

STAT 640: Homework 7

Due **Wednesday, March 23, 11:59pm MT** on the course Canvas webpage. Please follow the homework guidelines on the syllabus.

Name: **Hannah Butler**

Problem 1

For this question, use the same data and model as Problem 4 on Homework 6. Consider the null hypothesis that there is no difference in chick weight at age 6 between chicks on Diets 1, 2, and 4. (Note: this hypothesis does not involve Diet 3.)

a. Provide the form of the linear model for weight as a function of diet for the entire dataset of age 6 chicks. In other words, copy your answer to Problem 4a from Homework 6.

Answer:

$$Y = X\beta + \epsilon,$$

where

$$Y_{49 \times 1} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_{49} \end{bmatrix}, \quad X_{49 \times 4} = \begin{bmatrix} 1_{19} & 0_{19} & 0_{19} & 0_{19} \\ 0_{10} & 1_{10} & 0_{10} & 0_{10} \\ 0_{10} & 0_{10} & 1_{10} & 0_{10} \\ 0_{10} & 0_{10} & 0_{10} & 1_{10} \end{bmatrix}, \quad \beta_{4 \times 1} = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{bmatrix}, \quad \epsilon_{49 \times 1} = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_{49} \end{bmatrix}$$

and since we will be testing a hypothesis on β , we must make the distributional assumption

$$\epsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_{49 \times 49})$$

```
cw_6 <- ChickWeight %>%
  filter(Time == 6)

# Response
Y <- cw_6$weight

# Design
X <- cbind(as.numeric(cw_6$Diet == 1)
           , as.numeric(cw_6$Diet == 2)
           , as.numeric(cw_6$Diet == 3)
           , as.numeric(cw_6$Diet == 4)
           )
colnames(X) <- c("diet1"
                 , "diet2"
                 , "diet3"
                 , "diet4"
                 )
```

b. Find \mathbf{A} such that $\mathbf{A}\boldsymbol{\beta} = \mathbf{0}$ corresponds to this hypothesis.

Answer:

$$\mathbf{A} = \begin{bmatrix} 1 & -1 & 0 & 0 \\ 0 & 1 & 0 & -1 \end{bmatrix}$$

$\mathbf{A}\boldsymbol{\beta}$ is testable because each row of \mathbf{A} is estimable (ie, each row of \mathbf{A} is a linear combination of the rows of \mathbf{X}).

Note that with \mathbf{A} , we are testing the null hypothesis $H_0 : \beta_1 = \beta_2$ AND $\beta_2 = \beta_4$. If we fail to reject H_0 , then this is evidence to suggest that there is no difference between the 3 diets being tested. If we do reject H_0 , then further testing would be required to determine for which diets there is evidence of a difference. This is because there are three situations in which we should reject H_0 ; (1) $\beta_1 \neq \beta_2, \beta_2 = \beta_4$, (2) $\beta_1 = \beta_2, \beta_2 \neq \beta_4$, and (3) $\beta_1 \neq \beta_2, \beta_2 \neq \beta_4$. In the first two cases, we can infer that one diet differs from the other two which do not differ from each other, but in the third situation, it is not clear whether $\beta_1 = \beta_4$ or not and further testing is needed.

```
# Hypothesis: No difference between diets 1, 2, and 4
A <- rbind(c(1, -1, 0, 0)
           , c(1, 0, 0, -1)
           )
```

c. Compute $RSS_H - RSS$ using only $\hat{\boldsymbol{\beta}}$, \mathbf{A} , and \mathbf{X} .

Answer: By Proposition 5.4 in the notes, If $H : \mathbf{A}\boldsymbol{\beta} = \mathbf{0}$ is a testable hypothesis, then

$$RSS_H - RSS = (\mathbf{A}\hat{\boldsymbol{\beta}})^T (\mathbf{A}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{A}^T)^{-1} (\mathbf{A}\hat{\boldsymbol{\beta}}).$$

Where $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$. In R:

```
# (X'X)^-
XX_inv <- solve(t(X) %*% X)

# Find beta estimate
bh <- XX_inv %*% (t(X) %*% Y)

# Compute difference in residual SS
RSS_diff <- t(A %*% bh) %*% solve( A %*% XX_inv %*% t(A) ) %*% (A %*% bh)
RSS_diff

##           [,1]
## [1,] 1972.773
```

d. Conduct an F-test to test this hypothesis. Provide the test statistic, p-value, and a conclusion statement.

