

Week 9 - Homework

STAT 420, Summer 2018, Unger

Exercise 1 (longley Macroeconomic Data)

The built-in dataset `longley` contains macroeconomic data for predicting employment. We will attempt to model the `Employed` variable.

```
View(longley)
?longley
```

(a) What is the largest correlation between any pair of predictors in the dataset?

Solution:

```
sort(unique(as.numeric(cor(longley[-7]))), decreasing = TRUE)[2]
```

```
## [1] 0.9953
```

(b) Fit a model with `Employed` as the response and the remaining variables as predictors. Calculate and report the variance inflation factor (VIF) for each of the predictors. Which variable has the largest VIF? Do any of the VIFs suggest multicollinearity?

Solution:

```
employ_mod = lm(Employed ~ ., data = longley)
#summary(employ_mod)
library(car)
vif(employ_mod)
```

```
## GNP.deflator      GNP    Unemployed Armed.Forces  Population      Year
##      135.532      1788.513      33.619      3.589      399.151      758.981
```

```
vif(employ_mod)[which.max(vif(employ_mod))]
```

```
## GNP
## 1789
```

The VIFs for every predictor except for `Armed.Forces` are extremely large, suggesting multicollinearity. `GNP` has the largest VIF.

(c) What proportion of the observed variation in `Population` is explained by a linear relationship with the other predictors?

Solution:

```
partfit_1 = lm(Employed ~ . - Population, data = longley)
partfit_2 = lm(Population ~ . - Employed, data = longley)
summary(partfit_2)$r.squared
```

```
## [1] 0.9975
```

99.75% of the variation of `Population` is explained by a linear relationship with the other *predictors*.

(d) Calculate the partial correlation coefficient for `Population` and `Employed` with the effects of the other predictors removed.

Solution:

```
cor(resid(partfit_2), resid(partfit_1))
```

```
## [1] -0.07514
```

(e) Fit a new model with `Employed` as the response and the predictors from the model in (b) that were significant. (Use $\alpha = 0.05$.) Calculate and report the variance inflation factor for each of the predictors. Which variable has the largest VIF? Do any of the VIFs suggest multicollinearity?

Solution:

```
summary(employ_mod)
```

```
##
## Call:
## lm(formula = Employed ~ ., data = longley)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.4101 -0.1577 -0.0282  0.1016  0.4554
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -3.48e+03   8.90e+02  -3.91  0.00356 **
## GNP.deflator   1.51e-02   8.49e-02   0.18  0.86314
## GNP           -3.58e-02   3.35e-02  -1.07  0.31268
## Unemployed    -2.02e-02   4.88e-03  -4.14  0.00254 **
## Armed.Forces  -1.03e-02   2.14e-03  -4.82  0.00094 ***
## Population    -5.11e-02   2.26e-01  -0.23  0.82621
## Year           1.83e+00   4.55e-01   4.02  0.00304 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.305 on 9 degrees of freedom
## Multiple R-squared:  0.995, Adjusted R-squared:  0.992
## F-statistic: 330 on 6 and 9 DF, p-value: 4.98e-10
```

```
employ_mod_small = lm(Employed ~ Year + Armed.Forces + Unemployed, data = longley)
vif(employ_mod_small)
```

```
##           Year Armed.Forces  Unemployed
##          3.891          2.223          3.318
```

```
vif(employ_mod_small)[which.max(vif(employ_mod_small))]
```

```
## Year
## 3.891
```

None of these VIFs appear to be a problem. `Year` has the largest VIF. Note that we have fixed the multicollinearity, but that does not necessarily justify this model. This “procedure” for selecting variables is not justified in practice.

