

Homework 7

Yunting Chiu

2021-04-09

1. **(7.1)** State the number of degrees of freedom that are associated with each of the following extra sums of squares: $SSReg(X1 | X2)$, $SSReg(X2 | X1, X3)$, $SSReg(X1, X2 | X3, X4)$, $SSReg(X1, X2, X3 | X4, X5)$.

A note about the notation. $SSReg(A | B)$ is the extra sum of squares that appeared as a result of including variables A into the regression model that already had variables B in it. Thus, it is used to compare the full model with both A and B in it against the reduced model with only B.

Ans: We can calculate degrees of freedom by counting the number of variables to the left of the "|". - $SSReg(X1 | X2) = 1$ - $SSReg(X2 | X1, X3) = 1$ - $SSReg(X1, X2 | X3, X4) = 2$ - $SSReg(X1, X2, X3 | X4, X5) = 3$

2. **(7.2)** Explain in what sense the regression sum of squares $SSReg(X1)$ is an extra sum of squares.
 - Extra sum of squares uses extra sums of squares in tests for regression coefficients. For example, there is a response variable Y and 2 predictor variables X1 and X2:
 - The reduce model is $Y = \beta_0 + \beta_1 X1 + e_i$ and compute $SSE(X1)$
 - The full model is $Y = \beta_0 + \beta_1 X1 + \beta_2 X2 + e_i$ and compute $SSE(X1, X2)$
 - So the equation can be denoted as $SSE(X1) = SSE(X1, X2) + SS$? How can we define SS? As the extra sum of squares and denote it by $SSR(X2|X1)$ so we can write as

$$SSR(X2|X1) = SSE(X1) - SSE(X1, X2)$$

- $SSR(X2|X1)$ calculates the decrease in SSE when X2 is added to the regression model, given X1 is already present.

Reference: - <https://365datascience.com/tutorials/statistics-tutorials/sum-squares/> - https://www.stat.colostate.edu/~riczw/teach/STAT540_F15/Lecture/lec09.pdf

3. **(7.28b)** For a multiple regression model with five X variables, what is the relevant extra sum of squares for testing

The equation might be:

$$Y = \beta_0 + \beta_1 X1 + \beta_2 X2 + \beta_3 X3 + \beta_4 X4 + \beta_5 X5 + e_i$$

- (a) whether or not $\beta_5 = 0$?
 - $SSR(X5 | X2, X3, X4, X5)$
- (b) whether or not $\beta_2 = \beta_4 = 0$?
 - $SSR(X2, X4 | X1, X3, X5)$

4. **(7.28b, Stat-615 only)** Show that $SSReg(X1, X2, X3, X4) = SSReg(X2, X3) + SSReg(X1|X2, X3) + SSReg(X4 | X1, X2, X3)$

Reference: - https://www.stat.colostate.edu/~riczw/teach/STAT540_F15/Lecture/lec09.pdf - <https://www.math.arizona.edu/~piegorsch/571A/STAT571A.Ch07.pdf>

$$4. SS_{\text{Reg}}(X_1, X_2, X_3, X_4) = SS_{\text{Reg}}(X_2, X_3) + SS_{\text{Reg}}(X_1 | X_2, X_3) + SS_{\text{Reg}}(X_4 | X_1, X_2, X_3) \text{ Prove it!}$$

$$SS_{\text{Reg}}(X_1 | X_2, X_3) = SSE(X_2, X_3) - SSE(X_1, X_2, X_3) = SSR(X_1, X_2, X_3) - SSR(X_2, X_3)$$

$$SS_{\text{Reg}}(X_4 | X_1, X_2, X_3) = SSE(X_1, X_2, X_3) - SSE(X_1, X_2, X_3, X_4) = SSR(X_1, X_2, X_3, X_4) - SSR(X_1, X_2, X_3)$$

$$\begin{aligned} SS_{\text{Reg}}(X_1, X_2, X_3, X_4) &= \cancel{SS_{\text{Reg}}(X_2, X_3)} + \cancel{SSR(X_1, X_2, X_3)} - \cancel{SSR(X_2, X_3)} + SSR(X_1, X_2, X_3, X_4) - \cancel{SSR(X_1, X_2, X_3)} \\ &= SSR(X_1, X_2, X_3, X_4) \quad \# \end{aligned}$$

5. (7.3, 7.24, 7.30) Continue working with the Brand Preference data, which are available on our Blackboard, on <http://statweb.lsu.edu/EXSTWeb/StatLab/DataSets/NKNWData/CH06PR05.txt>, and in the previous homework.

