

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?

(a)

```
# Use rnorm() to generate a predictor X of length n = 100, and a noise vector of length n = 100
x <- rnorm(100)
noise <- rnorm(100)
```

(b)

```
# Generate a response vector Y of length n = 100 according to the model
y <- 1 + 2*x + 3*x^2 + 4*x^3 + noise
```

(c)

```
# Create the predictors x^2 through x^10
x2 <- x^2
x3 <- x^3
x4 <- x^4
x5 <- x^5
x6 <- x^6
x7 <- x^7
x8 <- x^8
x9 <- x^9
x10 <- x^10

# Use the regsubsets() function to perform best subset selection
new_data <- data.frame(x,x2,x3,x4,x5,x6,x7,x8,x9,x10,noise,y)
new_subset <- regsubsets(y ~ x+x2+x3+x4+x5+x6+x7+x8+x9+x10,data = new_data, nvmax = 10)
new_summary <- summary(new_subset)
new_summary$which
```

```
##      (Intercept)      x      x2      x3      x4      x5      x6      x7      x8      x9      x10
## 1             TRUE FALSE FALSE TRUE  FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## 2             TRUE FALSE  TRUE  TRUE  FALSE FALSE FALSE FALSE FALSE FALSE FALSE
```

```
## 3      TRUE TRUE TRUE TRUE FALSE FALSE FALSE FALSE FALSE FALSE
## 4      TRUE TRUE TRUE TRUE FALSE FALSE FALSE TRUE FALSE FALSE
## 5      TRUE TRUE TRUE TRUE TRUE FALSE TRUE FALSE FALSE FALSE
## 6      TRUE TRUE TRUE TRUE TRUE FALSE TRUE FALSE TRUE FALSE
## 7      TRUE TRUE TRUE TRUE TRUE FALSE TRUE FALSE TRUE TRUE
## 8      TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE FALSE
## 9      TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## 10     TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE

data.frame(variables=(1:length(new_summary$rsq)), r_squared=new_summary$rsq)
```

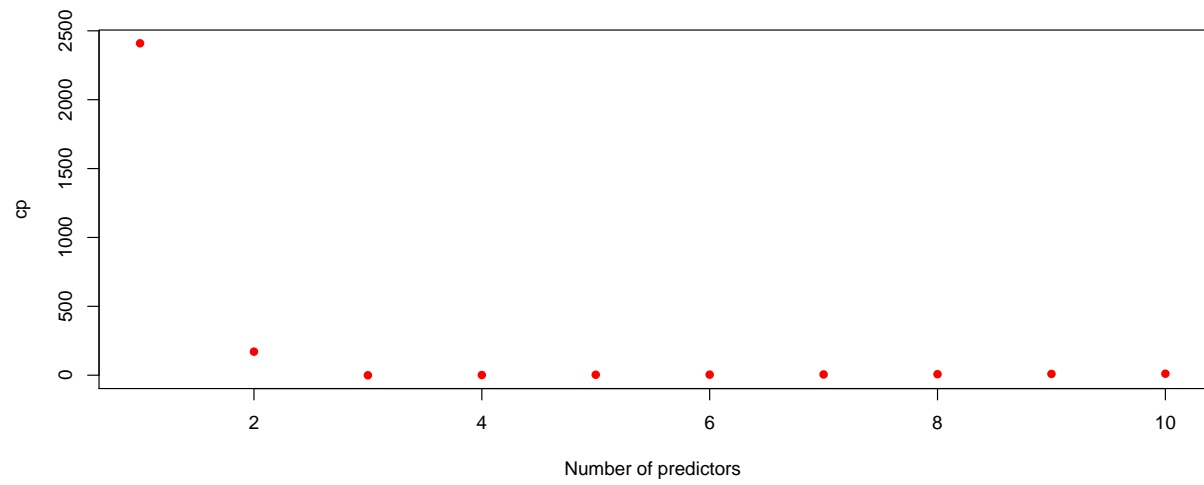
```
##      variables r_squared
## 1           1 0.9234130
## 2           2 0.9918994
## 3           3 0.9972016
## 4           4 0.9972040
## 5           5 0.9972072
## 6           6 0.9972490
## 7           7 0.9972625
## 8           8 0.9972679
## 9           9 0.9972740
## 10          10 0.9972794
```

```
# What is the best model obtained?
data.frame(variables = (1:length(new_summary$rsq)),
            r_squared = new_summary$rsq,
            rss = new_summary$rss,
            bic = new_summary$bic,
            cp = new_summary$cp)
```

