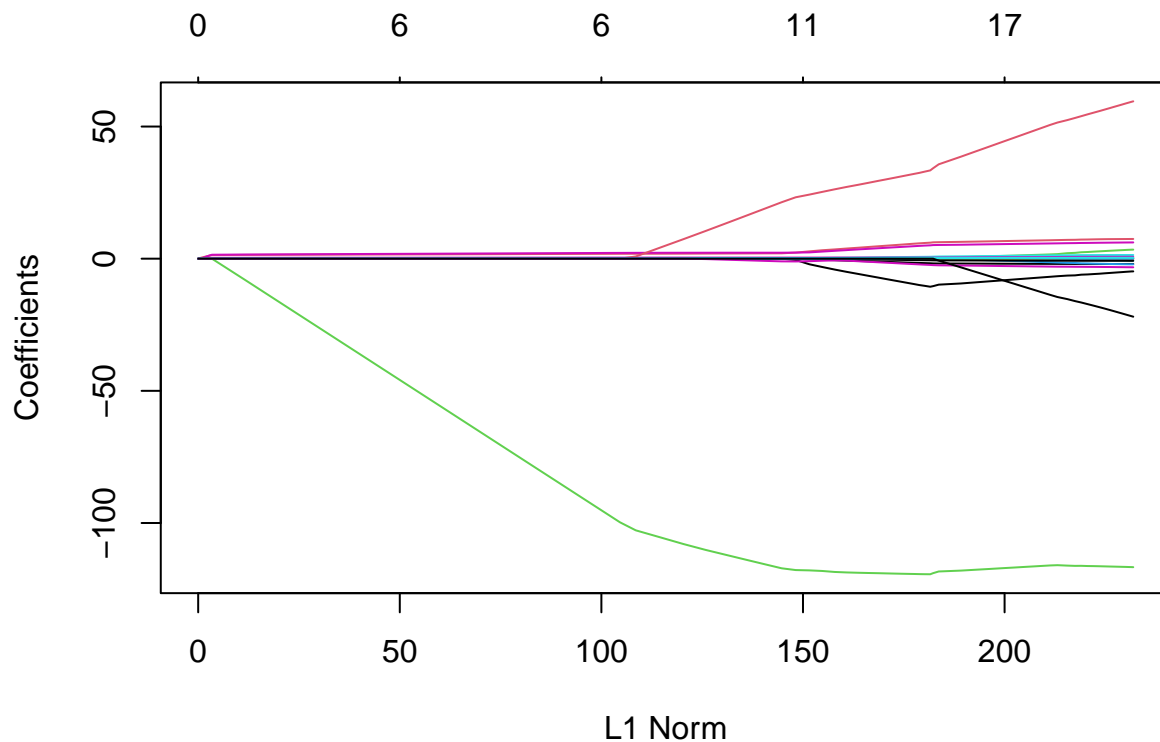
As we will see, we can also compute model fits for a particular value of  $\lambda$  that is not one of the original grid values. Note that by default, the `glmnet()` function standardizes the variables so that they are on the same scale. To turn off this default setting, use the argument `standardize = FALSE`.

Associated with each value of  $\lambda$  is a vector of ridge regression coefficients, stored in a matrix that can be accessed by `coef()`. In this case, it is a  $20 \times 100$  matrix, with 20 rows (one for each predictor, plus an intercept) and 100 columns (one for each value of  $\lambda$ ). Plot your coefficients from the ridge regression output with `plot(your.model.object)`.

```
dim(coef(ridge_mod))
```

```
## [1] 20 80
```

```
plot(ridge_mod)
```



We expect the coefficient estimates to be much smaller, in terms of  $l_2$  norm, when a large value of  $\lambda$  is used, as compared to when a small value of  $\lambda$  is used. Set  $\lambda$  to its 50th value. What are the coefficients at this value? What's their  $l_2$  norm (remove the intercept value)?

In contrast, what are the coefficients when  $\lambda$  is at its 60th value? Their  $l_2$  norm? Note the much larger  $l_2$  norm of the coefficients associated with this smaller value of  $\lambda$ .

Split the samples into a 80% training set and a 20% test set in order to estimate the test error of ridge regression and the lasso.

```
set.seed(lab.seed)
```

Next fit a ridge regression model on the training set, and evaluate its MSE (mean squared error) on the test set, using  $\lambda = 40$ . Note the use of the `predict()` function again: get predictions for a test set making sure to use `newx` argument.

```
set.seed(lab.seed)
```

Instead of arbitrarily choosing  $\lambda = 40$ , it would be better to use cross-validation to choose the tuning parameter  $\lambda$ . We can do this using the built-in cross-validation function, `cv.glmnet()`. By default, the function performs 10-fold cross-validation, though this can be changed using the argument `nfolds`. Set folds to 10 and calculate the  $\lambda$  that best minimizes the training MSE (`lambda.min` item in object returned from `cv.glmnet()`).

```
set.seed(lab.seed)
```

Plot the MSE as a function of  $\lambda$  by using `plot()` on our returned object from our call to `cv.glmnet`.

What is the test MSE associated with this value of  $\lambda$ ?

## Question 2. The Lasso

You just executed ridge regression with a wise choice of  $\lambda$ . Can lasso yield either a more accurate or a more interpretable model than ridge regression? Fit a lasso model, however, this time use the argument `alpha=1`.

Other than that change, proceed just as you did in fitting a ridge model. Plot your model object coefficients.

```
set.seed(lab.seed)
```

Now run the model again, this time performing cross-validation with folds equal to 10. Plot the model object. Then using the lambda that minimizes your cross validated training MSE, compute the associated testing MSE:

```
set.seed(lab.seed)
```

### Question 3. K-fold cross validation

Conduct K-fold cross validation for the ridge and lasso, with  $k \in \{5, 7, 9, 11, 13, 15\}$ . Assess the models at each  $k$  value by calculating the risk/prediction error in the test set and suggest the best  $k$  for each. Suggest the best final model.

```
set.seed(lab.seed)
```