(f) Use an F -test to compare the models in parts (b) and (e). Report the following:

- The null hypothesis
- The test statistic
- The distribution of the test statistic under the null hypothesis
- The p-value

- A decision
- Which model you prefer, (b) or (e)

Solution:

```
anova(employ_mod_small, employ_mod)
```

```
## Analysis of Variance Table
##
## Model 1: Employed ~ Year + Armed.Forces + Unemployed
## Model 2: Employed ~ GNP.deflator + GNP + Unemployed + Armed.Forces + Population +
##      Year
##   Res.Df    RSS Df Sum of Sq    F Pr(>F)
## 1      12  1.323
## 2       9  0.836  3      0.487 1.75   0.23
```

- Null: $\beta_{GNP.def} = \beta_{GNP} = \beta_{Pop} = 0$
- TS: $F = 1.75$
- Distribution: F with degrees of freedom 3 and 9
- p-value: 0.23
- Decision: Do NOT reject the null hypothesis.
- Prefer: The smaller model, based on the F test

(g) Check the assumptions of the model chosen in part (f). Do any assumptions appear to be violated?

Solution:

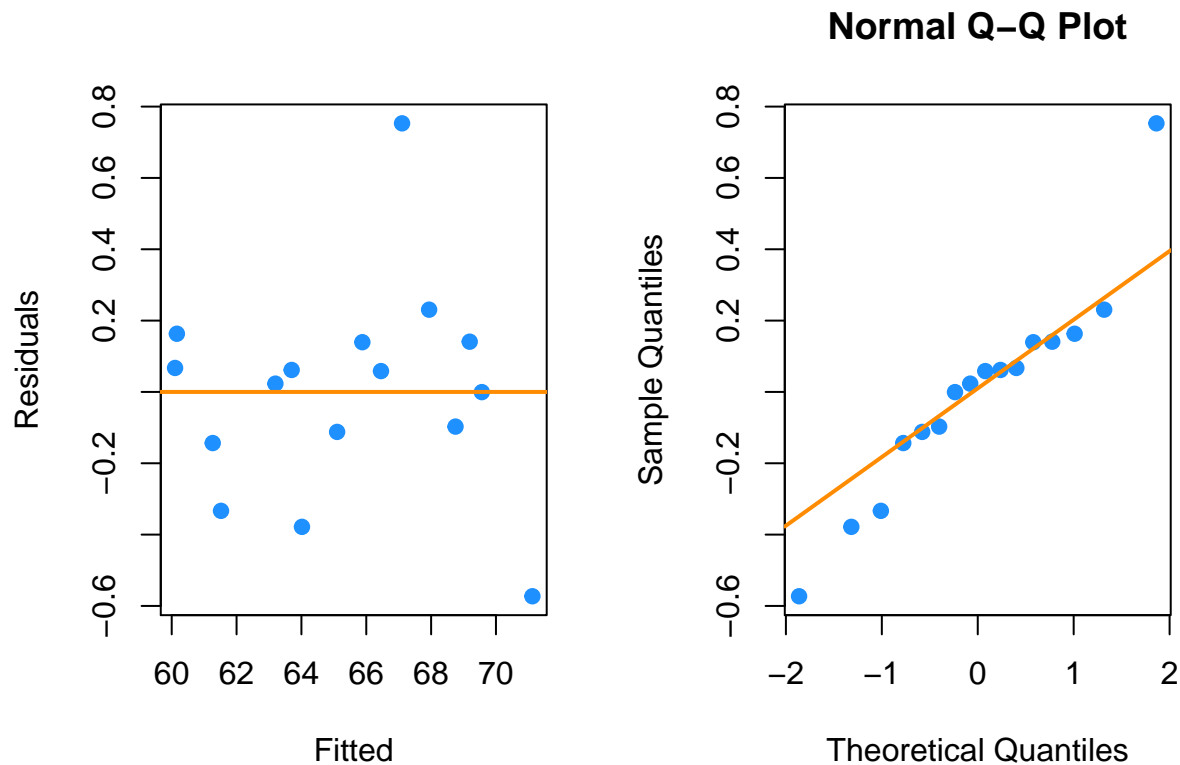
```
library(lmtest)
bptest(employ_mod_small)
```

```
##
## studentized Breusch-Pagan test
##
## data: employ_mod_small
## BP = 2.5, df = 3, p-value = 0.5
```

```
shapiro.test(resid(employ_mod_small))
```

```
##
## Shapiro-Wilk normality test
##
## data: resid(employ_mod_small)
## W = 0.93, p-value = 0.2
```

```
par(mfrow = c(1, 2))
plot_fitted_resid(employ_mod_small)
plot_qq(employ_mod_small)
```



