

S631 HW9

Erik Parker

October 29, 2017

1. ALR 5.14: Using the data file *BGSall*, consider the regression of *HT18* on *HT9* and the grouping factor *Sex*.

```
rm(list = ls())

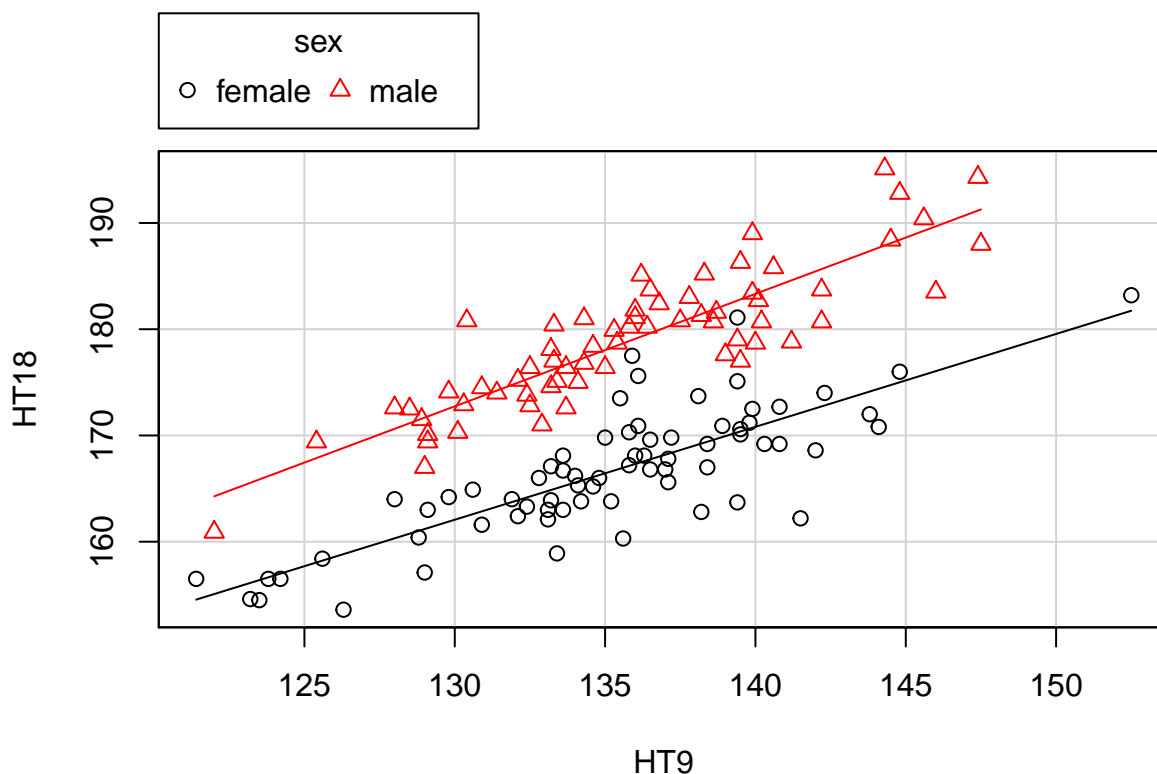
library(alr4)

Berkley <- BGSall

Berkley$sex <- ifelse(Berkley$Sex == "0", "male", "female")
```

5.14.1: Draw the scatterplot of *HT18* versus *HT9*, using a different symbol for males and females. Comment on the information in the graph about an appropriate mean function for these data.

```
scatterplot(HT18 ~ HT9 | sex, data = Berkley, smooth = FALSE, boxplots = FALSE)
```



From this plot, it seems pretty clear that there is real separation between the male and female groups in terms of their height. The intercepts of the two lines seem to be different, with the male one higher than the female one, but the slopes of the lines appear to be the same, or very close to the same. Furthermore, there is also a clear relationship in both sexes, that as the height at age 9 increases, so too does the height at age 18. This suggests to me that a proper mean function for

these data will be one with the continuous *HT9* and the categorical *sex* as predictors, but no interaction. So, it will be of the form: $HT18 \sim HT9 + sex$.

