

## Midterm Review

D1 (a) Since we are minimizing SSR,  $\sum_{i=1}^n E_i^2$

$$SS = S(A, B) = \sum_{i=1}^n E_i^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n (Y_i - A - BX_i)^2$$

$$(1) \quad \frac{\partial S(A, B)}{\partial A} = \sum_{i=1}^n (-1)(2) (Y_i - A - BX_i) = 0$$

$$\Rightarrow \sum_{i=1}^n (Y_i - A - BX_i) = 0$$

$$\Rightarrow \sum_{i=1}^n Y_i = \sum_{i=1}^n A + \sum_{i=1}^n BX_i$$

$$\Rightarrow \sum_{i=1}^n Y_i = nA + B \sum_{i=1}^n X_i$$

$$\Rightarrow \frac{1}{n} \sum_{i=1}^n Y_i = A + B \left( \frac{1}{n} \sum_{i=1}^n X_i \right) \Rightarrow \bar{Y} = A + b\bar{X}$$

$$\Rightarrow A = \bar{Y} - b\bar{X}$$

$$(b) \text{ From (1), } \sum_{i=1}^n (-1)(2) (Y_i - A - BX_i) = 0$$

$$\sum_{i=1}^n (Y_i - A - BX_i) = 0 \Rightarrow \sum_{i=1}^n (Y_i - \hat{Y}_i) = 0$$

$$\Rightarrow \sum_{i=1}^n Y_i = \sum_{i=1}^n \hat{Y}_i \Rightarrow \sum Y_i = \sum \hat{Y}_i$$

$$(c) \quad \frac{\partial S(A, B)}{\partial B} = \sum_{i=1}^n (-1)(2) X_i (Y_i - A - BX_i) = 0$$

$$\Rightarrow \sum_{i=1}^n X_i (Y_i - \bar{Y} + b\bar{X} - BX_i) = 0 \Rightarrow \sum_{i=1}^n X_i (Y_i - \bar{Y}) = B \sum_{i=1}^n X_i (X_i - \bar{X})$$

$$\Rightarrow B = \frac{\sum_{i=1}^n X_i (Y_i - \bar{Y})}{\sum_{i=1}^n X_i (X_i - \bar{X})} = \frac{\sum_{i=1}^n (X_i - \bar{X}) (Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\sum_{i=1}^n (X_i - \bar{X}) Y_i}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

$$\Rightarrow B = \sum_{i=1}^n k_i Y_i \quad \text{where } k_i = \frac{X_i - \bar{X}}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

Therefore B is linear combination of  $Y_i$

$$\begin{aligned}
 E[B] &= E\left[\sum_{i=1}^n k_i y_i\right] = \sum_{i=1}^n k_i \cdot E[y_i] = \sum_{i=1}^n k_i (\alpha + \beta x_i) \\
 &= \sum_{i=1}^n \alpha k_i + \sum_{i=1}^n \beta x_i k_i = \alpha \cdot \frac{\sum_{i=1}^n (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} + \beta \cdot \frac{\sum_{i=1}^n (x_i - \bar{x}) x_i}{\sum_{i=1}^n (x_i - \bar{x})^2} \\
 &= 0 \cdot \alpha + \beta \cdot \frac{\sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \beta \cdot 1 = \beta
 \end{aligned}$$

$\therefore B$  is unbiased estimator of  $\beta$

$$\begin{aligned}
 \text{cd) } \text{Reg SS} &= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y} + \bar{y} - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \\
 &= \sum_{i=1}^n (B x_i - B \bar{x})^2 = \sum_{i=1}^n B^2 (x_i - \bar{x})^2 \\
 &= \underbrace{\left(\sum_{i=1}^n (x_i - \bar{x})^2\right)}_{S_{xx}} B^2 = B^2 S_{xx}
 \end{aligned}$$