There do not appear to be any violations of assumptions.

Exercise 2 (Credit Data)

For this exercise, use the `Credit` data from the `ISLR` package. Use the following code to remove the `ID` variable which is not useful for modeling.

```
library(ISLR)
data(Credit)
Credit = subset(Credit, select = -c(ID))
```

Use `?Credit` to learn about this dataset.

(a) Find a “good” model for `balance` using the available predictors. Use any methods seen in class except transformations of the response. The model should:

- Reach a LOOCV-RMSE below 135
- Obtain an adjusted R^2 above 0.90
- Fail to reject the Breusch-Pagan test with an α of 0.01
- Use fewer than 10 β parameters

Store your model in a variable called `mod_a`. Run the two given chunks to verify your model meets the requested criteria. If you cannot find a model that meets all criteria, partial credit will be given for meeting at least some of the criteria.

Solution:

```
library(lmtest)

get_bp_decision = function(model, alpha) {
```

```

decide = unname(bptest(model)$p.value < alpha)
ifelse(decide, "Reject", "Fail to Reject")
}

get_sw_decision = function(model, alpha) {
  decide = unname(shapiro.test(resid(model))$p.value < alpha)
  ifelse(decide, "Reject", "Fail to Reject")
}

get_num_params = function(model) {
  length(coef(model))
}

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

get_adj_r2 = function(model) {
  summary(model)$adj.r.squared
}

mod_a = lm(Balance ~ log(Income) + (Limit) + (Cards) + Age + Education + Gender + Student,
           data = Credit)

get_loocv_rmse(mod_a)
get_adj_r2(mod_a)
get_bp_decision(mod_a, alpha = 0.01)
get_num_params(mod_a)

## [1] 131.4
## [1] 0.9205
## [1] "Fail to Reject"
## [1] 8

```

(b) Find another “good” model for **balance** using the available predictors. Use any methods seen in class except transformations of the response. The model should:

- Reach a LOOCV-RMSE below 125
- Obtain an adjusted R^2 above 0.91
- Fail to reject the Shapiro-Wilk test with an α of 0.01
- Use fewer than 25 β parameters

Store your model in a variable called **mod_b**. Run the two given chunks to verify your model meets the requested criteria. If you cannot find a model that meets all criteria, partial credit will be given for meeting at least some of the criteria.

Solution:

```

library(lmtest)

get_bp_decision = function(model, alpha) {
  decide = unname(bptest(model)$p.value < alpha)
  ifelse(decide, "Reject", "Fail to Reject")
}

```

```

get_sw_decision = function(model, alpha) {
  decide = unname(shapiro.test(resid(model))$p.value < alpha)
  ifelse(decide, "Reject", "Fail to Reject")
}

get_num_params = function(model) {
  length(coef(model))
}

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

get_adj_r2 = function(model) {
  summary(model)$adj.r.squared
}

mod_b = step(
  lm(Balance ~
      (log(Income) + Limit + Cards + Age + Education + Gender + Student + Married) ^ 2,
      data = Credit),
  trace = 0)

get_loocv_rmse(mod_b)
get_adj_r2(mod_b)
get_sw_decision(mod_b, alpha = 0.01)
get_num_params(mod_b)

## [1] 118.6
## [1] 0.9377
## [1] "Fail to Reject"
## [1] 20

```

Exercise 3 (Sacramento Housing Data)