5.14.2 Obtain the appropriate test for a parallel regression model.

```
mpar <- lm(HT18 ~ HT9 + sex, data = Berkley)
```

```
summary(mpar)
```

```
##
## Call:
## lm(formula = HT18 ~ HT9 + sex, data = Berkley)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.4694  -2.0952  -0.0136   1.7101  10.4467
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  36.82147    7.29177   5.05 1.43e-06 ***
## HT9           0.96006    0.05388  17.82 < 2e-16 ***
## sexmale      11.69584    0.59036  19.81 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.432 on 133 degrees of freedom
## Multiple R-squared:  0.8516, Adjusted R-squared:  0.8494
## F-statistic: 381.7 on 2 and 133 DF, p-value: < 2.2e-16
```

Looks good! Explains quite a bit of the variation, and the coefficients for both the continuous and categorical regressors are very significant.

Assuming the parallel regression model is adequate, estimate a 95% confidence interval for the difference between males and females. For the parallel regression model, this is the difference in the intercepts of the two groups.

```
confint(mpar, level = 0.95)
```

```
##              2.5 %    97.5 %
## (Intercept) 22.3986375 51.244301
## HT9         0.8534845  1.066628
## sexmale     10.5281335 12.863548
```

So, based on my previous model, we see here that we are 95% confident that the true coefficient obtained when we move the females to males, is within the interval of 10.528 to 12.864. This means that we are 95% confident that the true increase in height seen in 18 year olds in this study is between 10.528 cm and 12.864 cm when we move from females to males.

Assignment 9 Proofs

a) $HH_R = H_R$ and $HX_1 = H_R X_1 = X_1$

First, Want to show $HX_1 = X_1$

Know $X = (X_1 | X_2)$

$HX = X$ by HW 5, so $H(X_1 | X_2) = HX = X$

and $\therefore H(X_1 | X_2) = X$

then $H(X_1 | X_2) = (HX_1 | HX_2) = X$
 $= (HX_1 | HX_2) = (X_1 | X_2)$

Therefore, $HX_1 = X_1$

Now, Want to show $HH_R = H_R$

$H_R = X_1(X_1^T X_1)^{-1} X_1^T$

So, $H X_1 (X_1^T X_1)^{-1} X_1^T \xrightarrow{\text{by above}} X_1 (X_1^T X_1)^{-1} X_1^T = H_R$
 by above $= X_1$ so $HH_R = H_R$

Now, $H_R X_1 = X_1$ because $H_R X_1$ is essentially the same as HX_1 as shown above.

Also, $H_R X_1 = (X_1(X_1^T X_1)^{-1} X_1^T) X_1 \rightarrow X_1 (X_1^T X_1)^{-1} (X_1^T X_1) \rightarrow X_1 I \rightarrow X_1$
 $I, \text{ b/c } A^T A = I$

b) $H - H_R$ is Symmetric and idempotent.

Symmetric

$(H - H_R)^T = H - H_R$

$H^T - H_R^T = H - H_R$

$(X(X^T X)^{-1} X^T)^T = (X_1(X_1^T X_1)^{-1} X_1^T)^T$

$X((X^T X)^{-1})^T X^T = X_1((X_1^T X_1)^{-1})^T X_1^T$

$X(X^{-1} (X^T)^{-1})^T X^T = X_1(X_1^{-1} (X_1^T)^{-1})^T X_1^T$

$X(X^T X)^{-1} X^T = X_1(X_1^T X_1)^{-1} X_1^T = H - H_R$

Idempotent

$(H - H_R)(H - H_R) = H - H_R$

$H^2 - HH_R - H_R H + H_R^2$ by Prop 1, and class notes: $H_R H = H_R$

$H^2 - H_R - H_R + H_R^2$

$H^2 - H_R = (H - H_R)$

by HW 5

$X(X^T X)^{-1} X^T X(X^T X)^{-1} X^T$

I_p

$X I_p (X^T X)^{-1} X^T \rightarrow X(X^T X)^{-1} X^T \rightarrow H$

Figure 1: Handwritten proofs pg. 1

c) Under H_0 , $\frac{1}{\sigma^2} SS_{reg} \sim \chi^2_q$ with $SS_{reg} = RSS_R - RSS_F$

