

Modern Data Mining - HW 2

Anirudh Bajaj

Esther Shin

Matt LeBaron

Overview / Instructions

This is homework #2 of STAT 471/571/701. It will be **due on Oct, 10, 2018 by 11:59 PM** on Canvas. You can directly edit this file to add your answers. Submit the Rmd file, a PDF or word or HTML version with only 1 submission per HW team.

Problem 0

Review the code and concepts covered during lecture: multiple regression, model selection and penalized regression through elastic net.

Problem 1

Do ISLR, page 262, problem 8 only part (a) to (d) and write up the answer here. This question is designed to help us understanding model selections through simulations. (e) Describe as accurate as possible what Cp and BIC are estimating?

Problem 2:

This will be the last part of the Auto data from ISLR. The original data contains 408 observations about cars. It has some similarity as the data CARS that we use in our lectures. To get the data, first install the package ISLR. The data Auto should be loaded automatically. We use this case to go through methods learnt so far.

You can access the necessary data with the following code:

```
# check if you have ISLR package, if not, install it
if(!requireNamespace('ISLR')) install.packages('ISLR')
auto_data <- ISLR::Auto
```

Final modelling question: we want to explore the effects of each feature as best as possible. You may explore interactions, feature transformations, higher order terms, or other strategies within reason. The model(s) should be as parsimonious (simple) as possible unless the gain in accuracy is significant from your point of view. Use Mallows's Cp or BIC to select the model. * Describe the final model and its accuracy. Include diagnostic plots with particular focus on the model residuals. * Summarize the effects found. * Predict the mpg of a car that is: built in 1983, in US, red, 180 inches long, 8 cylinders, 350 displacement, 260 as horsepower and weighs 4000 pounds. Give a 95% CI.

Problem 3: Lasso

Crime data continuation: We use a subset of the crime data discussed in class, but only look at Florida and California. crimedata is available on Canvas; we show the code to clean here.

```
crime <- read.csv("CrimeData.csv", stringsAsFactors = F, na.strings = c("?"))
crime <- dplyr::filter(crime, state %in% c("FL", "CA"))
```

Our goal is to find the factors which relate to violent crime. This variable is included in crime as crime\$violentcrimes.perpop.

Use LASSO to choose a reasonable, small model. Fit an OLS model with the variables obtained. The final model should only include variables with p-values < 0.05 . Note: you may choose to use lambda 1st or lambda min to answer the following questions where apply.

1. What is the model reported by LASSO?
2. What is the model after running OLS?
3. What is your final model, after excluding high p-value variables?

Now, instead of Lasso, we want to consider how changing the value of alpha (i.e. mixing between Lasso and Ridge) will affect the model. Cross-validate between alpha and lambda, instead of just lambda. Note that the final model may have variables with p-values higher than 0.05; this is because we are optimizing for accuracy rather than parsimoniousness.

1. What is your final elastic net model? What were the alpha and lambda values? What is the prediction error?
2. Use the elastic net variables in an OLS model. What is the equation, and what is the prediction error.
3. Summarize your findings, with particular focus on the difference between the two equations.