Answer: Under the assumption of the null hypothesis $\mathbf{A}\boldsymbol{\beta} = \mathbf{0}$, we know, by Corollary 5.5.1, that the F statistic is computed as

$$F = \frac{(RSS_H - RSS)/q}{RSS/(n - r)},$$

and is distributed as $\mathcal{F}(q, n - r, 0)$. Using this, we can test against and make a decision regarding H_0 .

```
n <- length(Y)
r <- rankMatrix(X)[1]
q <- rankMatrix(A)[1]

RSS_full <- t(Y - (X %*% bh)) %*% (Y - (X %*% bh))

# F statistic
fstat <- (RSS_diff/q)/(RSS_full/(n-r))
cat(fstat)
```

```
## 25.17208
```

```
# P-value
cat(1 - pf(fstat, q, n-r))
```

```
## 4.604304e-08
```

With an F statistic = 25.17208 and a p -value of $\approx 4.6 \times 10^{-8}$, we reject the null hypothesis that there is no difference between the average effects of diets 1, 2, and 4 on the weight of 6-day old chicks. In other words, there is evidence in the data to suggest that at least one diet has a different effect on average on the weight of 6-day old chicks.

e. Check your answer to (c) by fitting a model that corresponds to the null hypothesis and calculating RSS_H .

Answer:

Under the null hypothesis $H_0 : \beta_1 = \beta_2 = \beta_4$, we would set

$$\mathbf{X}_{H_0} = \begin{bmatrix} \mathbf{1}_{19} & \mathbf{0}_{19} \\ \mathbf{1}_{10} & \mathbf{0}_{10} \\ \mathbf{0}_{10} & \mathbf{1}_{10} \\ \mathbf{1}_{10} & \mathbf{0}_{10} \end{bmatrix} \quad \text{and} \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_{1,2,4} \\ \beta_3 \end{bmatrix}.$$

```

# C1 is diets 1, 2, or 4; C2 is diet 3
X_H0 <- cbind(X[,1] + X[,2] + X[,4]
              , X[,3]
              )

# solve for b_124 and b_3
b_H0 <- solve( t(X_H0)%*%X_H0 ) %*% t(X_H0)%*%Y

# Calculate estimated response under H0
Y_H0 <- X_H0 %*% b_H0

# Calculate RSS under H0
RSS_H0 <- t(Y - Y_H0) %*% (Y - Y_H0)

# Check difference
cat(RSS_H0 - RSS_full)

```

```
## 1972.773
```

Problem 2

Prove Proposition 5.11. That is, under the conditions of that proposition, show that

$$\frac{RSS_H - RSS}{RSS} = \frac{R^2 - R_H^2}{1 - R^2}$$

Answer: (See Appendix for statement of Proposition 5.11)

$$\begin{aligned}
 \frac{RSS_H - RSS}{RSS} &= \frac{\mathbf{Y}^T(\mathbf{P}_X - \mathbf{P}_{X_H})\mathbf{Y}}{\mathbf{Y}^T(\mathbf{I} - \mathbf{P}_X)\mathbf{Y}} \\
 &= \frac{\mathbf{Y}^T(\mathbf{P}_X - \mathbf{J}_n)\mathbf{Y} - \mathbf{Y}^T(\mathbf{P}_{X_H} - \mathbf{J}_n)\mathbf{Y}}{\mathbf{Y}^T(\mathbf{I} - \mathbf{J}_n)\mathbf{Y} - \mathbf{Y}^T(\mathbf{P}_X - \mathbf{J}_n)\mathbf{Y}} \\
 &= \frac{RSS - RSS_H}{1 - \mathbf{Y}^T(\mathbf{P}_X - \mathbf{J}_n)\mathbf{Y} / \mathbf{Y}^T(\mathbf{I} - \mathbf{J}_n)\mathbf{Y}}
 \end{aligned}$$

Problem 3

Consider the two regression lines

$$Y_{ki} = \beta_k x_i + \epsilon_{ki}$$

for $k = 1, 2$ and $i = 1, \dots, n$. Assume uncorrelated, homoscedastic errors. Find the F-statistic for testing $H : \beta_1 = \beta_2$.

