

Homework 5 Solutions

Instructions

Due: 5:00pm on Wednesday, October 13th

1. Add your name between the quotation marks on the author line in the YAML above.
2. Compose your answer to each problem between the bars of red stars.
3. Commit your changes frequently.
4. Be sure to knit your .Rmd to a .pdf file.
5. Push both the .pdf and the .Rmd file to your repo on the class organization before the deadline.

Regression Competition II

Your objective is to revisit the data set from Homework 3 to build a new multiple regression model based on the feature selection techniques discussed recently in class.

As before, you will construct your model using a training data set with information on 66 variables recorded for 1808 houses. I've held back the data on 600 other houses which will serve as the test data set for assessing the predictive accuracy of your model.

You should record your answers in this .Rmd file. However, you are encouraged to use a separate .Rmd file for scratchwork. The assignment is divided into several **Components** to help organize your work. Put all work you want graded between the bars of red stars in the corresponding section.

Grading

This assignment restricts you to just using feature selection to build your model, so your overall score on this assignment will not be based on model accuracy (most models will have relatively similar accuracy). However, you will have optional opportunity at the end of this assignment to build a better model that synthesizes feature selection along with the other work you did on Homework 3, and I will run that model on test data as well and report the results.

The Data

The data set `house_train` can be found in the `hw_3` repo and can be loaded by running the following code.

```
house<-read_csv("house_train.csv")
```

Additionally, the `data_description.txt` file in the same repo gives a full description of the variables appearing in the data set.

There is one special column of note:

- `Sale_Price` is your response variable and should not be included as a predictor.

Components

Data Partitioning

In this section, create a training / validation split. We'll use the training set for model building, and the validation set for model assessment. (Ideally, we would not build **any** models using data from the validation set, but in this case, we've already peeked at the data in Homework 3.)

```
set.seed(10)
library(rsample)

## Warning: package 'rsample' was built under R version 3.6.2

house_part <- initial_split(house)
house_trn <- training(house_part)
house_tst <- testing(house_part)
```

Data Exploration

In this section, compute the correlations for all pairs of **quantitative** predictors. Which predictors appear to be most highly correlated? Create scatterplots for comparing these pairs of predictors. Explain what affect including highly correlated variables in the model would have on model accuracy (Consider the bias-variance tradeoff).

We see that the following pairs of variables are highly correlated:

- Garage_Cars and Garage_Area
- Gr_Liv_Area and TotRms_AbvGrd
- Total_Bsmt_SF and First_Flr_SF
- Second_Flr_SF and Gr_Liv_Area

Including many highly correlated predictors in MLR leads to model instability, and in particular, produces coefficient estimates that are highly variable from training set to training set, which ultimately leads to higher MSE.

```
house_trn %>% select(is.numeric) %>% #Selects just numeric variables
cor() %>% #Computes pairwise correlations
as.data.frame() %>% #Converts matrix to data frame
rownames_to_column("variable") %>% #uses row names as first column
pivot_longer(!variable) %>% #Converts to tidy frame with 1 column for correlation
filter(name != variable) %>% #Removes correlations between variable and itself
mutate(R2 = value^2) %>% #Computes R^2, which is always between 0 and 1
arrange(desc(R2)) #sorts in decreasing order of correlation strength
```

```
## Warning: Predicate functions must be wrapped in `where()`.
##
## # Bad
## data %>% select(is.numeric)
##
## # Good
## data %>% select(where(is.numeric))
##
## i Please update your code.
## This message is displayed once per session.
```

```
## # A tibble: 1,122 x 4
##   variable      name      value      R2
##   <chr>         <chr>    <dbl> <dbl>
## 1 Garage_Cars   Garage_Area 0.876 0.767
## 2 Garage_Area   Garage_Cars 0.876 0.767
## 3 Gr_Liv_Area   TotRms_AbvGrd 0.806 0.649
## 4 TotRms_AbvGrd Gr_Liv_Area 0.806 0.649
## 5 Total_Bsmt_SF First_Flr_SF 0.770 0.592
## 6 First_Flr_SF  Total_Bsmt_SF 0.770 0.592
## 7 Gr_Liv_Area   Sale_Price   0.751 0.564
## 8 Sale_Price    Gr_Liv_Area 0.751 0.564
## 9 Second_Flr_SF Gr_Liv_Area 0.707 0.499
## 10 Gr_Liv_Area   Second_Flr_SF 0.707 0.499
## # ... with 1,112 more rows
```

