

Biostat 250B HW1

Elvis Cui

Han Cui

Dept. of Biostat

UCLA



Reading material: Seber and Lee's Text, Chapter 4, section 4.4 and two papers posted on the class website entitled Reading Material 1 and Reading Material 2.

- 1 Recall that when we want to test whether the $p-1$ covariates are useful in a linear model with intercept, we showed in class that under the null hypothesis that all coefficients of the $p-1$ covariates are equal to 0, we have $ER^2=(p-1)/(n-1)$, where n is the sample size and R^2 is the square of the multiple correlation coefficient. This result was arrived at using the fact that if F is a F -variate with numerator degree of freedom a and denominator degree of freedom b , then (i) $\{aF/b\}/(1+aF/b)$ is a Beta distributed with parameters $a/2, b/2$ and (ii) the expectation of this distribution is $a/(a+b)$. Verify the latter two assertions.

Questions 2 - 3 concern pitfalls of the multiple correlation coefficient and Questions 4 - 5 review basic concepts of standardized regression coefficients and its properties along with its relationship with partial correlation with a covariate when there are two covariates in the model. Question 6 relates the partial correlation coefficient of the response with a covariate in a multiple regression model to the test statistic for testing if the coefficient of the covariate in the model is 0.

- 2 Review the paper entitled Reading Material 1, use the two data sets in the paper and see if you can confirm the reported results in Tables 1 and 2.
- 3 Write a $\frac{3}{4}$ page of your commentary on the paper entitled Reading Material 2 for your supervisor, who is not well trained in statistics, has an advanced degree in Public Health and wanted to know the main gist of the paper.
- 4 Refer to the Lecture Notes by Ernesto Amaral posted on the website. Use the data set on p.6 and verify all numerical results and the formula shown on p. 30. The last formula is for a linear regression on two explanatory variables with an intercept and expresses the square of the multiple correlation coefficient in terms of an ordinary correlation and a partial correlation coefficient.
- 5 Do Problem 2 in Ex. 4c on p.113 of Seber and Lee's text.
- 6 Consider the multiple regression of Y on x_1, \dots, x_{k-1}, x_k . Express the partial correlation coefficient of Y and x with the linear effects of x_1, \dots, x_{k-1} removed in terms of the test statistic for testing whether the coefficient of x_k in the multiple regression model is 0.

Suggested further reading: Reading Material 3 concerns use of square of the multiple correlation coefficient in nonlinear models and potential pitfalls. Review this paper when you find time to do so.

Q1: Let $F \sim F_{a,b}$. Show

$$(a) \frac{aF/b}{1+aF/b} \sim \text{Beta}\left(\frac{a}{2}, \frac{b}{2}\right)$$

$$(b) \mathbb{E} \text{Beta}\left(\frac{a}{2}, \frac{b}{2}\right) = \frac{a}{a+b}$$

Pf: (a): By definition F , we write

$$F = \frac{x_a^2}{x_a^2 + x_b^2} \frac{b}{a}, \quad x_a^2 \perp\!\!\!\perp x_b^2$$

where x_a^2, x_b^2 are χ^2 -variables with $df=a$ & b .

Since $P(F \leq 0) = 0$, we have

$$\begin{aligned} \frac{aF/b}{1+aF/b} &= \frac{F}{F+\frac{b}{a}} = \frac{1}{1+\frac{b}{a}F} \\ &= \frac{1}{1+\frac{x_b^2}{x_a^2}} = \frac{x_a^2}{x_a^2+x_b^2}. \quad (\Delta) \end{aligned}$$

Further, note $x_a^2 \sim \text{Gamma}\left(\frac{a}{2}, 2\right)$

So (Δ) is the definition of a Beta variable:

$$(\Delta) \sim \text{Beta}\left(\frac{a}{2}, \frac{b}{2}\right)$$

(b): Let $B \sim \text{Beta}\left(\frac{a}{2}, \frac{b}{2}\right)$, then

$$\mathbb{E} B = \int_{\mathbb{R}^+} b f_B(b) db$$

$$(*) = \int_{\mathbb{R}^+} b b^{\frac{a}{2}-1} (1-b)^{\frac{b}{2}-1} B\left(\frac{a}{2}, \frac{b}{2}\right) db$$

$$= \int_{\mathbb{R}^+} b^{\frac{a+2}{2}-1} (1-b)^{\frac{b}{2}-1} B\left(\frac{a}{2}, \frac{b}{2}\right) db$$

$$= B\left(\frac{a+2}{2}, \frac{b}{2}\right) / B\left(\frac{a}{2}, \frac{b}{2}\right)$$

where $B(\cdot, \cdot)$ is the Beta function, i.e.

$$B\left(\frac{a}{2}, \frac{b}{2}\right) = \frac{\Gamma\left(\frac{a}{2}\right) \Gamma\left(\frac{b}{2}\right)}{\Gamma\left(\frac{a+b}{2}\right)}$$

and (*) follows from the pdf of a Beta RV.