$$RSS_R = Y^T(I - H_R)Y \quad RSS_F = Y^T(I - H)Y$$

$$SS_{reg} = Y^T(I - H_R)Y - Y^T(I - H)Y$$

$$= (Y^T I - Y^T H_R)Y - (Y^T I - Y^T H)Y$$

$$= (Y^T - Y^T H_R)Y - (Y^T - Y^T H)Y$$

$$SS_{reg} = Y^T((I - H_R) - (I - H))Y = Y^T(I - H_R + H)Y = Y^T(H - H_R)Y$$

Now, by theorem 2: if $y \sim N(0, V)$, $q = y^T A y$ then $q \sim \chi^2_r$ with $\text{rank}(A) = r$, if AV is idempotent.

So, need to now prove $Y = (y - X\beta)$ because we know $E(y - X\beta) = 0$

Given here that $\beta_0 = 0$:

$$E(y|X) = X\beta \text{ and now } X\beta = X_i\beta_i$$

Want to show that $(I - H_R)(y - X_i\beta_i) = (I - H_R)y$ from SS_{reg} above.

$$(I - H_R)y - ((I - H_R)X_i\beta_i - (I - H_R)X_i\beta_i)$$

both = 0, from part a.

$$(I - H_R)y - X_i\beta_i + X_i\beta_i$$

So, $y = (y - X_i\beta_i)$. Then, by a similar argument with the transpose, $y^T = (y - X_i\beta_i)^T$

$$\text{So, can say } SS_{reg} = \underbrace{(y - X_i\beta_i)^T}_{y^T} \underbrace{(I - H_R)}_A \underbrace{(y - X_i\beta_i)}_y = q \text{ by theorem 2}$$

also by theorem 2: $y = (y - X_i\beta_i) \sim N(0, \sigma^2 I)$

- Now, to show AV is idempotent, need $A = \frac{(H - H_R)}{\sigma^2}$ to cancel σ^2 in V .

$$(AV)(AV) = AV$$

$$\left(\frac{1}{\sigma^2}(H - H_R)(\sigma^2 I)\right)^2 = \frac{1}{\sigma^2}(H - H_R)(\sigma^2 I)$$

$$(H - H_R)^2 = H - H_R$$

by part b

$$H - H_R = H - H_R$$

Now, $\text{rank}(A)$:

$$\text{rank}\left(\frac{H - H_R}{\sigma^2}\right) = \text{rank}(H - H_R) \text{ and from notes and because } H - H_R \text{ is symmetric and idempotent}$$

$$\begin{aligned} \text{trace}(H - H_R) &= \text{trace}(X(X^T X)^{-1}X^T - X_i(X_i^T X_i)^{-1}X_i^T) \\ &= \text{tr}(X^T X(X^T X)^{-1}) - \text{tr}(X_i^T X_i(X_i^T X_i)^{-1}) \\ &= \text{trace}(I_{p-1} - I_{(p-1)}) = p - p = 0 \end{aligned}$$

$$\text{So, by theorem 2: } q = y^T \frac{(H - H_R)}{\sigma^2} y = \frac{1}{\sigma^2} SS_{reg} \sim \chi^2_q$$

(2)

Figure 2: Handwritten proofs pg. 2

d) Show SS_{reg} and $\hat{\sigma}^2$ are independent.

$$SS_{reg} = Y^T (H - H_R) Y \quad \hat{\sigma}^2 = \frac{Y^T (I - H) Y}{n - p'}$$

Using theorem 3: If $Y \sim N(H|V)$, $Q_1 = Y^T A_1 Y$ and $Q_2 = Y^T A_2 Y$
then Q_1 and Q_2 are independent if $A_1 V A_2 = 0$

$$Q_1 = SS_{reg} = Y^T \underbrace{(H - H_R)}_{A_1} Y \quad Q_2 = Y^T \underbrace{(I - H)}_{A_2} Y \quad \Rightarrow Y \sim N(X\beta, \underbrace{\sigma^2 I}_V)$$