Answer: This is a paired test. Consider the model

$$\mathbf{Y}_1 - \mathbf{Y}_2 = (\beta_1 - \beta_2)\mathbf{x} + (\boldsymbol{\epsilon}_1 - \boldsymbol{\epsilon}_2).$$

Let $\alpha = \beta_1 - \beta_2$ and $\boldsymbol{\zeta} = \boldsymbol{\epsilon}_1 - \boldsymbol{\epsilon}_2$. Then we have $\boldsymbol{\zeta} = \boldsymbol{\epsilon}_1 - \boldsymbol{\epsilon}_2 \sim N(0, 2\sigma^2 \mathbf{I})$.

Now the null hypothesis can be reframed as $H_0 : \alpha = 0$ and for an F test we can find $RSS_H - RSS$ as

$$RSS_H - RSS = (\hat{\alpha} - 0)^T (\mathbf{x}^T \mathbf{x})^{-1} (\hat{\alpha} - 0) = \hat{\alpha}^2 (\mathbf{x}^T \mathbf{x})$$

and RSS as

$$RSS = ((\mathbf{Y}_1 - \mathbf{Y}_2) - \hat{\mathbf{Y}})^T ((\mathbf{Y}_1 - \mathbf{Y}_2) - \hat{\mathbf{Y}}).$$

where $\hat{\mathbf{Y}} = E[\alpha \mathbf{x} + \boldsymbol{\zeta}]$.

Since we are testing the value of 1 parameter, $q = 1$, and since we have only one predictor, $r = \text{rank}(\mathbf{X}) = 1$. So

$$F = \frac{\hat{\alpha}^2 (\mathbf{x}^T \mathbf{x})}{((\mathbf{Y}_1 - \mathbf{Y}_2) - \hat{\mathbf{Y}})^T ((\mathbf{Y}_1 - \mathbf{Y}_2) - \hat{\mathbf{Y}}) / (n - 1)}.$$

and under the assumption of $H_0 : \alpha = 0$, $F \sim \mathcal{F}(1, n - 1, 0)$.

Problem 4

Consider the linear model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ with $\mathbf{X} \in \mathbb{R}^{n \times p}$ and $\text{rank}(\mathbf{X}) = r$. Let $\mathbf{A}\boldsymbol{\beta} = \mathbf{0}$ be a testable hypothesis with $\mathbf{A} \in \mathbb{R}^{q \times p}$ and $q < r$. Prove that if $\text{rank}(\mathbf{A}) = q$, then $\text{rank}(\mathbf{A}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{A}^T) = q$. (Hint: recall that $\text{rank}(\mathbf{B}\mathbf{B}^T) = \text{rank}(\mathbf{B})$.)

Answer: We can first use the property $\text{rank}(\mathbf{AB}) \leq \min(\text{rank}(\mathbf{A}), \text{rank}(\mathbf{B}))$ to show that

$$\text{rank}(\mathbf{A}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{A}^T) \leq \min(\text{rank}(\mathbf{A}), \text{rank}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{A}^T)),$$

So $\text{rank}(\mathbf{A}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{A}^T) \leq q$. Then by Definition 5.1, we can write $\mathbf{A} = \mathbf{M}\mathbf{X}$. So we have

$$\begin{aligned} \text{rank}(\mathbf{A}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{A}^T) &= \text{rank}(\mathbf{M}\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{M}^T) \\ &= \text{rank}(\mathbf{M}\mathbf{P}_\mathbf{X} \mathbf{M}^T) \\ &= \text{rank}(\mathbf{M}\mathbf{P}_\mathbf{X} \mathbf{P}_\mathbf{X}^T \mathbf{M}^T) \\ &= \text{rank}(\mathbf{M}\mathbf{P}_\mathbf{X} (\mathbf{M}\mathbf{P}_\mathbf{X})^T) \\ &= \text{rank}(\mathbf{M}\mathbf{P}_\mathbf{X}). \end{aligned}$$

Then, again using the property $\text{rank}(\mathbf{AB}) \leq \min(\text{rank}(\mathbf{A}), \text{rank}(\mathbf{B}))$,

$$\begin{aligned} \text{rank}(\mathbf{A}) &= \text{rank}(\mathbf{M}\mathbf{X}) \\ &= \text{rank}(\mathbf{M}\mathbf{P}_\mathbf{X} \mathbf{X}) \\ &\leq \min(\text{rank}(\mathbf{M}\mathbf{P}_\mathbf{X}), \text{rank}(\mathbf{X})) \end{aligned}$$

So that $\text{rank}(\mathbf{M}\mathbf{P}_\mathbf{X}) = \text{rank}(\mathbf{A}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{A}^T) \geq q$. Therefore, $\text{rank}(\mathbf{A}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{A}^T) = q$.