Thus,

$$\begin{aligned} \mathbb{E} B &= \frac{\Gamma\left(\frac{a}{2}+1\right) / \Gamma\left(\frac{a+b}{2}+1\right)}{\Gamma\left(\frac{a}{2}\right) / \Gamma\left(\frac{a+b}{2}\right)} \\ &\stackrel{(O)}{=} \frac{a}{a+b} \end{aligned}$$

where (O) follows from the recursion formula

$$\Gamma(x+1) = x \Gamma(x)$$

Q2: See Codes & Results below.

```

8 ~ R_square <- function(y, y_hat) {
9   R1 <- 1 - sum((y - y_hat) ** 2) / sum((y - mean(y)) ** 2)
10  R2 <- sum((y_hat - mean(y)) ** 2) / sum((y - mean(y)) ** 2)
11  R3 <- sum((y_hat - mean(y_hat)) ** 2) / sum((y - mean(y)) ** 2)
12  e <- y - y_hat
13  R4 <- 1 - sum((e - mean(e)) ** 2) / sum((y - mean(y)) ** 2)
14  R6 <- cor(y, y_hat) ** 2
15  R7 <- 1 - sum(e ** 2) / sum(y ** 2)
16  R8 <- sum(y_hat ** 2) / sum(y ** 2)
17
18  output <- tibble("R1" = R1,
19                    "R2" = R2,
20                    "R3" = R3,
21                    "R4" = R4,
22                    "R6" = R6,
23                    "R7" = R7,
24                    "R8" = R8)
25
26  output
27
28 mod1 <- lm(y~x, data = data1)
29 mod2 <- lm(y~0 + x, data = data1)
30 mod3 <- lm(log(y)~log(x), data = data1)
31
32 R_square(data1$y, mod1$fitted.values)
33 R_square(data1$y, mod2$fitted.values)
34 R_square(data1$y, exp(mod3$fitted.values))
> R_square(data1$y, mod1$fitted.values)
# A tibble: 1 x 7
      R1     R2     R3     R4     R6     R7     R8
    <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 0.981 0.981 0.981 0.981 0.981 0.997 0.997
> R_square(data1$y, mod2$fitted.values)
# A tibble: 1 x 7
      R1     R2     R3     R4     R6     R7     R8
    <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 0.978 1.08  1.08  0.978 0.981 0.996 0.996
> R_square(data1$y, exp(mod3$fitted.values))
# A tibble: 1 x 7
      R1     R2     R3     R4     R6     R7     R8
    <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 0.978 1.10  1.10  0.978 0.981 0.996 1.02

```

Q4: See results below.

```

38 # Q4 # formula on page 30
39
40 y <- c(1,2,3,5,3,1,5,0,6,3,7,4)
41 x <- c(1,1,1,1,2,2,3,3,4,4,5,5)
42 z <- c(12,14,16,16,18,16,12,12,10,12,10,16)
43
44 mod_q4 <- lm(y~x + z)
45 summary(mod_q4) # 0.2539
46 ry1 <- cor(y, x) ** 2
47 library(ppcor)
48 ry2.1 <- pcor(t(rbind(y,x,z)))$estimate[1,3]
49
50 ry1 + ry2.1 ** 2 * (1 - ry1) # 0.2538

```

$$S_0 \approx 0.2539 \approx 0.2538$$

\uparrow definition. \uparrow formula.

Q3: Commentary

1 Commentary

Recall the classical linear regression model:

$$\begin{aligned}\mathbf{y} &= \mathbf{X}\beta + \epsilon \\ \epsilon &\sim \mathcal{N}(0, \mathbf{W})\end{aligned}$$

As practitioners, we would like to estimate β based on observed data and address the goodness-of-fit based on some measurements such as **coefficient of determinants**.

There are 2 common ways to estimate β :

- Ordinary Least Squares (OLS): $\hat{\beta}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$.
- Weighted Least Squares (WLS): $\hat{\beta}_{WLS} = (\mathbf{X}^T \mathbf{W}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}^{-1} \mathbf{y}$.