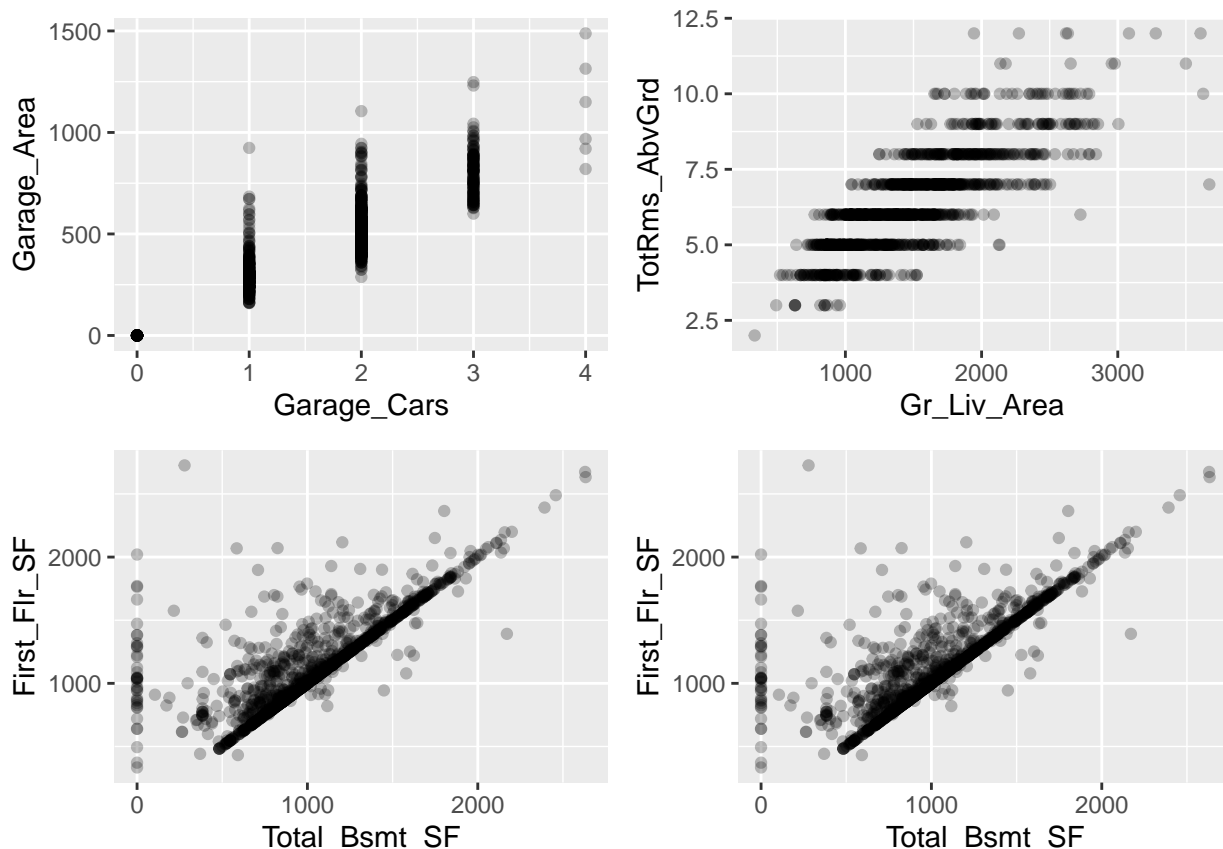
```
library(gridExtra)
```

```
##
## Attaching package: 'gridExtra'

## The following object is masked from 'package:dplyr':
##
##   combine

g1<- house_trn %>% ggplot(aes(x = Garage_Cars, y= Garage_Area))+geom_point(alpha = .25)
g2<- house_trn %>% ggplot(aes(x = Gr_Liv_Area, y= TotRms_AbvGrd))+geom_point(alpha = .25)
g3<- house_trn %>% ggplot(aes(x = Total_Bsmt_SF, y= First_Flr_SF))+geom_point(alpha = .25)
g4<- house_trn %>% ggplot(aes(x = Second_Flr_SF, y= Gr_Liv_Area))+geom_point(alpha = .25)

grid.arrange(g1,g2,g3,g3, nrow =2)
```



