

Part A.

$$\begin{matrix} b & d & c & a & b \\ a & d & d & b & c \end{matrix}$$

$$\begin{matrix} a & a & a & a & c \\ b & a & c & c & b \end{matrix}$$

2 ii)  $\begin{pmatrix} \bar{Y} \\ \vdots \\ \bar{Y} \end{pmatrix}_{n \times 1}$

iii)  $RSS = \sum_{i=1}^n (Y_i - \bar{Y})^2 = (Y - HY)'(Y - HY) = (Y' - Y'H')(Y - HY) = Y'Y - Y'HY - Y'H'Y + Y'H'HY$   
 $\frac{H' = H}{H'H = H} \quad Y'(I - H)Y$

iv)  $J'J = \begin{bmatrix} \frac{1}{n} & \cdots & \frac{1}{n} \\ \vdots & \ddots & \vdots \\ \frac{1}{n} & \cdots & \frac{1}{n} \end{bmatrix}^2 = J \quad J' = J$

$$SYY = \sum_{i=1}^n (Y_i - \bar{Y})^2 = (Y - JY)'(Y - JY) = Y'Y - Y'JY - Y'J'Y + Y'J'JY = Y'(I - J)Y$$

v)  $HX = X \Rightarrow H \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}_{n \times 1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}_{n \times 1} \Rightarrow H \begin{pmatrix} 1 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \cdots & 1 \end{pmatrix}_{n \times n} = \begin{pmatrix} 1 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \cdots & 1 \end{pmatrix}_{n \times n} \Rightarrow HJ = J$

Similarly

$$X'H = X' \Rightarrow JH = J$$

vi).  $\hat{Y} = HY \quad J\hat{Y} = JHY = JY = \begin{pmatrix} \bar{Y} \\ \vdots \\ \bar{Y} \end{pmatrix}_{n \times 1}$

also  $J\hat{Y} = \begin{pmatrix} \frac{1}{n} \sum \hat{Y}_i / n \\ \vdots \\ \frac{1}{n} \sum \hat{Y}_i / n \end{pmatrix}_{n \times 1} \Rightarrow \bar{\hat{Y}} = \frac{1}{n} \sum \hat{Y}_i$

vii).  $SS_{reg} = \sum_{i=1}^n (\hat{Y}_i - \bar{\hat{Y}})^2 = (HY - JY)'(HY - JY) = Y'HY - Y'HJY - Y'JHY + Y'JY$   
 $= Y'HY - Y'JY - Y'JY + Y'JY = Y'(H - J)Y$

viii)  $RSS + SS_{reg} = Y'(I - H)Y + Y'(H - J)Y = Y'(I - J)Y = SYY$

19. For a simple linear regression model,

$$Y = \beta_0 + \beta_1 x + e \quad e_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$$

Given the following  $X$  observations  $\mathcal{A}_X = \{5, 3, 5, 4\}$  and  $Y$  observations  $\mathcal{A}_Y = \{9, 3, 9, 4\}$ . If Helen wants to test the null hypothesis against the alternative hypothesis at 10% significance level,

$$H_0 : \beta_0 = 0 \quad H_1 : \beta_0 < 0$$

then what are the correct p-value and conclusion?

- a)  $p = 0.2994$  and cannot reject Null
  - b)  $p = 0.1472$  and cannot reject Null
  - c)  $p = 0.0736$  and reject Null
  - d)  $p = 0.1472$  and reject Null
20. For the multiple regression model  $Y = X\beta + e$ ,  $e \sim N(0, \sigma^2)$ , Let  $\hat{\beta} = (X'X + kI)^{-1}X'Y$ . What is  $E(\hat{\beta})$ ?
- a)  $\beta$
  - b)  $\beta - k(X'X + kI)^{-1}\beta$
  - c)  $\beta - k(X'X + kI)\beta$
  - d)  $\beta - kX'X\beta$

Part B: Write down the calculation steps and answers on the blank area. (50 Marks)

1. (5 marks) In a multiple regression analysis,  $MSE = 30$ , number of observation = 58, number of predictor = 4, and  $TSS = 2000$ . What is the  $R^2$  of the regression?

$$R^2 = \frac{S_{reg}}{S_{YY}} = \frac{S_{YY} - MSE \cdot (n - (p+1))}{S_{YY}} = 0.205$$

