# Homework 7

## Yunting Chiu

## 2021-04-06

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 aresult 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 $Y = \beta 0 + \beta 1 X1 + ei$ and compute SSE(X1)
- The full model is $Y = \beta 0 + \beta 1 X1 + \beta 2 X2 + ei$ 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:

$$Y = \beta 0 + \beta 1 X1 + \beta 2 X2 + \beta 3 X3 + \beta 4 X4 + \beta 5 X5 + ei$$

(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. $SSReg(X_1, X_2, X_3, X_4) = SSReg(X_2, X_3) + SSReg(X_1|X_2, X_3) + SSReg(X_4|X_1, X_3, X_3)$  Prove it!

$$SSReg(X_1|X_2, X_3) = SSE(X_2, X_3) - SSE(X_1, X_2, X_3) = SSR(X_1, X_2, X_3) - SSR(X_2, X_3)$$

$$SSReg(X_4|X_1, X_2, X_3) = SSE(X_1, X_2, X_3) - SSE(X_1, X_2, X_3, X_4) = SSR(X_1, X_2, X_3, X_4) - SSR(X_1, X_2, X_3)$$

$$SSReg(X_1, X_2, X_3, X_4) = SSReg(X_2, X_3) + SSR(X_1, X_2, X_3) - SSR(X_2, X_3) + SSR(X_1, X_2, X_3, X_4) - SSR(X_1, X_2, X_3)$$

$$= SSR(X_1, X_2, X_3, X_4) \#$$

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*.

```
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 H0: $\beta2 = 0$ vs $\beta2 \neq 0$. According to the analysis of variance table above, the p-value of X2 is 2.011e-05, indicating that there is evidence that $\beta2 \neq 0$, 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 SSreg(X1) equal SSreg(X1|X2) here? Is the difference substantial?

- There are no different between sum of squares of X1. The first model SSReg(X1) is 1566.45, and the second model SSReg(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 X1 is a numerical variable, and X2 is a dummy variable. Sketch the response curves (the graphs of E(Y) as a function of X1 for different values of X2), if $\beta 0 = 25$, $\beta 1 = 0.2$, and $\beta 2 = -12$.

- The blue line indicates the association between E(Y) and X1 when X2 = 0
- The green line indicates the association between E(Y) and X1 when X2 = 1

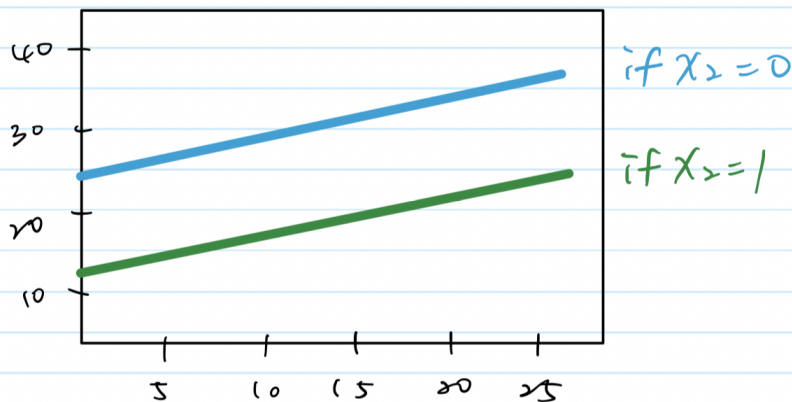6. $Y = 25 + 0.2 X_1 - 12 X_2 + \varepsilon$

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

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

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

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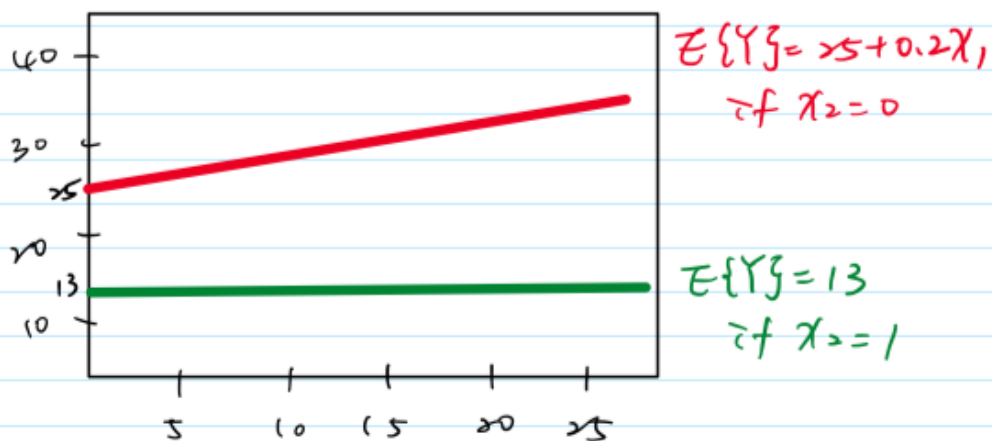


if $X_2 = 0$

if $X_2 = 1$

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 X1 when X2 = 0
- The green line indicates the association between E(Y) and X1 when X2 = 1

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

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

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

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

$= 25 - 12 = 13$



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

$E\{Y\} = 13$ if $X_2 = 1$

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 | $X_2$ | $X_3$ |
|---|---|---|
| Commerical | 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 (X1) and type of bank (X2, X3).

$$Yi = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + ei$$

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

- In this data, we can see the X2 and X3 are dummy variables. Also, Y represents profit or loss, X1 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) $\beta_2$, (2) $\beta_3$, (3) $\beta_2 - beta_3$.

1. $\beta_2$: The difference between the commercial bank's and the savings and loan bank's expected profit or loss.

2. $\beta_3$: 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

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