Recall the variables: It was collected to study the relation between degree of brand liking (Y) and moisture content (X1) and sweetness (X2) of the product.

- (a) Obtain the ANOVA table that decomposes the regression sum of squares into extra sum of squares associated with X1 and with X2, given X1.

- $SSR(X1) = 1566.45$
- $SSR(X2|X1) = 306.25$

```
brand <- read.table("./data/CH06PR05.txt")
brand %>%
  rename(Y = V1, X1 = V2, X2 = V3) -> brand

# SSR(X1)
X1 <- lm(Y ~ X1, data = brand)

# SSR(X2|X1)
X2givenX1 <- lm(Y ~ X1 + X2, data = brand)

anova(X1)

## Analysis of Variance Table
##
## Response: Y
##          Df Sum Sq Mean Sq F value    Pr(>F)
## X1          1 1566.45  1566.45   54.751 3.356e-06 ***
## Residuals 14   400.55    28.61
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

anova(X2givenX1)

## Analysis of Variance Table
##
## Response: Y
##          Df Sum Sq Mean Sq F value    Pr(>F)
## X1          1 1566.45  1566.45  215.947 1.778e-09 ***
## X2          1   306.25   306.25   42.219 2.011e-05 ***
## Residuals 13    94.30    7.25
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

(b) Test whether X2 can be dropped from the model while X1 is retained.

Consider dropping X2, the hypothesis is $H_0: \beta_2 = 0$ vs $\beta_2 \neq 0$. According to the analysis of variance table above, the p-value of X2 is 2.011e-05, indicating that there is evidence that $\beta_2 \neq 0$, so X2 cannot be removed from the model.

(c) Fit first-order simple linear regression for relating brand liking (Y) to moisture content (X1).

```
summary(X1)$coefficients[, 1]
```

```
## (Intercept)      X1
##      50.775      4.425
```

$$\hat{Y} = 50.775 + 4.425X_1$$

(d) Compare the estimated regression coefficient for X1 with the corresponding coefficient obtained in (a).

- In the X2givenX1 model, the estimated regression coefficient for X1 is 4.425.
- In the X1 model, the estimated regression coefficient for X1 is 4.425, too.

```
summary(X2givenX1)$coefficients[2,1]
```

```
## [1] 4.425
```

```
summary(X1)$coefficients[2,1]
```

```
## [1] 4.425
```

(e) Does $SS_{\text{reg}}(X1)$ equal $SS_{\text{reg}}(X1|X2)$ here? Is the difference substantial?

- There are no difference between sum of squares of X1. The first model $SS_{\text{reg}}(X1)$ is 1566.45, and the second model $SS_{\text{reg}}(X1|X2)$ is 1566.45.

```
# SSReg(X1)
anova(X1)
```

```
## Analysis of Variance Table
##
## Response: Y
##      Df Sum Sq Mean Sq F value    Pr(>F)
## X1      1 1566.45  1566.45   54.751 3.356e-06 ***
## Residuals 14   400.55    28.61
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
# SSReg(X1|X2)
X1givenX2 <- lm(Y ~ X2 + X1, data = brand)
anova(X1givenX2)
```

```
## Analysis of Variance Table
##
## Response: Y
##      Df Sum Sq Mean Sq F value    Pr(>F)
## X2      1  306.25   306.25  42.219 2.011e-05 ***
## X1      1 1566.45  1566.45 215.947 1.778e-09 ***
## Residuals 13    94.30     7.25
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

(f)

- Regress Y on X2 and obtain the residuals.

```
# residuals(lm(Y ~ X2, data = brand))
```

- Regress X1 on X2 and obtain the residuals.
 - Regress residuals from the model “Y on X2” on residuals from the model “X1 on X2”; compare the estimated slope, error sum of squares with #1. What about R^2 ?
6. (8.13) Consider a regression model $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + e$, where X_1 is a numerical variable, and X_2 is a dummy variable. Sketch the response curves (the graphs of $E(Y)$ as a function of X_1 for different values of X_2), if $\beta_0 = 25$, $\beta_1 = 0.2$, and $\beta_2 = -12$.
- The blue line indicates the association between $E(Y)$ and X_1 when $X_2 = 0$
 - The green line indicates the association between $E(Y)$ and X_1 when $X_2 = 1$