Model Building

In this section, you should perform one of the following algorithms: best-subset, forward-selection, or backwards-elimination, using the `regsubsets` package. Briefly explain why you choose the algorithm you did. Do not perform any data processing, or include any transformations or interaction terms (i.e just do feature selection via `regsubsets`)

We restrict our attention just to quantitative predictors, for two reasons:

- 1) `regsubsets` treats each level of a categorical variable as a separate variable, so adding in the 32 categorical variables actually adds a very large number of predictors, and even forward selection can consume too much memory.
- 2) Similarly, `regsubsets` looks at models that don't use all levels of a categorical variable. However, if a final model includes 1 level, it should include all, which can lead to inaccurate error measurements from adjusted R^2 , C_p and BIC .

Since the full model performed well in the initial trial, we will use backwards elimination, which tends to favor models with more variables, compared to forward selection. We will not use best subset, given the large number of predictors.

```
library(leaps)
house_trn_num <- house_trn %>% select(is.numeric)

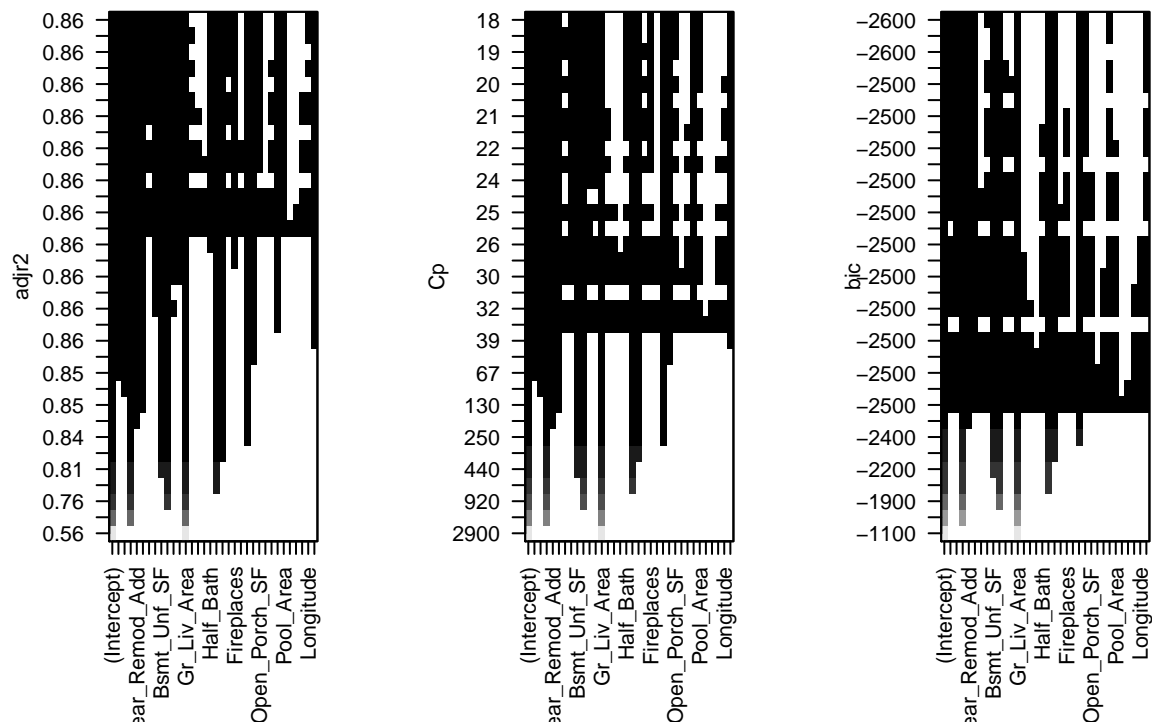
back_elim <- regsubsets(Sale_Price ~., data= house_trn_num, really.big=T, method = "backward", nvmax = 32)
```

Model Diagnostics

In this section, compare the results of your models by computing appropriate metrics on training data for models of **each** predictor size. Display these metrics both graphically, and then explicitly compute optimal values. Which model size appears to be most accurate?