3. (25 marks) The results of thirteen countries in the 2014 Asian Olympic game is given below:

Country	Gold (G)	Silver (S)	Bronze (B)
China	151	108	83
Korea	79	71	84
Japan	47	76	77
Kazakhstan	28	23	33
Iran	21	18	18
Thailand	12	7	28
North Korea	11	11	14
India	11	10	36
Taiwan	10	18	23
Qatar	10	0	4
Uzbekistan	9	14	21
Bahrain	9	6	4
Hong Kong	6	12	24

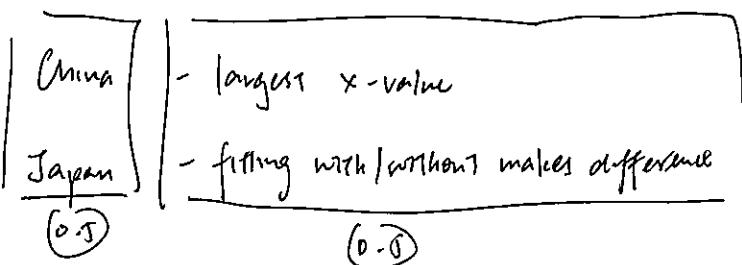
i) (5 marks) Plot a scatterplot of  $G$  against  $S$ . Which country would have the most influence on the regression fit?

P101

$$X\text{-axis} = S \quad \textcircled{1}$$

$$Y\text{-axis} = G \quad \textcircled{1}$$

points  $\quad \textcircled{2}$



ii) (5 marks) It is found that  $\sum XY = 15614.76$ ,  $\sum X^2 = 12520.78$ ,  $\sum Y^2 = 20686.73$

$$\frac{1}{n} \sum G_i = 31.1, \frac{1}{n} \sum S_i = 28.8, \frac{1}{n} \sum G_i^2 = 2558.5, \frac{1}{n} \sum S_i^2 = 1869.5, \frac{1}{n} \sum S_i G_i = 2096.8$$

Use these results to estimate the regression model  $G = \beta_0 + \beta_1 S + e$ ,  $e \sim N(0, \sigma^2)$ . Plot this line on the above scatter plot diagram. Is this fitting appropriate?

$$\hat{\beta}_1 = 1.1549 \quad \textcircled{1}$$

$$\hat{\beta}_0 = -2.1599 \quad \textcircled{1}$$

$$\hat{G} = 241.29642 \quad \textcircled{1}$$

P101  
Fitting Appropriate?  
 $\textcircled{1}$

- iii) (5 marks) Is there linear relationship between the number of gold and silver medals at 5% significance level?

$$F = \frac{\text{SSreg}/1}{\text{MSres}/(n-2)} = 74.73390 > 4.84 = F_{0.05}(1, 11)$$

$\therefore$  Regression is effective / H<sub>0</sub> is reg.

Or  $|t| = 8.64488 > 2.201 = t_{0.05}(11)$

$\wedge$   
formula

- iv) (5 marks) Write down the equation governing the 95% confidence ellipse for the parameter  $(\beta_0, \beta_1)$ .

$$\frac{(\hat{\beta} - \beta)'(X'X)(\hat{\beta} - \beta)}{2\hat{\sigma}^2} \leq 3.982 = F_{0.05}(2, 11)$$

$$\beta_0^2 + 1869.5\beta_1^2 + 57.6\beta_0\beta_1 - 62.20\beta_0 - 4193.76\beta_1 + 432.86 \leq 0$$

- v) (5 marks) If Country A got 20 silver medals, what is the 90% prediction interval of its number of gold medal?

$$\hat{\sigma} = 15.53353, \text{ sepred}(y|x^*) = \hat{\sigma} \sqrt{1 + \frac{1}{13} + \frac{(x^* - \bar{x})^2}{SSS}} = 1.0405 \cdot 15.53353$$

$$= 16.16271 \quad \textcircled{1}$$

$$\tilde{Y}_x = -2.1599 + 1.1549 \cdot 20 = 20.9381 \quad \textcircled{1}$$

$$\tilde{Y}_x \pm t_{0.05}(11) \text{ sepred}(y|x^*) \quad \textcircled{1}$$

$$= 20.9381 \pm 1.796 \cdot 16.16271.$$

$$= 20.9381 \pm 29.02823$$

$$= [ \quad , \quad 49.9663 ] \quad \textcircled{1}$$

End of paper