Problem 5

Consider the linear model

$$\begin{aligned} Y_1 &= \theta_1 + \theta_2 + \epsilon_1 \\ Y_2 &= 2\theta_2 + \epsilon_2 \\ Y_3 &= -\theta_1 + \theta_2 + \epsilon_3 \end{aligned}$$

where $E[\epsilon] = \mathbf{0}$ and $\text{Var}(\epsilon) = \sigma^2 \mathbf{I}$.

a. Show that $H : \theta_1 = 2\theta_2$ is a testable hypothesis.

Answer: $H : \theta_1 - 2\theta_2$

the vector $\mathbf{g}^T = [1 \quad -2]$ is estimable (and hence $\mathbf{a} = \mathbf{g}^T$ is testable) because

$$\mathbf{g}^T = [0 \quad -1/2 \quad -1] \begin{bmatrix} 1 & 1 \\ 0 & 2 \\ -1 & 1 \end{bmatrix} = [1 \quad -2],$$

where $\begin{bmatrix} 1 & 1 \\ 0 & 2 \\ -1 & 1 \end{bmatrix}$ is the design matrix \mathbf{X} .

b. Derive the form of the F-statistic for testing H .

Answer:

First finding the components of the F -statistic, we have

$$(\mathbf{X}^T \mathbf{X})^{-1} = \begin{bmatrix} 2 & 0 \\ 0 & 6 \end{bmatrix}^{-1} = \begin{bmatrix} 1/2 & 0 \\ 0 & 1/6 \end{bmatrix}.$$

Then

$$(\mathbf{A}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{A}^T)^{-1} = [1 \quad -2] \begin{bmatrix} 1/2 & 0 \\ 0 & 1/6 \end{bmatrix} \begin{bmatrix} 1 \\ -2 \end{bmatrix} = \frac{6}{7}$$

Computing $\hat{\beta}$ next, we have

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \begin{bmatrix} 1/2 & 0 \\ 0 & 1/6 \end{bmatrix} \begin{bmatrix} 1 & 0 & -1 \\ 1 & 2 & 1 \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \end{bmatrix} = \begin{bmatrix} 1/2 & 0 \\ 0 & 1/6 \end{bmatrix} \begin{bmatrix} Y_1 - Y_3 \\ Y_1 + 2Y_2 + Y_3 \end{bmatrix} = \begin{bmatrix} \frac{1}{2}(Y_1 - Y_3) \\ \frac{1}{6}(Y_1 + 2Y_2 + Y_3) \end{bmatrix}.$$

So

$$\mathbf{a}\hat{\beta} = [1 \quad -2] \begin{bmatrix} \frac{1}{2}(Y_1 - Y_3) \\ \frac{1}{6}(Y_1 + 2Y_2 + Y_3) \end{bmatrix} = \frac{1}{2}(Y_1 - Y_3) - \frac{1}{3}(Y_1 + 2Y_2 + Y_3) = \frac{1}{6}(Y_1 - 4Y_2 - 5Y_3).$$

We can now compute $RSS_H - RSS$:

$$RSS_H - RSS = \frac{6}{7} \cdot \frac{1}{36}(Y_1 - 4Y_2 - 5Y_3)^2 = \frac{1}{42}(Y_1 - 4Y_2 - 5Y_3)^2$$

To compute RSS , we need to find \mathbf{P}_X :

$$\begin{aligned}\mathbf{P}_X &= \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X} \\ &= \frac{1}{3} \begin{bmatrix} 2 & 1 & -1 \\ 1 & 2 & 1 \\ -1 & 1 & 2 \end{bmatrix},\end{aligned}$$

Then

$$\begin{aligned}RSS &= \mathbf{Y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{Y} \\ &= \frac{1}{3} \begin{bmatrix} Y_1 & Y_2 & Y_3 \end{bmatrix} \begin{bmatrix} 1 & -1 & 1 \\ -1 & 1 & -1 \\ 1 & -1 & 1 \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \end{bmatrix} \\ &= \frac{1}{3} (Y_1^2 + Y_2^2 + Y_3^2 - 2(Y_1 Y_2 - Y_1 Y_3 + Y_2 Y_3))\end{aligned}$$

Then we get

$$F = \frac{\frac{1}{42}(Y_1 - 4Y_2 - 5Y_3)^2/1}{\frac{1}{3}(Y_1^2 + Y_2^2 + Y_3^2 - 2(Y_1 Y_2 - Y_1 Y_3 + Y_2 Y_3))/(3-2)}$$

c. If we assume the errors are normally distributed and the null hypothesis is true, what is the distribution of F ?