For this exercise, use the `Sacramento` data from the `caret` package. Use the following code to perform some preprocessing of the data.

```

library(caret)
library(ggplot2)
data(Sacramento)
sac_data = Sacramento
sac_data$limits = factor(ifelse(sac_data$city == "SACRAMENTO", "in", "out"))
sac_data = subset(sac_data, select = -c(city, zip))

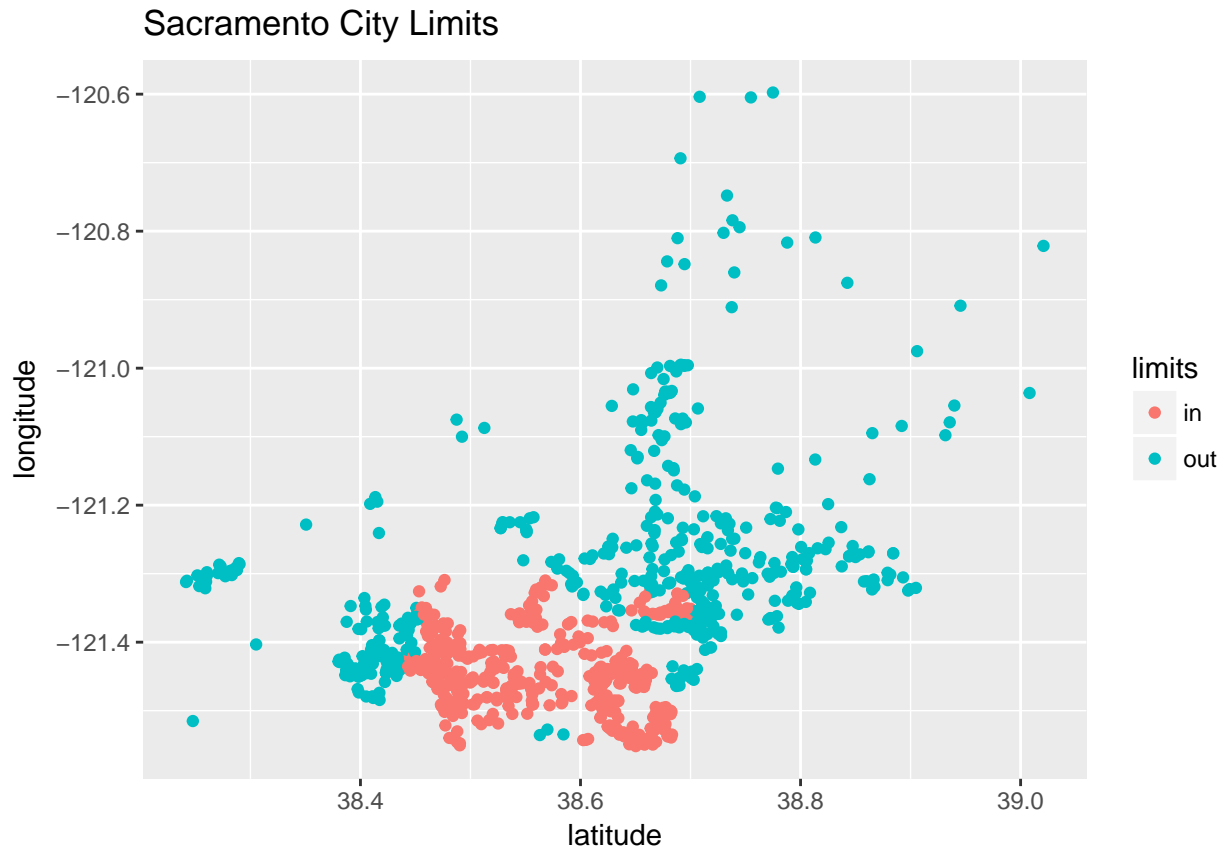
```

Instead of using the `city` or `zip` variables that exist in the dataset, we will simply create a variable (`limits`) indicating whether or not a house is technically within the city limits of Sacramento. (We do this because they would both be factor variables with a **large** number of levels. This is a choice that is made due to laziness, not necessarily because it is justified. Think about what issues these variables might cause.)