```
##         V1 V2
## 1    3.897 21
## 2    3.885 14
## 3    3.778 28
## 4    2.540 22
## 5    3.028 21
## 6    3.865 31
## 7    2.962 32
## 8    3.961 27
## 9    0.500 29
## 10   3.178 26
## 11   3.310 24
## 12   3.538 30
## 13   3.083 24
## 14   3.013 24
## 15   3.245 33
## 16   2.963 27
## 17   3.522 25
## 18   3.013 31
## 19   2.947 25
## 20   2.118 20
## 21   2.563 24
## 22   3.357 21
## 23   3.731 28
## 24   3.925 27
## 25   3.556 28
## 26   3.101 26
## 27   2.420 28
## 28   2.579 22
## 29   3.871 26
## 30   3.060 21
## 31   3.927 25
## 32   2.375 16
## 33   2.929 28
## 34   3.375 26
## 35   2.857 22
## 36   3.072 24
## 37   3.381 21
## 38   3.290 30
## 39   3.549 27
## 40   3.646 26
## 41   2.978 26
## 42   2.654 30
## 43   2.540 24
## 44   2.250 26
```

```
## 45   2.069 29
## 46   2.617 24
## 47   2.183 31
## 48   2.000 15
## 49   2.952 19
## 50   3.806 18
## 51   2.871 27
## 52   3.352 16
## 53   3.305 27
## 54   2.952 26
## 55   3.547 24
## 56   3.691 30
## 57   3.160 21
## 58   2.194 20
## 59   3.323 30
## 60   3.936 29
## 61   2.922 25
## 62   2.716 23
## 63   3.370 25
## 64   3.606 23
## 65   2.642 30
## 66   2.452 21
## 67   2.655 24
## 68   3.714 32
## 69   1.806 18
## 70   3.516 23
## 71   3.039 20
## 72   2.966 23
## 73   2.482 18
## 74   2.700 18
## 75   3.920 29
## 76   2.834 20
## 77   3.222 23
## 78   3.084 26
## 79   4.000 28
## 80   3.511 34
## 81   3.323 20
## 82   3.072 20
## 83   2.079 26
## 84   3.875 32
## 85   3.208 25
## 86   2.920 27
## 87   3.345 27
## 88   3.956 29
## 89   3.808 19
## 90   2.506 21
## 91   3.886 24
## 92   2.183 27
## 93   3.429 25
## 94   3.024 18
## 95   3.750 29
## 96   3.833 24
## 97   3.113 27
## 98   2.875 21
```

```
## 99   2.747 19
## 100  2.311 18
## 101  1.841 25
## 102  1.583 18
## 103  2.879 20
## 104  3.591 32
## 105  2.914 24
## 106  3.716 35
## 107  2.800 25
## 108  3.621 28
## 109  3.792 28
## 110  2.867 25
## 111  3.419 22
## 112  3.600 30
## 113  2.394 20
## 114  2.286 20
## 115  1.486 31
## 116  3.885 20
## 117  3.800 29
## 118  3.914 28
## 119  1.860 16
## 120  2.948 28
```

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.

(a) Fit the regression model $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + e$, where X1 is the entrance test score and X2 = 1 if a student has indicated a major at the time of application, otherwise X2 = 0. State the estimated regression function.
(b) Test whether X2 can be dropped from the model, using $\alpha = 0.05$.
(c) Fit the regression model $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + e$ and state the estimated regression function. Interpret $\beta_3$. Test significance of the interaction term.