Answer: By Corollary 5.5.1, $F \sim \mathcal{F}_{1,1,0}$ under the assumption that H_0 is true.

Appendix

List of R packages used:

tidyverse
ggplot2
Matrix

Relevant Definitions

Definition 5.1. The hypothesis $H : \mathbf{A}\beta = \mathbf{0}$ is **testable** if $\mathbf{a}_i^T \beta$ is an estimable function, for each row \mathbf{a}_i^T of \mathbf{A} .

Definition 5.2. The **residual sum of squares** for a model is $RSS = \sum (Y_i - \hat{Y})^2 = (\mathbf{Y} - \hat{\mathbf{Y}})^T (\mathbf{Y} - \hat{\mathbf{Y}})$.

Definition 5.4. The **sample multiple correlation coefficient** is the correlation between the pairs (Y_i, \hat{Y}_i) . That is:

$$R = \frac{\sum_i (Y_i - \bar{Y})(\hat{Y}_i - \bar{\hat{Y}})}{\sqrt{\sum_i (Y_i - \bar{Y})^2 \sum_i (\hat{Y}_i - \bar{\hat{Y}})^2}}$$

Definition 5.5. The square of the sample multiple correlation coefficient is called the **coefficient of determination**.

Relevant Propositions

Proposition 5.3 Consider the two models $\mathbf{Y} = \mathbf{X}_1 \beta_1 + \epsilon_1$ and $\mathbf{Y} = \mathbf{X}_2 \beta_2 + \epsilon_1$. The difference in RSS between these models is

$$RSS_1 - RSS_2 = \mathbf{Y}^T (\mathbf{P}_{\mathbf{X}_2} - \mathbf{P}_{\mathbf{X}_1}) \mathbf{Y}$$

Proposition 5.4 If $H : \mathbf{A}\beta = \mathbf{0}$ is a testable hypothesis, then

$$RSS_H - RSS = (\mathbf{A}\hat{\beta})^T (\mathbf{A}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{A}^T)^{-1} (\mathbf{A}\hat{\beta}).$$

Proposition 5.5 If $\mathbf{Y} \sim N(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$ with $\boldsymbol{\mu} = \mathbf{X}\beta$ and $H : \mathbf{A}\beta = \mathbf{0}$ is a testable hypothesis with $\text{rank}(\mathbf{A}) = q$, then:

$$F = \frac{(RSS_H - RSS)/q}{RSS/(n-r)} \sim \mathcal{F}\left(q, n-r, \frac{1}{2\sigma^2} \boldsymbol{\mu}^T (\mathbf{P}_{\mathbf{X}} - \mathbf{P}_{\mathcal{W}}) \boldsymbol{\mu}\right),$$

where $\mathcal{W} = \mathcal{N}(\mathbf{M}) \cap \mathcal{C}(\mathbf{X})$ and \mathbf{M} is a matrix such that $\mathbf{A} = \mathbf{M}\mathbf{X}$.

Proposition 5.11 In the linear model $\mathbf{Y} = \mathbf{X}\beta + \epsilon$, assume $\mathbf{x}_1 = \mathbf{1}$. Then the test statistic for a hypothesis of the form $H : [\mathbf{0} \quad \mathbf{A}'] \beta = \mathbf{0}$ (that is, a hypothesis that does not involve the intercept β_1) can be written:

$$F = \frac{R^2 - R_H^2}{1 - R^2} \frac{n-p}{q},$$

where R^2 and R_H^2 are the coefficients of determination for the full and reduced model, respectively.

Relevant Corollaries

Corollary 5.5.1 *If $\mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$ and $H : \mathbf{A}\boldsymbol{\beta} = \mathbf{0}$ is a testable hypothesis with $\text{rank}(\mathbf{A}) = q$, then when H is true,*

$$F = \frac{(RSS_H - RSS)/q}{RSS/(n-r)} \sim \mathcal{F}(q, n-r, 0).$$