So, need to show $A_1 V A_2 = 0$.

$$(H - H_R) \sigma^2 I \left(\frac{I - H}{n - p'} \right) \rightarrow \frac{\sigma^2}{n - p'} ((H - H_R) I (I - H)) \rightarrow \frac{\sigma^2}{n - p'} ((H I - H_R I) (I - H))$$

$$\downarrow$$

$$\frac{\sigma^2}{n - p'} (H I^2 - H^2 I - H_R I^2 + H_R H I)$$

$$\downarrow$$

$$\frac{\sigma^2}{n - p'} (H - H^2 - H_R + H_R H)$$

by part and HW 5

$$\frac{\sigma^2}{n - p'} (H - H - H_R + H_R) = 0 \quad \checkmark$$

So, SS_{reg} and $\hat{\sigma}^2$ are independent \blacksquare

e) Show $\frac{SS_{reg}}{\frac{RSS}{n - p'}} \sim F_{q, n - p'}$

by theorem 5: If $W_1 \sim \chi^2_q$, $W_2 \sim \chi^2_{n - p'}$ and W_1 and W_2 are independent,
then $\frac{W_1/q}{W_2/(n - p')} \sim F_{q, n - p'}$

here $W_1 = SS_{reg}$ and $W_2 = RSS$, need to show they are independent

Already know $\frac{1}{\sigma^2} SS_{reg} \sim \chi^2_q$ from part c.

Need to show now that $\frac{1}{\sigma^2} RSS \sim \chi^2_{n - p'}$ so $\frac{\frac{1}{\sigma^2} SS_{reg}}{\frac{1}{\sigma^2} \frac{RSS}{n - p'}} \sim F_{q, n - p'}$
 $\frac{1}{\sigma^2}$ terms cancel

First: $\frac{1}{\sigma^2} RSS = \frac{1}{\sigma^2} Y^T (I - H) Y$

need to show $(I - H)Y = (I - H)(Y - X\beta)$ like before.

$$= (I - H)Y - (I - H)X\beta$$

$$= (I - H)Y - X\beta + X\beta$$

$$= (I - H)Y, \text{ and like before can show same for } Y^T \text{ via transpose.}$$

So, $\frac{1}{\sigma^2} (Y - X\beta)^T (I - H) (Y - X\beta)$

and using LM theorem 2: $A = \frac{1}{\sigma^2} (I - H)$ $AV = I - H$ which is idempotent by HW 5.
 $V = \sigma^2 I$

Now, $r = \text{rank}(A) = \text{rank}(I - H) = \text{trace}(I - H) = \text{trace}(I) - \text{trace}(H) = n - p'$

So, $\frac{1}{\sigma^2} RSS \sim \chi^2_{n - p'}$ and know that $\frac{1}{\sigma^2} SS_{reg} \sim \chi^2_q$

- Now, need to show W_1 and W_2 are independent.

Know SS_{reg} and $\hat{\sigma}^2 = \frac{RSS}{n - p'}$ are independent by part d.

Ans, because $\hat{\sigma}^2$ and RSS are only different by a constant $(\frac{1}{n - p'})$, SS_{reg} is also independent from RSS .

(3)

Figure 3: Handwritten proofs pg. 3

Therefore: W_1 and W_2 are independent, and

$$\frac{W_1/J}{W_2/K} = \frac{SS_{\text{reg}}/q}{RSS/(n-p)} \sim F_{q, n-p}$$

(4)

Figure 4: Handwritten proofs pg. 4

ALR 6.4

6.4: With the UN data, consider testing $NH: \text{lifeExpF} \sim \log(\text{ppgdp}) + \text{group}:\log(\text{ppgdp})$ $AH: \text{lifeExpF} \sim \text{group} + \log(\text{ppgdp}) + \text{group}:\log(\text{ppgdp})$ The AH model is the most general model given at (6.10), but the NH was not given previously.

```
un <- UN11

m1 <- lm(lifeExpF ~ log(ppgdp) + group:log(ppgdp), data = un)
m2 <- lm(lifeExpF ~ group + log(ppgdp) + group:log(ppgdp), data = un)
```

6.4.1: What is the meaning of the NH model?