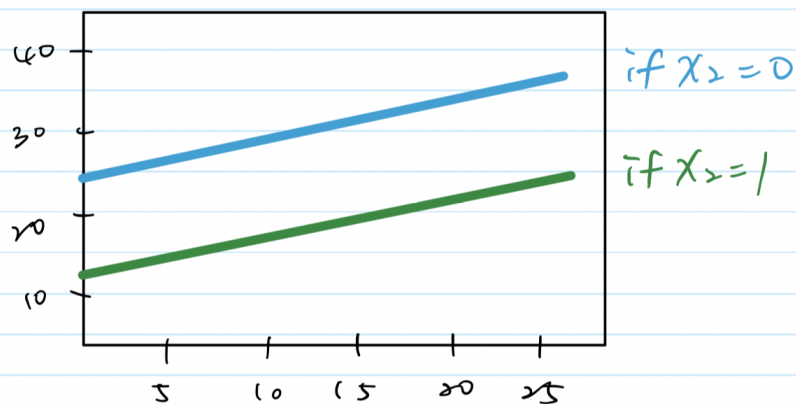
b. $Y = 25 + 0.2X_1 - 12X_2 + e$

$$E\{Y\} = 25 + 0.2X_1 + (-12)X_2$$

As X_2 is a dummy variable, so the equation can be denoted as:

if $X_2 = 0$ $E\{Y\} = 25 + 0.2X_1$

if $X_2 = 1$ $E\{Y\} = 25 + 0.2X_1 - 12 = 13 + 0.2X_1$



7. Continue the previous exercise. Sketch the response curves for the model with interaction, $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + e$, given that $\beta_3 = -0.2$

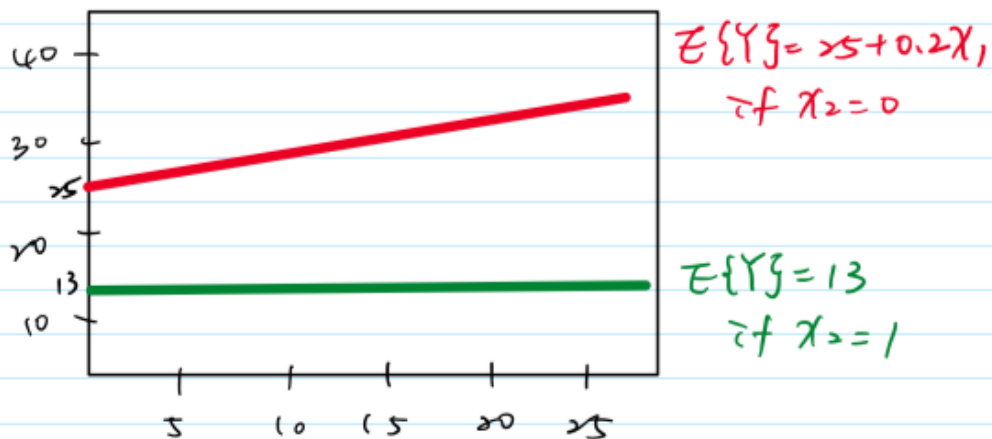
- The red line indicates the association between $E(Y)$ and X_1 when $X_2 = 0$
- The green line indicates the association between $E(Y)$ and X_1 when $X_2 = 1$

$$7. Y = 25 + 0.2X_1 + (-12)X_2 + (-0.2)X_1X_2 + \varepsilon$$

$$E\{Y\} = 25 + 0.2X_1 + (-12)X_2 + (-0.2)X_1X_2$$

$$\text{if } X_2 = 0 \Rightarrow E\{Y\} = 25 + 0.2X_1$$

$$\begin{aligned} \text{if } X_2 = 1 \Rightarrow E\{Y\} &= 25 + 0.2X_1 - 12 - 0.2X_1 \\ &= 25 - 12 = 13 \end{aligned}$$



8. (8.34) In a regression study, three types of banks were involved, namely, (1) commercial, (2) mutual savings, and (3) savings and loan. Consider the following dummy variables for the type of bank:

Type of Bank	X_2	X_3
Commercial	1	0
Mutual Saving	0	1
Saving and loan	0	0

- (a) Develop the first-order linear regression model (no interactions) for relating last year's profit or loss (Y) to size of bank (X_1) and type of bank (X_2, X_3).

$$Y_i = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + e_i$$

- (b) State the response function for the three types of banks.

