

Name - NetID

## Homework 12

### Homework Instructions

For questions that require code, please create a code chunk directly below the question and type your code there. Your knitted pdf will show both your code and your output. You are encouraged to knit your file as you work to check that your coding and formatting is done so appropriately.

For written responses or multiple choice questions, please bold your (selected) answer.

### Grading Details

All questions will be graded full credit (1 point), half credit (0.5 point) or no credit (0 points).

Full credit responses should have the correct response and appropriate code (if applicable). Half credit responses will have a reasonable attempt (typically no more than one small error or oversight), and no credit responses will be either non-attempts or attempts with significant errors.

#### Exercise 1

For exercises 1 - 8, use the the built-in R dataset `mtcars`. Use `mpg` as the response variable. Do not modify any of the data. (An argument could be made for `cyl`, `gear`, and `carb` to be coerced to factors, but for simplicity, we will keep them numeric.)

```
# starter
mtcars
```

Fit an additive linear model with all available variables as predictors. Call this model `full_mod`. Which predictor seems to be most explained by the other predictors in the model?

*Hint: We learned a function in chapter 15 that can help here!*

```
# solution
full_mod = lm(mpg ~ ., data = mtcars)
library(faraway)
vif(full_mod)
```

```
##      cyl      disp      hp      drat      wt      qsec      vs      am
## 15.373833 21.620241  9.832037  3.374620 15.164887  7.527958  4.965873  4.648487
##      gear      carb
##  5.357452  7.908747
```

`disp` has the largest vif.

#### Exercise 2

Create a function that takes a model name as input and calculates the LOOCV-RMSE of that model.

Once created, calculate the LOOCV-RMSE of the model fit in Exercise 1.

```
# solution
full_mod = lm(mpg ~ ., data = mtcars)

calc_loocv_rmse = function(model) {
  sqrt(mean((resid(model) / (1 - hatvalues(model))) ^ 2))
}

calc_loocv_rmse(full_mod)
```

```
## [1] 3.490209
```

### Exercise 3

Start with the full model, and then perform variable selection using backwards AIC. Call this model `selected`. Then print the model coefficients of the selected model.

*Tip: use `trace = 0` as an argument to suppress lengthy output*

```
# solution
full_mod = lm(mpg ~ ., data = mtcars)
selected = step(full_mod, trace = FALSE)
coef(selected)
```

```
## (Intercept)          wt          qsec          am
##      9.617781     -3.916504      1.225886      2.935837
```

### Exercise 4

What is the LOOCV-RMSE of the “selected” model?

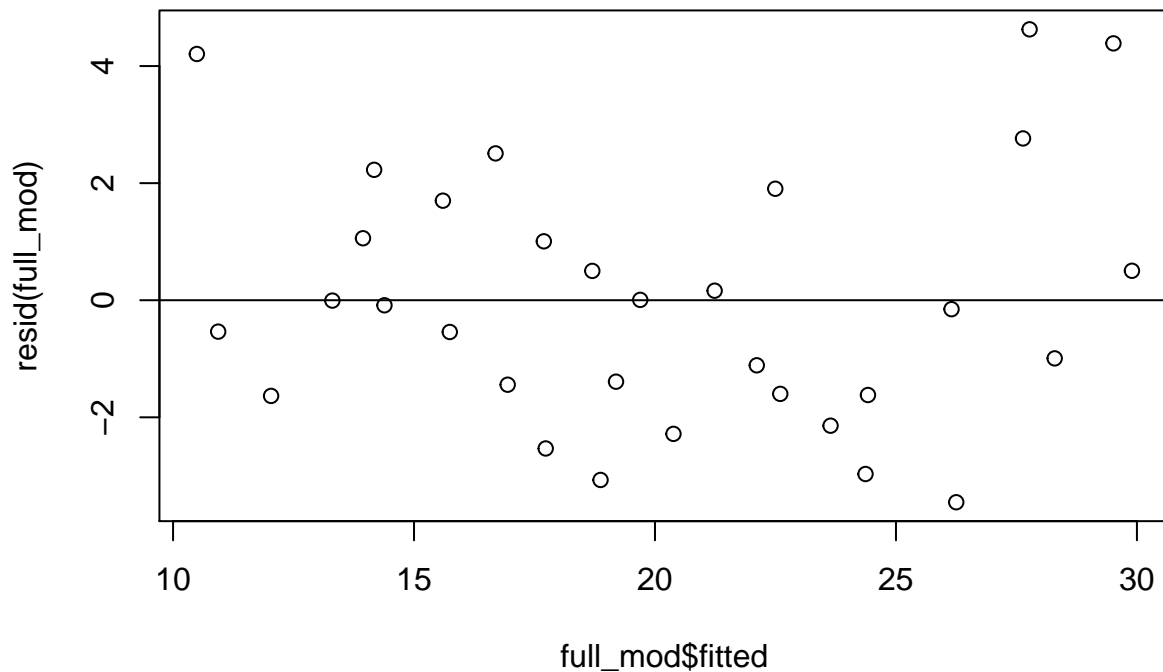
```
# solution
calc_loocv_rmse(selected)
```

```
## [1] 2.688538
```