Below are the predictor inclusion plot, ranked by size of each metric. Not surprisingly, we see that BIC tends to favor models with fewer variables than either Cp or Adjusted R^2 .

```
par(mfrow=c(1,3))
plot(back_elim, scale = "adjr2")
plot(back_elim, scale = "Cp")
plot(back_elim, scale = "bic")
```

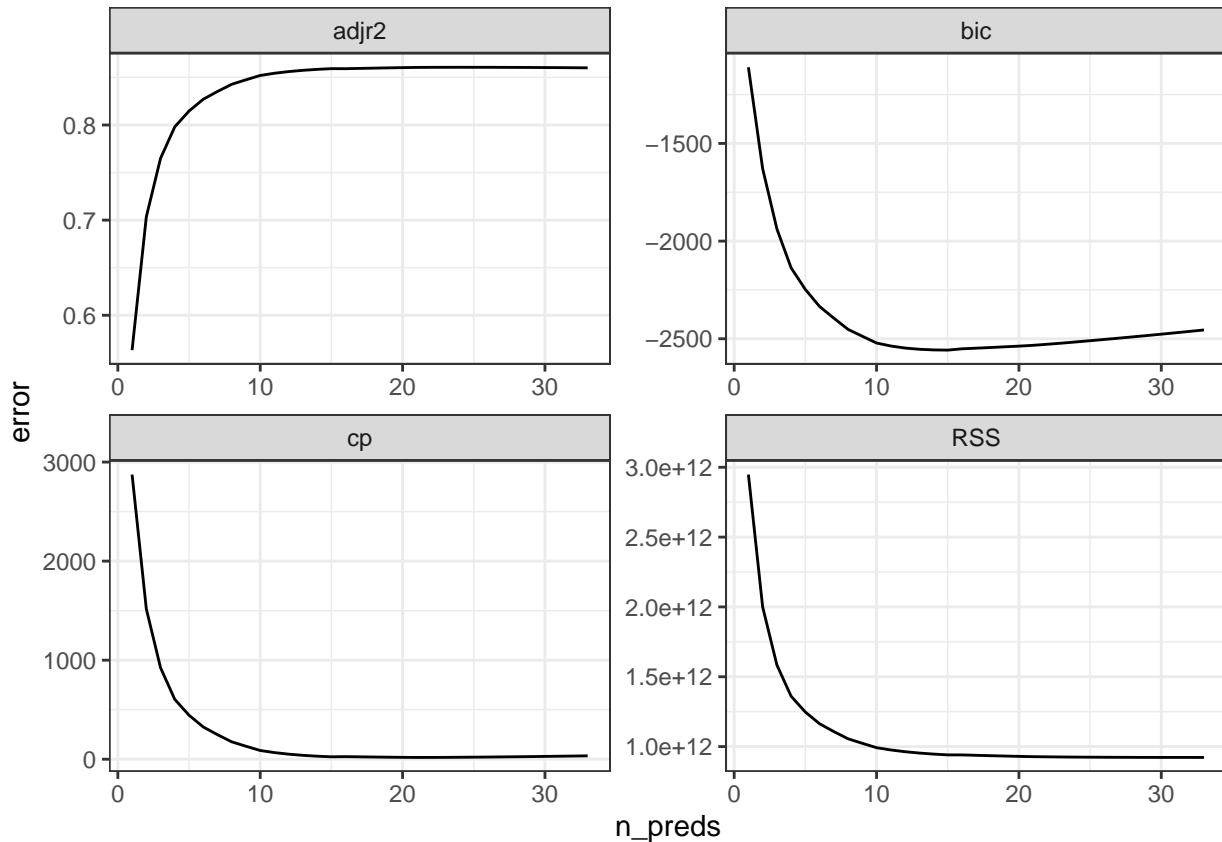


We create a data frame with the metric values (RSS, Adjusted R^2 , Cp, and BIC) for the best model with each number of predictors.

```
back_elim_metrics <- data.frame(
  RSS = summary(back_elim)$rss,
  adjr2 = summary(back_elim)$adjr2,
  cp = summary(back_elim)$cp,
  bic = summary(back_elim)$bic,
  n_preds = 1:33
)
```