Use `?Sacramento` to learn more about this dataset.

A plot of longitude versus latitude gives us a sense of where the city limits are.

```
qplot(y = longitude, x = latitude, data = sac_data,  
      col = limits, main = "Sacramento City Limits")
```



After these modifications, we test-train split the data.

```
set.seed(420)  
sac_trn_idx = sample(nrow(sac_data), size = trunc(0.80 * nrow(sac_data)))  
sac_trn_data = sac_data[sac_trn_idx, ]  
sac_tst_data = sac_data[-sac_trn_idx, ]
```

The training data should be used for all model fitting. Our goal is to find a model that is useful for predicting home prices.

(a) Find a “good” model for **price**. Use any methods seen in class. The model should reach a LOOCV-RMSE below 77,500 in the training data. Do not use any transformations of the response variable.

Solution:

```
sac_mod_lm = step(lm(price ~ (. - baths) ^ 2, data = sac_trn_data), trace = 0)  
get_loocv_rmse(sac_mod_lm)
```

```
## [1] 76645
```

(b) Is a model that achieves a LOOCV-RMSE below 77,500 useful in this case? That is, is an average error of 77,500 low enough when predicting home prices? To further investigate, use the held-out test data and your model from part (a) to do two things:

- Calculate the average percent error:

$$\frac{1}{n} \sum_i \frac{|\text{predicted}_i - \text{actual}_i|}{\text{predicted}_i} \times 100$$

- Plot the predicted versus the actual values and add the line $y = x$.

Based on all of this information, argue whether or not this model is useful.

Solution:

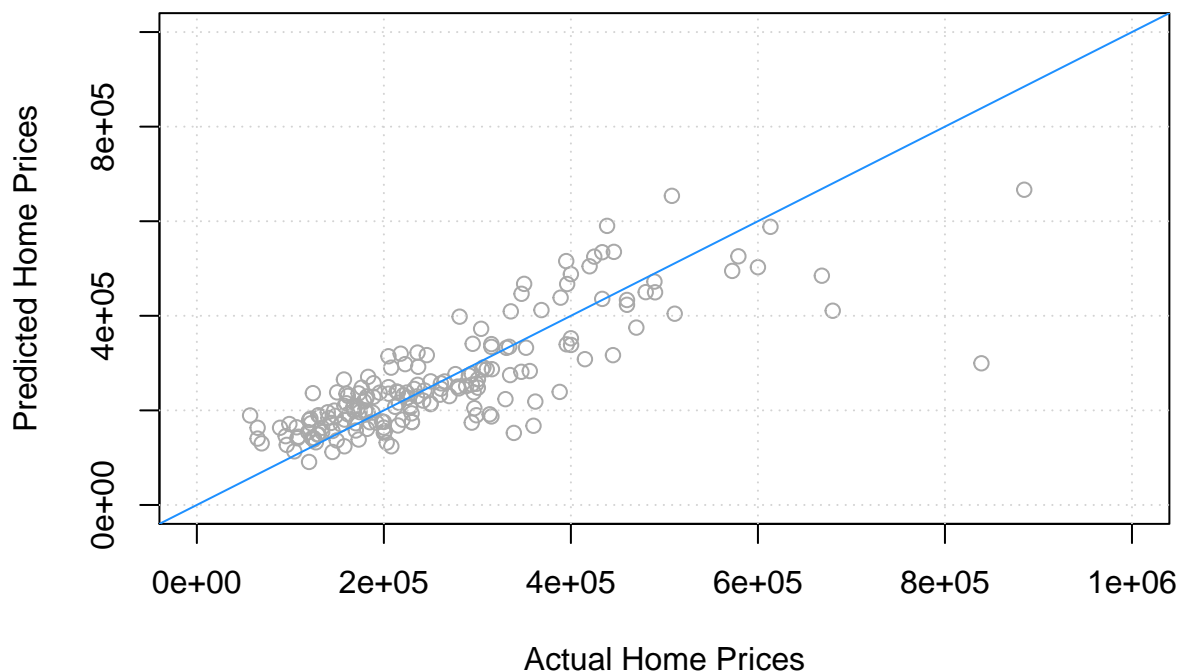
```
get_perc_err = function(actual, predicted) {
  100 * mean((abs(actual - predicted)) / actual)
}

get_perc_err(sac_tst_data$price, predict(sac_mod_lm, sac_tst_data))
```

```
## [1] 24.61
```

```
plot(sac_tst_data$price, predict(sac_mod_lm, sac_tst_data),
     xlim = c(0, 1000000), ylim = c(0, 1000000), col = "darkgrey",
     xlab = "Actual Home Prices",
     ylab = "Predicted Home Prices",
     main = "Sacramento Home Prices")
grid()
abline(0, 1, col = "dodgerblue")
```