D2

$x_i$ # of obs class	$y_i$ # of stats class	$x_i y_i$	$x_i^2$	$\hat{y}$	$E_i = y_i - \hat{y}$
1	1	1	1	0.6	0.4
2	1	2	4	1.3	-0.3
3	2	6	9	2	0
4	2	8	16	2.7	-0.7
5	4	20	25	3.4	0.6

$$\bar{x} = \frac{(1+2+3+4+5)}{5} = 3$$

$$\bar{y} = \frac{(1+1+2+2+4)}{5} = 2$$

$$\begin{aligned}
 \text{(a) } \text{Sme } B &= \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2} = \frac{(1+4+9+16+25) - 5 \cdot 3 \cdot 2}{(1+4+9+16+25) - 5 \cdot (3)^2} \\
 &= \frac{37 - 5(12)}{55 - 5(9)} = 0.7
 \end{aligned}$$

$$A = \bar{y} - B \bar{x} = 2 - 0.7 \cdot 3 = 2 - 2.1 = -0.1$$

$$\text{Model is } \hat{y} = -0.1 + 0.7x$$

We calculate  $\hat{y} = -0.1 + 0.7x_i$  and fill in the table.

$$\text{cb) } \text{RSS} = \sum (y_i - \hat{y}_i)^2 = 0.4^2 + 0.3^2 + 0^2 + 0.7^2 + 0.6^2 = 1.1$$

- (c) A: the # of stats class a person took is  $-0.1$  when the person does not take any CS class  
 B: with every unit increase in # of CS class the student take, (one more class in CS)

the person will have  $0.7$  unit increase in Stats classes taken.

A does not make sense since a person cannot take negative number of classes

- (d) Conduct t-test

Null Hypothesis:  $\beta = 0$ , there is no linear relationship

Alternative Hypothesis:  $\beta \neq 0$ , there exist a linear relationship

$$RMS = \frac{RSS}{n-2} = \frac{1.1}{5-2} = \frac{1.1}{3} = 0.3667$$

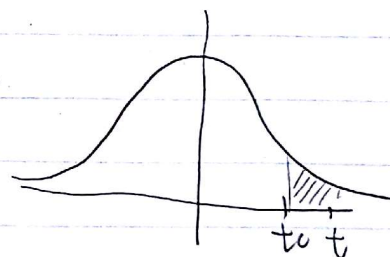
$$t = \frac{B}{\sqrt{\frac{RMS}{S_{xx}}}}$$

$$S_{xx} = \sum_{i=1}^5 x_i^2 - 5\bar{x}^2 = 10$$

$$t = \frac{0.7}{\sqrt{\frac{0.3667}{10}}} = \frac{0.7}{\sqrt{0.03667}} = 3.65546$$

By  $R_1$  we know  $t_{c, \alpha=0.025, df=3} = 3.18245$

Since  $t = 3.65546 > t_c = 3.18245$



Therefore we reject the null hypothesis, there is a linear relationship between # of CS class taken and # of Stats classes taken

- (e) ANOVA Table

sources of variation	sum of squares	df	mean square	F
Regression	$RegSS = 4.9$	1	$RegMS = \frac{RegSS}{1} = 4.9$	$F = \frac{RegMS}{RMS} = \frac{4.9}{0.3667} = 13.3624$
Residual	$RSS = 1.1$	3	$RMS = 0.3667$	
Total	$TSS = 6$			

$$\begin{aligned} 1) RegSS &= B^2 S_{xx} \\ &= 0.7^2 \cdot 10 \\ &= 4.9 \end{aligned}$$



(f)  $\therefore F_{\text{crit}} = 10.13$  for  $\alpha = 0.05$

Since  $F = 13.3624 > F_c = 10.13$

Thus, we reject  $H_0$  and there is a linear relationship between  
# of CS classes taken and # of stats classes taken

(g) Yes, same result from d and f

(h)  $t^2 = (3.165546)^2 = 10.0206 = F$

They are approximately equal

(i)  $r^2 = \frac{\text{RegSS}}{\text{TSS}} = \frac{4.9}{6} = 0.81667 \approx 0.82$

Therefore, 82% of the variation in the # stats classes taken  
can be explained by the # of CS classes taken by our  
model.