Using this data frame, we plot the values of each metric as a function of the number of predictors. Note that BIC reaches a local minimum around 15 predictors, which is around the point CP, Adjusted R^2 and RSS reach a point of diminishing returns.

```
back_elim_metrics %>%
  pivot_longer(!n_preds, names_to = "metric", values_to = "error") %>%
  ggplot(aes(x = n_preds, y = error)) +
  geom_line() +
  facet_wrap(~metric, scales = "free")+
  theme_bw()
```



The following table summarizes the number of predictors that attain the lowest / highest value of the metric. In particular, BIC selects a model with 15 variables, while Adjusted R^2 selects a model with 24 variables, and CP a model with 21 variables. RSS selects the full model, since RSS is minimized on training data for models of greatest complexity.

```
back_elim_metrics %>%
  summarize(
    min_rss = which.min(RSS),
    max_adjr2 = which.max(adjr2),
    min_cp = which.min(cp),
    min_bic = which.min(bic)
  )
```

	min_rss	max_adjr2	min_cp	min_bic
## 1	33	24	21	15

Model Assessment

Choose 3 models based on the information in the previous component, and then compute rMSE for each model on the *validation* set. Which model performed best?

For comparison purposes, we'll create the models selected by each of the 4 metrics. This does present a coding challenge, since `regsubsets` can output the vector of coefficients used in the model, as well as the model coefficients, but it does not actually output an `lm` object.

One option is to use the `reformulate` function, which converts a vector of character strings to an R formula

```
adjr2_names <- names(coef(back_elim, 24)[-1])
adjr2_formula <- reformulate(adjr2_names, response = "Sale_Price")
house24 <- lm(adjr2_formula, data = house_trn_num)

cp_names <- names(coef(back_elim, 21)[-1])
cp_formula <- reformulate(cp_names, response = "Sale_Price")
house21 <- lm(cp_formula, data = house_trn_num)

bic_names <- names(coef(back_elim, 15)[-1])
bic_formula <- reformulate(bic_names, response = "Sale_Price")
house15 <- lm(bic_formula, data = house_trn_num)

house_full <- lm(Sale_Price ~ ., house_trn_num)
```

We'll also write a function to get rMSE from a model.

```
house_tst_num <- house_tst %>% select(is.numeric)
get_rmse <- function(mod){
  preds <- predict(mod, house_tst_num)
  rmse <- sqrt(mean((house_tst_num$Sale_Price - preds)^2))
  rmse
}
```

And then apply our function to each of the models:

```
data.frame(n_preds = c(15, 21, 24, 33), rMSE = c(
  get_rmse(house15),
  get_rmse(house21),
  get_rmse(house24),
  get_rmse(house_full)
)) %>%
  arrange(rMSE)
```

```
##   n_preds    rMSE
## 1      33 25418.64
## 2      24 25429.70
## 3      21 25548.64
## 4      15 25831.04
```

All models performed relatively similarly, with rMSE within 400 of each other (and the difference in rMSE between the full model and the model with 24 predictors was only 11.06!). However, it does seem the full model performed best on the validation set.

With that said, given how close the error was on the validation set, I would favor a simpler and more interpretable model over the full model.

Your model

Choose 1 of the 3 models from the previous part that you think will be most accurate on my test set.

The following two functions will help me assess your model accuracy. Copy the following templates and modify to create R functions for your model. Be sure to change the name of the functions to your own first and last names.

These functions should be self-contained, so include any packages you need or data processing you use. I will input the training data and run in a separate .Rmd, so it is important it can stand alone.

Because you are just doing feature selection, there is no need to perform data processing.