Sacramento Home Prices



Whether or not \$77,500 is an acceptable error is highly dependent on the price of a home. That is probably a reasonable error for a \$2,000,000 home, but not for a \$150,000 home. This is why we look at the percent error instead of simply looking at RMSE. The calculated percent error here seems rather high at roughly 25%. However, when we look at the predicted versus actual plot, we see that this may also be misleading. This plot indicates that we are severely underestimating the price of extremely expensive homes, which contributes

significantly to the overall percent error. This suggests the need for further analysis, in particular, the need for a more flexible model, or perhaps some careful evaluation of “outliers.”

Exercise 4 (Does It Work?)

In this exercise, we will investigate how well backwards AIC and BIC actually perform. For either to be “working” correctly, they should result in a low number of both **false positives** and **false negatives**. In model selection,

- **False Positive**, FP: Incorrectly including a variable in the model. Including a *non-significant* variable
- **False Negative**, FN: Incorrectly excluding a variable in the model. Excluding a *significant* variable

Consider the **true** model

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \beta_5 x_5 + \beta_6 x_6 + \beta_7 x_7 + \beta_8 x_8 + \beta_9 x_9 + \beta_{10} x_{10} + \epsilon$$

where $\epsilon \sim N(0, \sigma^2 = 4)$. The true values of the β parameters are given in the R code below.

```
beta_0 = 1
beta_1 = -1
beta_2 = 2
beta_3 = -2
beta_4 = 1
beta_5 = 1
beta_6 = 0
beta_7 = 0
beta_8 = 0
beta_9 = 0
beta_10 = 0
sigma = 2
```

Then, as we have specified them, some variables are significant, and some are not. We store their names in R variables for use later.

```
not_sig = c("x_6", "x_7", "x_8", "x_9", "x_10")
signif = c("x_1", "x_2", "x_3", "x_4", "x_5")
```

We now simulate values for these x variables, which we will use throughout part (a).

```
set.seed(420)
n = 100
x_1 = runif(n, 0, 10)
x_2 = runif(n, 0, 10)
x_3 = runif(n, 0, 10)
x_4 = runif(n, 0, 10)
x_5 = runif(n, 0, 10)
x_6 = runif(n, 0, 10)
x_7 = runif(n, 0, 10)
x_8 = runif(n, 0, 10)
x_9 = runif(n, 0, 10)
x_10 = runif(n, 0, 10)
```

We then combine these into a data frame and simulate y according to the true model.

```
sim_data_1 = data.frame(x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9, x_10,
  y = beta_0 + beta_1 * x_1 + beta_2 * x_2 + beta_3 * x_3 + beta_4 * x_4 +
    beta_5 * x_5 + rnorm(n, 0, sigma)
)
```