- In this data, we can see the X_2 and X_3 are dummy variables. Also, Y represents profit or loss, X_1 represents the size of bank.

(b) Method: $Y_i = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \varepsilon_i$
 $E\{Y\} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3$

Commercial $E\{Y\} = \beta_0 + \beta_1 X_1 + \beta_2$

Mutual saving $E\{Y\} = \beta_0 + \beta_1 X_1 + \beta_3$

Saving and loan $E\{Y\} = \beta_0 + \beta_1 X_1$

- (c) Interpret each of the following quantities: (1) β_2 , (2) β_3 , (3) $\beta_2 - \beta_3$.

- β_2 : The difference between the commercial bank's and the savings and loan bank's expected profit or loss.

2. β_3 : The difference between the mutual saving bank's and the savings and loan bank's expected profit or loss.
3. $\beta_2 - \beta_3$: The difference between the mutual saving bank's and the commercial bank's expected profit or loss.
4. (8.16, 8.20) Refer to our old GPA data

An assistant to the director of admissions conjectured that the predictive power of the model could be improved by adding information on whether the student had chosen a major field of concentration at the time the application was submitted. Suppose that the first 10 students chose their major when they applied.

```
GPA <- read.table("./data/CH01PR19.txt")
```

```
GPA %>%
  rename(Y = "V1", X1 = "V2") -> GPA
# Suppose that the first 10 students chose their major when they applied.
GPA %>%
  mutate(X2 = 0) -> GPA
GPA$X2[1:10] = 1
GPA
```

```
##      Y X1 X2
## 1  3.897 21  1
## 2  3.885 14  1
## 3  3.778 28  1
## 4  2.540 22  1
## 5  3.028 21  1
## 6  3.865 31  1
## 7  2.962 32  1
## 8  3.961 27  1
## 9  0.500 29  1
## 10 3.178 26  1
## 11 3.310 24  0
## 12 3.538 30  0
## 13 3.083 24  0
## 14 3.013 24  0
## 15 3.245 33  0
## 16 2.963 27  0
## 17 3.522 25  0
## 18 3.013 31  0
## 19 2.947 25  0
## 20 2.118 20  0
## 21 2.563 24  0
## 22 3.357 21  0
## 23 3.731 28  0
## 24 3.925 27  0
## 25 3.556 28  0
## 26 3.101 26  0
## 27 2.420 28  0
## 28 2.579 22  0
## 29 3.871 26  0
## 30 3.060 21  0
## 31 3.927 25  0
## 32 2.375 16  0
## 33 2.929 28  0
```

##	34	3.375	26	0
##	35	2.857	22	0
##	36	3.072	24	0
##	37	3.381	21	0
##	38	3.290	30	0
##	39	3.549	27	0
##	40	3.646	26	0
##	41	2.978	26	0
##	42	2.654	30	0
##	43	2.540	24	0
##	44	2.250	26	0
##	45	2.069	29	0
##	46	2.617	24	0
##	47	2.183	31	0
##	48	2.000	15	0
##	49	2.952	19	0
##	50	3.806	18	0
##	51	2.871	27	0
##	52	3.352	16	0
##	53	3.305	27	0
##	54	2.952	26	0
##	55	3.547	24	0
##	56	3.691	30	0
##	57	3.160	21	0
##	58	2.194	20	0
##	59	3.323	30	0
##	60	3.936	29	0
##	61	2.922	25	0
##	62	2.716	23	0
##	63	3.370	25	0
##	64	3.606	23	0
##	65	2.642	30	0
##	66	2.452	21	0
##	67	2.655	24	0
##	68	3.714	32	0
##	69	1.806	18	0
##	70	3.516	23	0
##	71	3.039	20	0
##	72	2.966	23	0
##	73	2.482	18	0
##	74	2.700	18	0
##	75	3.920	29	0
##	76	2.834	20	0
##	77	3.222	23	0
##	78	3.084	26	0
##	79	4.000	28	0
##	80	3.511	34	0
##	81	3.323	20	0
##	82	3.072	20	0
##	83	2.079	26	0
##	84	3.875	32	0
##	85	3.208	25	0
##	86	2.920	27	0
##	87	3.345	27	0


```
## 88 3.956 29 0
## 89 3.808 19 0
## 90 2.506 21 0
## 91 3.886 24 0
## 92 2.183 27 0
## 93 3.429 25 0
## 94 3.024 18 0
## 95 3.750 29 0
## 96 3.833 24 0
## 97 3.113 27 0
## 98 2.875 21 0
## 99 2.747 19 0
## 100 2.311 18 0
## 101 1.841 25 0
## 102 1.583 18 0
## 103 2.879 20 0
## 104 3.591 32 0
## 105 2.914 24 0
## 106 3.716 35 0
## 107 2.800 25 0
## 108 3.621 28 0
## 109 3.792 28 0
## 110 2.867 25 0
## 111 3.419 22 0
## 112 3.600 30 0
## 113 2.394 20 0
## 114 2.286 20 0
## 115 1.486 31 0
## 116 3.885 20 0
## 117 3.800 29 0
## 118 3.914 28 0
## 119 1.860 16 0
## 120 2.948 28 0
```