This NH model here is stating that the mean function is fully described by this first set of regressors, $\log(\text{ppgdp})$ and the interaction between group and $\log(\text{ppgdp})$, and that there is no influence of the main effects from group .

6.4.2 Perform the test and summarize the results

```
anova(m1, m2)

## Analysis of Variance Table
##
## Model 1: lifeExpF ~ log(ppgdp) + group:log(ppgdp)
## Model 2: lifeExpF ~ group + log(ppgdp) + group:log(ppgdp)
##   Res.Df    RSS Df Sum of Sq    F Pr(>F)
## 1     195 5232.0
## 2     193 5077.7  2     154.31 2.9326 0.05564 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

From this test, we can see that we got a p-value of about 0.06 when comparing the reduced ($m1$) and full ($m2$) models outlined above. This means that, with an alpha of 0.05, we are unable to reject the null hypothesis that the mean function is fully described by the first reduced model, containing the regressors $\log(\text{ppgdp})$ and the interaction $\text{group}:\log(\text{ppgdp})$. So, we can't reject the hypothesis that the estimated coefficient associated with the main effect for group is equal to zero, meaning that regressor does not have an effect on lifeExpF when $\log(\text{ppgdp})$ and $\text{group}:\log(\text{ppgdp})$ are already present in the model.

So, here we would (just barely with this p-value) conclude that the reduced null model is sufficient to explain the response, and that there is no need to include group in the model.

In addition, using the full model, perform the test

$$H_0: \beta_{02} - \beta_{03} = 14 \text{ and } \beta_{12} + \beta_{13} = 0.2$$

with H_A : at least one equality doesn't hold. Show your work. In addition, how could you interpret this test?

```
summary(m2)

##
## Call:
## lm(formula = lifeExpF ~ group + log(ppgdp) + group:log(ppgdp),
##     data = un)
```

```
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -18.634  -2.089   0.301   2.255  14.489
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      59.2137    15.2203   3.890 0.000138 ***
## groupother      -11.1731    15.5948  -0.716 0.474572
## groupafrica     -22.9848    15.7838  -1.456 0.146954
## log(ppgdp)        2.2425     1.4664   1.529 0.127844
## groupother:log(ppgdp)  0.9294     1.5177   0.612 0.540986
## groupafrica:log(ppgdp) 1.0950     1.5785   0.694 0.488703
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5.129 on 193 degrees of freedom
## Multiple R-squared:  0.7498, Adjusted R-squared:  0.7433
## F-statistic: 115.7 on 5 and 193 DF,  p-value: < 2.2e-16

# here beta_02 corresponds to groupother, beta_03 to groupafrica, beta_12 to
# groupother:log(ppgdp) and beta_13 to groupafrica:log(ppgdp)

# To compare level means like this, can use the corresponding code from
# class on 10.26.17:

L = matrix(c(0, 1, -1, 0, 0, 0, 0, 0, 0, 0, 1, 1), byrow = TRUE, nrow = 2)

c.vector = c(14, 0.2)
ht3 = linearHypothesis(m2, hypothesis.matrix = L, rhs = c.vector)
ht3
```

```
## Linear hypothesis test
##
## Hypothesis:
## groupother - groupafrica = 14
## groupother:log(ppgdp) + groupafrica:log(ppgdp) = 0.2
##
## Model 1: restricted model
## Model 2: lifeExpF ~ group + log(ppgdp) + group:log(ppgdp)
##
##      Res.Df    RSS Df Sum of Sq    F Pr(>F)
## 1      195 5092.6
## 2      193 5077.7  2      14.906 0.2833 0.7536
```

So, with this large p-value (F-value?) from the test performed above, I would interpret this as meaning that there is not enough support to reject the null hypothesis that both $\beta_{02} - \beta_{03} = 14$ and $\beta_{12} + \beta_{13} = 0.2$, and so it appears that both equalities hold.