We do a quick check to make sure everything looks correct.

```
head(sim_data_1)
```

```
##      x_1  x_2  x_3  x_4  x_5  x_6  x_7  x_8  x_9  x_10      y
## 1 6.055 4.088 8.7894 1.8180 0.8198 8.146 9.7305 9.6673 6.915 4.5523 -11.627
## 2 9.703 3.634 5.0768 5.5784 6.3193 6.033 3.2301 2.6707 2.214 0.4861 -0.147
## 3 1.745 3.899 0.5431 4.5068 1.0834 3.427 3.2223 5.2746 8.242 7.2310 15.145
## 4 4.758 5.315 7.6257 0.1287 9.4057 6.168 0.2472 6.5325 2.102 4.5814 2.404
## 5 7.245 7.225 9.5763 3.0398 0.4194 5.937 9.2169 4.6228 2.527 9.2349 -7.910
## 6 8.761 5.177 1.7983 0.5949 9.2944 9.392 1.0017 0.4476 5.508 5.9687 9.764
```

Now, we fit an incorrect model.

```
fit = lm(y ~ x_1 + x_2 + x_6 + x_7, data = sim_data_1)
coef(fit)
```

```
## (Intercept)      x_1      x_2      x_6      x_7
##    -1.3758    -0.3572     2.1040     0.1344    -0.3367
```

Notice, we have coefficients for `x_1`, `x_2`, `x_6`, and `x_7`. This means that `x_6` and `x_7` are false positives, while `x_3`, `x_4`, and `x_5` are false negatives.

To detect the false negatives, use:

```
# which are false negatives?
!(signif %in% names(coef(fit)))
```

```
## [1] FALSE FALSE TRUE TRUE TRUE
```

To detect the false positives, use:

```
# which are false positives?
names(coef(fit)) %in% not_sig
```

```
## [1] FALSE FALSE FALSE TRUE TRUE
```

Note that in both cases, you could `sum()` the result to obtain the number of false negatives or positives.

(a) Set a seed equal to your birthday; then, using the given data for each `x` variable above in `sim_data_1`, simulate the response variable `y` 300 times. Each time,

- Fit an additive model using each of the `x` variables.
- Perform variable selection using backwards AIC.
- Perform variable selection using backwards BIC.
- Calculate and store the number of false negatives for the models chosen by AIC and BIC.
- Calculate and store the number of false positives for the models chosen by AIC and BIC.

Calculate the rate of false positives and negatives for both AIC and BIC. Compare the rates between the two methods. Arrange your results in a well formatted table.

Solution:

First some setup:

```
num_sims = 300
fp_aic = rep(0, num_sims)
```

```
fn_aic = rep(0, num_sims)
fp_bic = rep(0, num_sims)
fn_bic = rep(0, num_sims)
```

Then, running the simulation:

```
set.seed(420)
for (i in 1:num_sims) {

  sim_data_1$y = beta_0 + beta_1 * x_1 + beta_2 * x_2 + beta_3 * x_3 +
    beta_4 * x_4 + beta_5 * x_5 + rnorm(n, 0, sigma)

  fit = lm(y ~ ., data = sim_data_1)
  fit_back_aic = step(fit, direction = "backward", trace = 0)
  fit_back_bic = step(fit, direction = "backward", trace = 0, k = log(nrow(sim_data_1)))

  fn_aic[i] = sum(!(signif %in% names(coef(fit_back_aic))))
  fp_aic[i] = sum(names(coef(fit_back_aic)) %in% not_sig)
  fn_bic[i] = sum(!(signif %in% names(coef(fit_back_bic))))
  fp_bic[i] = sum(names(coef(fit_back_bic)) %in% not_sig)
}
```

Finally, checking the results:

```
results = data.frame(
  FP = c(mean(fp_aic), mean(fp_bic)),
  FN = c(mean(fn_aic), mean(fn_bic))
)

rownames(results) = c("AIC", "BIC")
knitr::kable(results)
```

	FP	FN
AIC	0.9000	0
BIC	0.2133	0

We see that with both AIC and BIC, no false negatives are produced, only false positives. That means, on average, both methods are returning models that are slightly too big. Also, BIC returns fewer false positives, which matches our intuition, as BIC will always return a smaller model for this sample size.