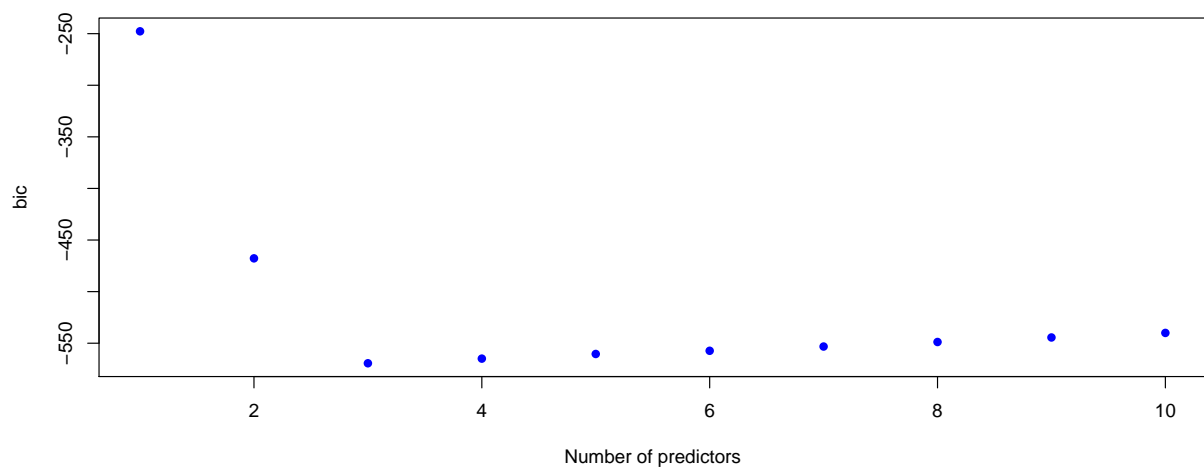
```
##      variables r_squared      rss      bic      cp
## 1           1 0.9234130 2393.35184 -247.7225 2409.3733748
## 2           2 0.9918994  253.14385 -467.7666  170.9923174
## 3           3 0.9972016   87.45188 -569.4483   -0.4549067
## 4           4 0.9972040   87.37544 -564.9306    1.4650777
## 5           5 0.9972072   87.27596 -560.4393    3.3609395
## 6           6 0.9972490   85.96853 -557.3435    3.9923170
## 7           7 0.9972625   85.54607 -553.2310    5.5500748
## 8           8 0.9972679   85.37705 -548.8236    7.3731487
## 9           9 0.9972740   85.18836 -544.4397    9.1756283
## 10          10 0.9972794   85.02059 -540.0317   11.0000000
```

The best model obtained has three variables: x, x2, and x3. This makes sense because these three variables were used to generate Y in the first place, so they should be most predictive, although they aren't perfect because we added in noise as well.

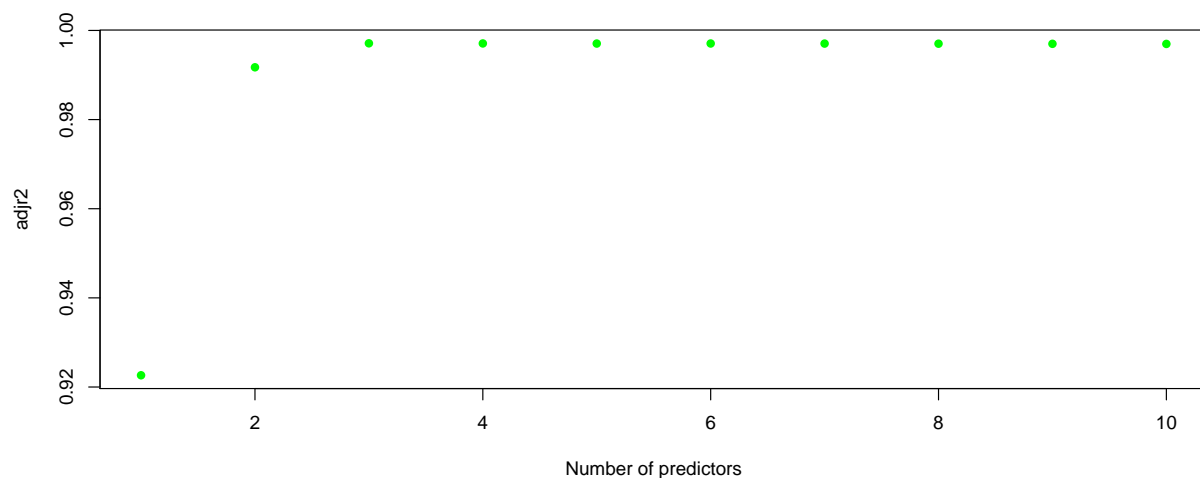
```
# Comparing the plots for Cp, bic, and r2
plot(new_summary$cp, xlab="Number of predictors",
     ylab="cp", col="red", type="p", pch=16)
```



```
plot(new_summary$bic, xlab="Number of predictors",  
     ylab="bic", col="blue", type="p", pch=16)
```



```
plot(new_summary$adjr2, xlab="Number of predictors",  
     ylab="adjr2", col="green", type="p", pch=16)
```



```
# Find the optimal model size for Cp
opt.size <- which.min(new_summary$cp)
opt.size
```

```
## [1] 3
```

```
# Find the optimal model size for bic
opt.size <- which.min(new_summary$bic)
opt.size
```

```
## [1] 3
```

```
# Find the optimal model size for r2
opt.size <- which.min(new_summary$rsq)
opt.size
```

```
## [1] 1
```

```
# Report the coefficients of the best model obtained
coef(new_summary, 3)
```

```
## NULL
```

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.