For the second one, usually we write it in the **transformed** version, i.e.,

$$\hat{\beta}_{WLS} = (\mathbf{X}_*^T \mathbf{X}_*)^{-1} \mathbf{X}_*^T \mathbf{y}_*$$

where

$$\mathbf{X}_* = \mathbf{W}^{-1/2} \mathbf{X}, \quad \mathbf{y}_* = \mathbf{W}^{-1/2} \mathbf{y}$$

Then, the coefficient of determination can be calculated in three ways:

- $R_{OLS}^2 = 1 - \frac{\|\mathbf{y} - \mathbf{X}\hat{\beta}_{OLS}\|_2^2}{\mathbf{y}^T \mathbf{y} - n\bar{y}^2}$
- $R_{WLS}^2 = 1 - \frac{\|\mathbf{y}_* - \mathbf{X}_*\hat{\beta}_{WLS}\|_2^2}{\mathbf{y}_*^T \mathbf{y}_* - n\bar{y}_*^2}$
- pseudo $R_{WLS}^2 = 1 - \frac{\|\mathbf{y} - \mathbf{X}\hat{\beta}_{WLS}\|_2^2}{\mathbf{y}^T \mathbf{y} - n\bar{y}^2}$

The first two are more intuitive than the last one. While usually the second one (R_{WLS}^2) is higher than the first one (R_{OLS}^2), the third one (pseudo R_{WLS}^2) is always less than the first one. Therefore, "sole reliance on the **coefficient of determination** may fail to reveal important data characteristics and model inadequacies".

Q5: General full rank LRM:

R^2 & F for testing $H: \beta_j = 0, j \neq 0$

are indep. of the units in which Y_i & X_{ij} are measured.

$$\text{Sol. } F = \frac{(RSS_0 - RSS)}{RSS} \frac{n-k}{k}, \quad = \{j: \beta_j \neq 0\}$$

$$R^2 = 1 - \frac{RSS}{TSS}, \quad TSS = (RSS)_0$$

$$(RSS)_0 = y^T (I - \frac{1}{n} J) y, \quad RSS = y^T (I - P) y$$

where $y = (y_1, \dots, y_n)$ & $P = X(X^T X)^{-1} X^T$

① if $y \leftarrow \alpha y$, $\alpha \neq 0$.

$$\text{then } (RSS)_0 - RSS = \alpha^2 y^T (P - \frac{1}{n} J) y.$$

$$RSS = \alpha^2 y^T (I - P) y.$$

$$\text{So } [(RSS)_0 - RSS] / RSS \text{ and}$$

RSS_0 / RSS stay the same

since α^2 cancels out.

② if $X \leftarrow \alpha X$, $\alpha \neq 0$.

$$\text{then } P = \alpha X (\alpha^2 X^T X)^{-1} \alpha X^T$$

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

is unchanged.

$$\text{So } [(RSS)_0 - RSS] / RSS \text{ is the same}$$

To sum up, since R^2 & F are functions of RSS_0 / RSS and $(RSS_0 - RSS) / RSS$, so both of them stay the same.

Q6: Let $y = \beta_0 + \sum_{i=1}^k \beta_i X_i$

Show

$r_{yx_k, x_1, \dots, x_{k-1}}$ is a function of the test statistic for $H_0: \beta_k = 0$.

Sol. Define the following:

$$X = [x_1 \ \dots \ x_{k-1}]$$

$$\tilde{X} = [x_1 \ \dots \ x_k]$$

$$P_{x_k^\perp} = \frac{x_k (x_k^T)^T}{\|x_k\|_2^2}$$

$$P_X = X (X^T X)^{-1} X^T; \quad P_{\tilde{X}} = \tilde{X} (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T$$

$$Y^\perp = (I - P_X) Y; \quad \hat{Y} = P_X Y$$

$$X_k^\perp = (I - P_X) X_k; \quad \hat{X}_k = P_X X_k$$

$$\bullet \quad a = Y^T P_{x_k^\perp} Y; \quad b = Y^T (I - P_X) Y$$

Then we have:

$$\text{① } r_{yx_k, x_1, \dots, x_{k-1}}^2 = \frac{\langle Y^\perp, X_k^\perp \rangle^2}{\|Y^\perp\|_2^2 \|X_k^\perp\|_2^2}$$

$$= \frac{Y^\perp P_{x_k^\perp} Y^\perp}{Y^T (I - P_X) Y}$$

$$= \frac{Y^T P_{x_k^\perp} Y - Y^T P_{x_k^\perp} \hat{Y}}{Y^T (I - P_X) Y}$$

$$\stackrel{*}{=} \frac{a}{b} \quad (1)$$

(*) is because $\hat{Y} \in V(X)$ & $P_{x_k^\perp} \hat{Y} \in V(X)^\perp$

② Test stat for $\beta_k = 0$ is

$$F = (n-k) Y^T (P_{\tilde{X}} - P_X) Y / Y^T (I - P_X) Y$$

$$\stackrel{(\Delta)}{=} (n-k) Y^T P_{x_k^\perp} Y / Y^T (I - P_X) Y$$

$$= (n-k) a / [Y^T (I - P_X) Y - Y^T (P_{\tilde{X}} - P_X) Y]$$

$$\stackrel{(\Delta)}{=} (n-k) a / (b - a) \quad (2)$$

(Δ) are due to $(P_{\tilde{X}} - P_X) Y \in V(\tilde{X})$

Combine (1) & (2):

$$r^2 = F / (n-k+F)$$