#This function creates your linear model. I will apply it to the results of FirstName_LastName_processing

```
Nate_Wells_model <- function(training_data){
  library(tidyverse)      ## Load whatever packages you need
  my_mod <- lm(Sale_Price ~ Lot_Frontage + Lot_Area + Year_Built + Year_Remod_Add +
    Mas_Vnr_Area + BsmtFin_SF_2 + Bsmt_Unf_SF + Total_Bsmt_SF +
    Gr_Liv_Area + Bedroom_AbvGr + Kitchen_AbvGr + Garage_Area +
    Wood_Deck_SF + Screen_Porch + Latitude, data = training_data)    ## Create your model. Replace 1 wi
  my_mod                  ##return your model as output
}
```

This function makes predictions for the Sale_Price of houses.

I will apply it to the results of FirstName_LastName_processing(house_test) and FirstName_LastName_mo

```
Nate_Wells_predictions <- function(model, test_data){
  library(tidyverse)      ## Load whatever packages you need
  my_preds <- predict(model, test_data)    ## Make predictions based on your model. Don't change this li
  my_preds                  ##return your predictions as output
}
```

To verify that your functions are working as desired, open a new .Rmd file, load the `house_train` data, copy the code for your 3 functions over to the new .Rmd, and then run the following code:

```
library(dplyr)
B <- sample_n(house, size = 100) # This creates a test set of 100 observations
mod <- FirstName_LastName_model(house) # Change to your First and Last Name
FirstName_LastName_predictions(mod,B) # Change to your First and Last Name
```

If everything is working correctly, the result of the code should be 100 predicted sale prices.

Model Interpretation

Consider the highly correlated predictors you identified in the first component. Did any pairs of these predictors appear in your final model? If so, why do you think this is? If not, why not?

Of the pairs of predictors identified in the first component, only `Gr_Liv_Area` and `Second_Flr_SF` both appeared in the final model (among other pairs, at most one of the predictors appeared). It's not surprising that pairs do not appear, since `C_p`, adjusted `R2`, and `BIC` all penalize additional predictors which do not appreciably decrease `RSS` (which would happen if a variable is added which is highly correlated with another).

However, `Gr_Liv_Area` was the single most important predictor, and so `Second_Flr_SF` may add enough information to the model to warrant its inclusion alongside `Gr_Liv_Area`.

Optional Model Enhancement

Optionally, you may use this space to build a model that improves on the one from Homework by incorporating feature selection, along with transformations, interaction terms, and feature engineering. I will assess your optional model on test data, alongside the model you made in the previous part. Use the following functions to create this model. **NOTE that the number 2 should immediately follow your last name**

```
#This function performs data processing. Anything you do to the training set must be repeated on the test set.  
# I will apply this function to the test and training data
```

```
FirstName_LastName2_processing <- function(my_data){  
  library(tidyverse) ## Load whatever packages you need  
  processed_data <- my_data %>% mutate(Sale_Price = Sale_Price) ## Include all relevant processing steps  
  processed_data ## returns the processed data as output  
}
```

```
#This function creates your linear model. I will apply it to the results of FirstName_LastName_processing
```

```
FirstName_LastName2_model <- function(training_data){  
  library(tidyverse) ## Load whatever packages you need  
  my_mod <- lm(Sale_Price ~ 1, data = training_data) ## Create your model. Replace 1 with your actual formula  
  my_mod ##return your model as output  
}
```

```
# This function makes predictions for the Sale_Price of houses.  
# I will apply it to the results of FirstName_LastName_processing(house_test) and FirstName_LastName_model  
# If you performed any transformations on the response variable, you must transform the predicted value.
```

```
FirstName_LastName2_predictions <- function(model, test_data){  
  library(tidyverse) ## Load whatever packages you need  
  my_preds <- predict(model, test_data) ## Make predictions based on your model. Don't change this line  
  my_preds <- my_preds*1 ## Transform your model predictions back to the original units for Sale_Price,  
  my_preds ##return your predictions as output  
}
```

To verify that your functions are working as desired, open a new .Rmd file, load the house_train data, copy the code for your 3 functions over to the new .Rmd, and then run the following code:

```
A <- FirstName_LastName2_processing(house) # Change to your First and Last Name  
library(dplyr)  
B <- FirstName_LastName2_processing(sample_n(house, size = 100)) # This creates a test set of 10 observations  
mod <- FirstName_LastName2_model(A) # Change to your First and Last Name  
FirstName_LastName2_predictions(mod,B) # Change to your First and Last Name
```