### Exercise 5

Create a residual plot for the **full model**, and also run a Breusch-Pagan test and Shapiro-Wilk test.

```
# solution
plot(full_mod$fitted, resid(full_mod))
abline(h = 0)
```



```
library(lmtest)
```

```
## Warning: package 'lmtest' was built under R version 4.0.4
```

```
## Loading required package: zoo
```

```
## Warning: package 'zoo' was built under R version 4.0.4
```

```
##
```

```
## Attaching package: 'zoo'
```

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

```
##
```

```
## as.Date, as.Date.numeric
```

```
bptest(full_mod)
```

```
##
```

```
## studentized Breusch-Pagan test
```

```
##
```

```
## data: full_mod
```

```
## BP = 14.914, df = 10, p-value = 0.1352
```

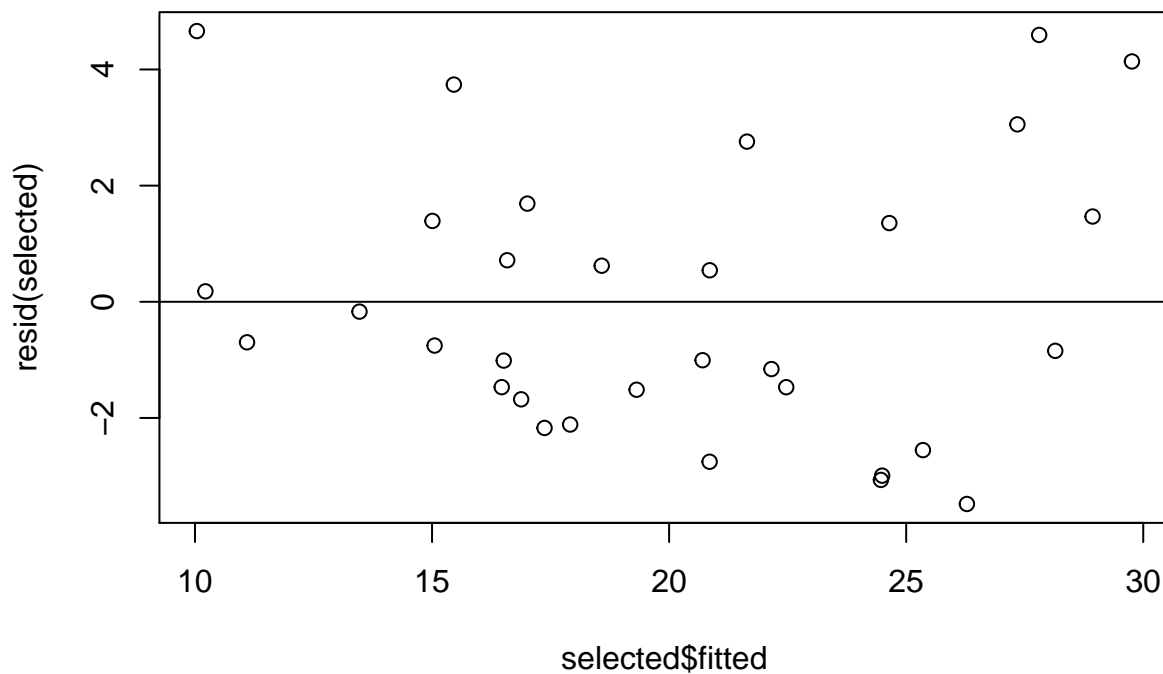
```
shapiro.test(resid(full_mod))
```

```
##  
## Shapiro-Wilk normality test  
##  
## data:  resid(full_mod)  
## W = 0.95694, p-value = 0.2261
```