(j) Since  $t_{0.025, df=3} = 3.18245$

Since confidence interval for  $\beta$  is  $[B - t_{\alpha/2, df} \sqrt{\frac{\text{PMS}}{\text{Sxx}}}, B + t_{\alpha/2, df} \sqrt{\frac{\text{PMS}}{\text{Sxx}}}]$

$t_{\alpha/2, df} \sqrt{\frac{\text{PMS}}{\text{Sxx}}} = t_{0.025, df=3} \sqrt{\frac{0.3667}{10}} = 3.18245 \sqrt{\frac{0.3667}{10}} = 0.60942$

$\Rightarrow 0.7 - 0.60942 \leq \beta \leq 0.7 + 0.60942$

$\Rightarrow 0.09058 \leq \beta \leq 1.30942 \Rightarrow 0.1 \leq \beta \leq 1.3$

$\Rightarrow$  95% confidence interval for  $\beta$  is  $[0.1, 1.3]$

3) F-test on poverty & minority for crime.  $\alpha = 0.05$

(a)  $H_0: \beta_1 = \beta_2 = 0$

$H_A$ : at least one of the  $\beta_1, \beta_2$  is not 0

from the code output given,  $p\text{-value} = 4.249 \times 10^{-9} < \alpha$

Then p-value is significant. We reject  $H_0$ , and conclude that at least one of poverty and minority will affect crime rate.

(b) t-test on poverty & crime  $\alpha = 0.05$

$$H_0: \beta_1 = 0$$

$$H_A: \beta_1 \neq 0$$

From code output, we have p-value is  $0.075 > \alpha$ .

p-value is not significant. We fail to reject  $H_0$  and conclude that poverty itself does not significantly affect crime.

(c) t-test on minority & crime  $\alpha = 0.05$

$$H_0: \beta_2 = 0$$

$$H_A: \beta_2 \neq 0$$

From code output, we have p-value is  $6.86 \times 10^{-8} < \alpha$

p-value is significant, we reject the  $H_0$  and conclude that minority itself significantly affects crime.

(d) We can see that adjusted R-square =  $0.4405$ , and t-stat value is fairly low, thus there is not a strong linear relationship between minority & poverty on crime. Also, we can see that residual squared error is  $18.62$  and the fitted model is not very accurate.

(e)  $\hat{R}^2 = 0.4405 = 44.05\%$

Thus  $44.05\%$  variation in the crime can be attributed to poverty and minority.

(f) I would use R-adjusted square instead of R-square since we  $\hat{R}^2$  can eliminate the artificial inflation in  $R^2$  when the explanatory variables are co-linear correlated.



$$(g) \hat{y} = 58.670 + 1.195x$$

$\uparrow$  Crime                       $\uparrow$  minority

No poverty since poverty is not significant.

(h) When  $x = \text{minority} = 25$

$$\hat{y} = 58.670 + 1.195(25) = 88.545$$

(i) Since calculated from the original dataset, we can calculate the table below (Covariance table)

	Crime	Poverty	Minority
Crime	1	0.56	0.65
Poverty	0.36	1	0.73
Minority	0.65	0.73	1

Yes, since we can see that

① The Covariance between Minority and Poverty is very high

② Since Covariance between poverty and Crime is very low, they are not very correlated, however since Crime and minority has high covariance, thus the high covariance between poverty and minority may cause the increase in the correlation between poverty and Crime

Here: Minority  $\xrightarrow{\text{affect}}$  Poverty  $\xrightarrow{\text{affect}}$  Crime

Thus, there is pair-wise correlation between minority and poverty.