(b) Set a seed equal to your birthday; then, using the given data for each x variable below in `sim_data_2`, simulate the response variable y 300 times. Each time,

- Fit an additive model using each of the x variables.
- Perform variable selection using backwards AIC.
- Perform variable selection using backwards BIC.
- Calculate and store the number of false negatives for the models chosen by AIC and BIC.
- Calculate and store the number of false positives for the models chosen by AIC and BIC.

Calculate the rate of false positives and negatives for both AIC and BIC. Compare the rates between the two methods. Arrange your results in a well formatted table. Also compare to your answers in part (a) and suggest a reason for any differences.

```
set.seed(420)
x_1 = runif(n, 0, 10)
x_2 = runif(n, 0, 10)
```

```

x_3 = runif(n, 0, 10)
x_4 = runif(n, 0, 10)
x_5 = runif(n, 0, 10)
x_6 = runif(n, 0, 10)
x_7 = runif(n, 0, 10)
x_8 = x_1 + rnorm(n, 0, 0.1)
x_9 = x_1 + rnorm(n, 0, 0.1)
x_10 = x_2 + rnorm(n, 0, 0.1)

sim_data_2 = data.frame(x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9, x_10,
  y = beta_0 + beta_1 * x_1 + beta_2 * x_2 + beta_3 * x_3 + beta_4 * x_4 +
    beta_5 * x_5 + rnorm(n, 0, sigma)
)

```

Solution:

First some setup:

```

num_sims = 300
fp_aic = rep(0, num_sims)
fn_aic = rep(0, num_sims)
fp_bic = rep(0, num_sims)
fn_bic = rep(0, num_sims)

```

Running the simulations on the new x data:

```

set.seed(420)
for (i in 1:num_sims) {

  sim_data_2$y = beta_0 + beta_1 * x_1 + beta_2 * x_2 + beta_3 * x_3 +
    beta_4 * x_4 + beta_5 * x_5 + rnorm(n, 0, sigma)

  fit = lm(y ~ ., data = sim_data_2)
  fit_back_aic = step(fit, direction = "backward", trace = 0)
  fit_back_bic = step(fit, direction = "backward", trace = 0, k = log(length(resid(fit))))

  fn_aic[i] = sum(!(signif %in% names(coef(fit_back_aic))))
  fp_aic[i] = sum(names(coef(fit_back_aic)) %in% not_sig)
  fn_bic[i] = sum(!(signif %in% names(coef(fit_back_bic))))
  fp_bic[i] = sum(names(coef(fit_back_bic)) %in% not_sig)
}

```

Lastly, checking the results:

```

results = data.frame(
  FP = c(mean(fp_aic), mean(fp_bic)),
  FN = c(mean(fn_aic), mean(fn_bic))
)

rownames(results) = c("AIC", "BIC")
knitr::kable(results)

```

	FP	FN
AIC	1.497	0.7767
BIC	1.000	0.8367

Again, we see that BIC returns fewer false positives than AIC. However, both now return false negatives! This is a result of the higher correlations within this set of predictors.

```
cor(sim_data_2[, c(1, 2, 8, 9, 10)])
```

```
##           x_1      x_2      x_8      x_9      x_10
## x_1  1.00000 0.02989 0.99945 0.99946 0.02751
## x_2  0.02989 1.00000 0.03183 0.02290 0.99940
## x_8  0.99945 0.03183 1.00000 0.99890 0.02942
## x_9  0.99946 0.02290 0.99890 1.00000 0.02047
## x_10 0.02751 0.99940 0.02942 0.02047 1.00000
```