### Exercise 6

Create a residual plot for the **selected model**, and also run a Breusch-Pagan test and Shapiro-Wilk test.

```
# solution  
plot(selected$fitted, resid(selected))  
abline(h = 0)
```



```
library(lmtest)  
bptest(selected)
```

```
##  
## studentized Breusch-Pagan test  
##  
## data:  selected  
## BP = 6.1871, df = 3, p-value = 0.1029
```

```
shapiro.test(resid(selected))
```

```
##  
## Shapiro-Wilk normality test  
##  
## data: resid(selected)  
## W = 0.9411, p-value = 0.08043
```

### Exercise 7

Based on the previous exercises, which model is probably better for predicting?

- **The selected model: It has better accuracy, and the model diagnostics are reasonable**
- The selected model: The model diagnostics for the full model are problematic and will lead to unreliable estimates
- The full model: It has better accuracy, and the model diagnostics are reasonable
- The full model: The model diagnostics for the selected model are problematic and will lead to unreliable estimates

### Exercise 8

```
# starter  
LifeCycleSavings
```

For exercises 8 - 12, use the built-in R dataset `LifeCycleSavings`. Use `sr` as the response variable.

Fit a model with all available predictors as well as their two-way interactions. What is the Adjusted  $R^2$  of this model?

*Hint: You can use the “ $\wedge 2$ ” notation to do this concisely—see 16.3 from the book for an example.*

```
# solution  
full_mod2 = lm(sr ~ . ^ 2, data = LifeCycleSavings)  
summary(full_mod2)$adj
```

```
## [1] 0.261233
```

### Exercise 9

Start with the model fit in Exercise 8, then perform variable selection using backwards **BIC**. Print the coefficients of the selected model.

```
# solution  
selected_BIC = step(full_mod2, k = log(nrow(LifeCycleSavings)), trace = FALSE)  
coef(selected_BIC)
```

```
## (Intercept)      pop15      dpi      ddp      dpi:ddpi  
## 16.5287997486 -0.2023668518 -0.0027411096 0.0462478987 0.0008170652
```

## Exercise 10

Start with the model fit in Exercise 8, then perform variable selection using backwards AIC. Print the coefficients of the selected model.

```
# solution
full_mod2 = lm(sr ~ . ^ 2, data = LifeCycleSavings)
selected_AIC = step(full_mod2, trace = FALSE)
coef(selected_AIC)
```

```
##      (Intercept)      pop15      dpi      ddpi      dpi:ddpi
## 16.5287997486 -0.2023668518 -0.0027411096  0.0462478987  0.0008170652
```

*If done correctly, you'll find that AIC and BIC in this case actually produce the same model! It doesn't always happen that way, but sometimes it does.*

## Exercise 11

Consider the full model in Exercise 8 and the BIC model (or AIC) model fit next. Based on LOOCV-RMSE, which of these models is likely more accurate at prediction?

```
# solution
calc_loocv_rmse(full_mod2)
```

```
## [1] 4.338482
```

```
calc_loocv_rmse(selected_BIC)
```

```
## [1] 3.833628
```

```
summary(selected_BIC)$adj
```

```
## [1] 0.3188961
```

- The full model is likely more accurate at prediction
- **The selected model is likely more accurate at prediction**

## Exercise 12

Which **one** of these statements best distinguishes AIC vs. BIC for model selection? Bold your answer.

- AIC will generally favor more interaction terms and fewer solo predictors
- BIC will generally favor more interaction terms and fewer solo predictors
- AIC will generally favor smaller models with fewer parameters
- **BIC will generally favor smaller models with fewer parameters**