- (a) Fit the regression model $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + e$, where X_1 is the entrance test score and $X_2 = 1$ if a student has indicated a major at the time of application, otherwise $X_2 = 0$. State the estimated regression function.

```
lm.fit <- lm(Y ~ X1 + X2, data = GPA)
summary(lm.fit)
```

```
##
## Call:
## lm(formula = Y ~ X1 + X2, data = GPA)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.81035 -0.33271  0.02987  0.44702  1.15523
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   2.11062    0.32220   6.551 1.6e-09 ***
## X1             0.03871    0.01282   3.018 0.00312 **
## X2             0.07728    0.20663   0.374 0.70910
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.6254 on 117 degrees of freedom
## Multiple R-squared:  0.07373,    Adjusted R-squared:  0.05789
## F-statistic: 4.656 on 2 and 117 DF,  p-value: 0.01133
```

State the Estimated Regression Function

$$\hat{Y} = 2.11062 + 2.11062X_1 + 0.07728X_2$$

(b) Test whether X_2 can be dropped from the model, using $\alpha = 0.05$.

Significance of the whole model is tested by $H_0 : \beta_2 = 0$ vs $H_1 : \beta_2 \neq 0$. With a large p-value 0.7091 and a small test statistic $F = 0.1399$, we fail to reject the null hypothesis, meaning that we have no evidence to conclude that X_2 is significant so X_2 may be removed from the model.

```
lm.fit.droppedX2 <- lm(Y ~ X1, data = GPA)
anova(lm.fit, lm.fit.droppedX2)
```

```
## Analysis of Variance Table
##
## Model 1: Y ~ X1 + X2
## Model 2: Y ~ X1
##   Res.Df    RSS Df Sum of Sq    F Pr(>F)
## 1     117 45.763
## 2     118 45.818 -1 -0.054703 0.1399 0.7091
```

(c) Fit the regression model $Y = \beta_0 + \beta_1X_1 + \beta_2X_2 + \beta_3X_1X_2 + e$ and state the estimated regression function. Interpret β_3 . Test significance of the interaction term.

```
# interaction term
lm.fit.interaction <- lm(Y ~ X1 * X2, data = GPA)
summary(lm.fit.interaction)
```

```
##
## Call:
## lm(formula = Y ~ X1 * X2, data = GPA)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.47832 -0.31337  0.04355  0.45001  1.07374
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   1.83364    0.33492   5.475 2.57e-07 ***
## X1             0.04992    0.01336   3.738 0.00029 ***
## X2             2.49114    1.00135   2.488 0.01428 *
## X1:X2         -0.09635    0.03915  -2.461 0.01531 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.6123 on 116 degrees of freedom
## Multiple R-squared:  0.1197, Adjusted R-squared:  0.09694
## F-statistic: 5.258 on 3 and 116 DF,  p-value: 0.001947
```

State the Estimated Regression Function

$$\hat{Y} = 1.83364 + 0.04992X_1 + 2.49114X_2 - 0.09635X_1X_2$$

- If $X_2 = 0$: $\hat{Y} = 1.83364 + 0.04992X_1$
- If $X_2 = 1$: $\hat{Y} = 1.83364 + 0.04992X_1 + 2.49114 - 0.09635X_1 = 4.32478 - 0.04645X_1$
- As previously stated, $X_2 = 1$ is the student has indicated a major at the time of application, otherwise X_2 is 0. The estimated value of *beta3* is -0.09635, indicating that there is an expected difference value on GPA between the students in these two groups.