

DS 303 HOMEWORK 5
DUE: OCT. 11, 2021 on Canvas by 11:59 pm (CT)

Instructions: Homework is to be submitted on Canvas by the deadline stated above. Please clearly print your name and student ID number on your HW.

Show your work (including calculations) to receive full credit. Please work hard to make your submission as readable as you possibly can - **this means no raw R output or code** (unless it is asked for specifically or needed for clarity).

Code should be submitted with your homework as a separate file (for example, a .R file, text file, Word file, or .Rmd are all acceptable). You should mark sections of the code that correspond to different homework problems using comments (e.g. ##### Problem 1 #####).

Problem 1: Concept Review

- a. Subset selection will produce a collection of $p+1$ models $M_0, M_1, M_2, \dots, M_p$. These represent the ‘best’ model of each size (where ‘best’ here is defined as the model with the smallest RSS). Is it true that the model identified as M_{k+1} must contain a subset of the predictors found in M_k ? In other words, is it true that if $M_1 : Y \sim X_1$, then M_2 must also contain X_1 . And if M_2 contains X_1 and X_2 , then M_3 must also contain X_1 and X_2 ? Explain your answer.
- b. Same question as part (a) but instead of subset selection, we now carry out forward stepwise selection.
- c. Suppose we perform subset, forward stepwise, and backward stepwise selection on a single data set. For each approach, again we can obtain $p+1$ models containing $0, 1, 2, \dots, p$ predictors. As we know, best subset will give us a best model with k predictors. Call this $M_{k,subset}$. Forward stepwise selection will give us a best model with k predictors. Call this $M_{k,forward}$. Backward stepwise selection will give us a best model with k predictors. Call this $M_{k,backward}$. Which of these three models has the smallest training MSE? Explain your answer. Hint: Consider the case for $k=0$ and $k=p$ first. Then the case for $k=1$. Then the case for $k=2, \dots, p-1$.
- d. Same setup as part (c). Which of these three models has the smallest test MSE? Explain your answer.

Problem 2: Simulation Studies

- a. Use the `rnorm()` function to generate a predictor X of length $n=100$, as well as error vector ϵ of length $n=100$. Assume that ϵ has variance 1.
- b. Generate a response vector Y of length $n=100$ according to the model:

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \epsilon,$$

where $\beta_0, \beta_1, \beta_2$, and β_3 are constants of your choice.

- c. Use best subset selection in order to choose the best model containing the predictors X, X^2, \dots, X^{10} . What is the best model obtained according to BIC and adjusted R^2 ? Report the coefficients of the best model obtained. Note: you will need to use the `data.frame()` function to create a single data set containing both X and Y .
- d. Repeat (c) using forward selection and also using backward selection. Report the best models obtained according to BIC and adjusted R^2 for both approaches. How does your answer compared to the results in part (c)?
- e. Now fit a lasso model to the simulated data, again using X, X^2, \dots, X^{10} as predictors. Use 10-fold cross-validation to select the optimal value of λ . Present a plot of the cross-validation error as a function of λ . Report the resulting coefficient estimates and discuss the results obtained.
- f. Now generate a response vector Y according to the model

$$Y = \beta_0 + \beta_7 X^7 + \epsilon,$$

and perform best subset selection and the lasso (again using predictors X, X^2, \dots, X^{10}). Discuss the results obtained.

Problem 3: Ridge Regression

For this problem, we will use the `College` data set in the `ISLR2` R package. Our aim is to predict the number of applications (`Apps`) received using the other variables in the dataset.

- a. Split the data set into a training and a test set. Please `set.seed(12)` so that we can all have the same results.
- b. Fit a least squares linear model (using all predictors) on the training set, and report the test MSE obtained.
- c. Fit a ridge regression model (using all predictors) on the training set. The function `glmnet`, by default, internally scales the predictor variables so that they will have standard deviation 1. Explain why this scaling is necessary when implementing regularized models.
- d. Find an optimal λ for the ridge regression model on the training set by using 5-fold cross-validation. Report the optimal λ here.
- e. Using that optimal λ , evaluate your trained ridge regression model on the test set. Report the test MSE obtained. Is there an improvement over the model from part (b)?
- f. Fit a lasso regression model on the training set. Find the optimal lambda using 5-fold cross-validation. Report the optimal λ and the test MSE obtained.
- g. Comment on your results. How accurately can we predict the number of college applications received? Is there much difference among the test errors resulting from these 3 approaches?

Problem 4: Regularized Regression Models

For this problem, we will continue with the `Hitters` example from lecture. Our aim is to predict the salary of baseball players based on their career statistics.

- a. We will start with a little data cleaning. We'll also split the data into a training and test set. So that we all get the same results, please use the following code:

```
library(ISLR2)
Hitters = na.omit(Hitters)
n = nrow(Hitters) #there are 263 observations
x = model.matrix(Salary ~ ., data=Hitters)[,-1] #19 predictors
Y = Hitters$Salary
set.seed(1)
train = sample(1:nrow(x), nrow(x)/2)
test=(-train)
Y.test = Y[test]
```

- b. Fit a ridge regression model. Replicate the example we had in class to obtain the optimal λ using 10-fold CV. Present a plot of the cross-validation error as a function of λ . Report that value here and call it $\lambda_{\min}^{\text{ridge}}$.
- c. Naturally, if we had taken a different training/test set or a different set of folds to carry out cross-validation, our optimal λ and therefore test error would change. An alternative is to select λ using the *one-standard error rule*. The idea is, instead of picking the λ that produces the smallest CV error, we pick the model whose CV error is within one standard error of the lowest point on the curve you produced in part (b). The intention is to produce a more **parimonious** model. The `glmnet` function does all of this hard work for you and we can extract the λ based on this rule using the following code: `cv.out$lambda.1se` (assuming your `cv.glmnet` object is named `cv.out`). Report your that λ here and call it $\lambda_{1se}^{\text{ridge}}$.
- d. Fit a lasso regression model. Replicate the example we had in class to obtain the optimal λ using 10-fold CV. Report that value here and call it $\lambda_{\min}^{\text{lasso}}$. Also report the optimal λ using the smallest standard error rule and called it $\lambda_{1se}^{\text{lasso}}$.
- e. You now have 4 values for the tuning parameter:

$$\lambda_{\min}^{\text{ridge}}, \lambda_{1se}^{\text{ridge}}, \lambda_{\min}^{\text{lasso}}, \lambda_{1se}^{\text{lasso}}.$$

Now evaluate the ridge regression models on your test set using $\lambda = \lambda_{\min}^{\text{ridge}}$ and $\lambda = \lambda_{1se}^{\text{ridge}}$. Evaluate the lasso models on your test set using $\lambda_{\min}^{\text{lasso}}$ and $\lambda_{1se}^{\text{lasso}}$. Compare the obtained test errors and report them here. Which model performs the best in terms of prediction? Do you have any intuition as to why?

- f. Report the coefficient estimates coming from ridge using $\lambda_{\min}^{\text{ridge}}$ and $\lambda_{1se}^{\text{ridge}}$ and likewise for the lasso models. How do the ridge regression estimates compare to those from the lasso? How do the coefficient estimates from using λ_{\min} compare to those from the one-standard error rule?
- g. If you were to make a recommendation to an upcoming baseball player who wants to make it big in the major leagues, what handful of features would you tell this player to focus on?