

《迴歸分析》

試題評析

本年度迴歸分析只考四題，第一題為簡單迴歸，應屬送分題。第二題為複迴歸；缺截距項，求OLSE與 $E(\hat{\beta}_1)$ 、 $V(\hat{\beta}_1)$ 的結果，只要將迴歸模型含截距項的結果中， \bar{X} 、 \bar{Y} 代0即可。第三題是比較複迴歸之ANOVA F-test與偏F-test之差異性。中上程度同學，二、三題應該都可以拿到分數。至於第四題為確認離群值之後，檢測離群值是否對迴歸模型具有影響力，屬於比較冷僻題型，公式難記、不易拿分。

一、今對一組樣本資料 (x_i, y_i) ， $i=1, \dots, 20$ 適配一簡單線性迴歸模型 $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$ ，其中 ε_i 為 i.i.d.N(0, σ^2)。已知 $F_{1,18,0.05} = 3.01$ ， $\bar{x} = 6$ ， $\bar{y} = 15$ ， $\sum(x_i - \bar{x})^2 = 25$ ， $\sum(y_i - \bar{y})^2 = 208$ ， $\sum(x_i - \bar{x})(y_i - \bar{y}) = 40$ 。

- (一) 試寫此模型之變異數分析 (Analysis of Variance) 表。(10分)
- (二) 試求此一迴歸線之斜率與截距。(10分)
- (三) 試求此一迴歸線斜率 β_1 之 90% 信賴區間。(5分)

答：

(一)

$$SSTO = S_{YY} = 208$$

$$SSR = \frac{S_{XY}^2}{S_{XX}} = \frac{40^2}{25} = 64$$

$$SSE = SSTO - SSR = 144$$

ANOVA table

變異來源	平方和	自由度	均方	F值
迴歸	64	1	64	8
誤差	144	18	8	
總和	208	19		

(二)

$$(1) \text{ 斜率 } \hat{\beta}_1 = \frac{S_{XY}}{S_{XX}} = \frac{40}{25} = 1.6$$

$$(2) \text{ 截距 } \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 15 - 1.6 \times 6 = 5.4$$

(三)

β_1 之 90% C.I.

$$\begin{aligned} &= \hat{\beta}_1 \pm t_{0.1} \left(\frac{20-2}{2} \right) \cdot \sqrt{\frac{MSE}{S_{XX}}} \\ &= 1.6 \pm \sqrt{3.01} \times \sqrt{\frac{8}{25}} = 1.6 \pm 0.9814 = (0.6186, 2.5814) \end{aligned}$$

二、若 (x_{1i}, x_{2i}, y_i) ， $i=1, \dots, n$ 彼此獨立且來自截距項為零的線性迴歸模型 $Y_i = \beta_1 X_{1i} + \beta_2 X_{2i} + \varepsilon_i$ ，其中 ε_i 為 i.i.d.N(0, σ^2)。

-- 1 --

(一) 試求 β_1 與 β_2 之最小平方估計量 $\hat{\beta}_1$ 與 $\hat{\beta}_2$ 。 (15分)

(二) 試求 $E(\hat{\beta}_1)$ 與 $\text{Var}(\hat{\beta}_2)$ 。 (10分)

答：

(一)

$$\text{令 } Q = \sum_{i=1}^n (Y_i - \hat{\beta}_1 X_{1i} - \hat{\beta}_2 X_{2i})^2$$

$$\begin{cases} \frac{\partial Q}{\partial \hat{\beta}_1} = -2 \sum_{i=1}^n [(Y_i - \hat{\beta}_1 X_{1i} - \hat{\beta}_2 X_{2i}) \cdot X_{1i}] = 0 \\ \frac{\partial Q}{\partial \hat{\beta}_2} = -2 \sum_{i=1}^n [(Y_i - \hat{\beta}_1 X_{1i} - \hat{\beta}_2 X_{2i}) \cdot X_{2i}] = 0 \end{cases}$$

$$\begin{cases} \sum_{i=1}^n X_{1i}^2 \cdot \hat{\beta}_1 + \sum_{i=1}^n X_{1i} X_{2i} \cdot \hat{\beta}_2 = \sum_{i=1}^n X_{1i} Y_i \\ \sum_{i=1}^n X_{1i} X_{2i} \cdot \hat{\beta}_1 + \sum_{i=1}^n X_{2i}^2 \cdot \hat{\beta}_2 = \sum_{i=1}^n X_{2i} Y_i \end{cases}$$

$$\hat{\beta}_1 = \frac{\begin{vmatrix} \sum_{i=1}^n X_{1i} Y_i & \sum_{i=1}^n X_{1i} X_{2i} \\ \sum_{i=1}^n X_{2i} Y_i & \sum_{i=1}^n X_{2i}^2 \end{vmatrix}}{\begin{vmatrix} \sum_{i=1}^n X_{1i}^2 & \sum_{i=1}^n X_{1i} X_{2i} \\ \sum_{i=1}^n X_{1i} X_{2i} & \sum_{i=1}^n X_{2i}^2 \end{vmatrix}}$$

$$= \frac{(\sum_{i=1}^n X_{2i}^2)(\sum_{i=1}^n X_{1i} Y_i) - (\sum_{i=1}^n X_{1i} X_{2i})(\sum_{i=1}^n X_{2i} Y_i)}{(\sum_{i=1}^n X_{1i}^2)(\sum_{i=1}^n X_{2i}^2) - (\sum_{i=1}^n X_{1i} X_{2i})^2}$$

同理

$$\hat{\beta}_2 = \frac{(\sum_{i=1}^n X_{1i}^2)(\sum_{i=1}^n X_{2i} Y_i) - (\sum_{i=1}^n X_{1i} X_{2i})(\sum_{i=1}^n X_{1i} Y_i)}{(\sum_{i=1}^n X_{1i}^2)(\sum_{i=1}^n X_{2i}^2) - (\sum_{i=1}^n X_{1i} X_{2i})^2}$$

(二)

$$(1) E\hat{\beta}_1 = E \left[\frac{(\sum_{i=1}^n X_{2i}^2)(\sum_{i=1}^n X_{1i} Y_i) - (\sum_{i=1}^n X_{1i} X_{2i})(\sum_{i=1}^n X_{2i} Y_i)}{(\sum_{i=1}^n X_{1i}^2)(\sum_{i=1}^n X_{2i}^2) - (\sum_{i=1}^n X_{1i} X_{2i})^2} \right]$$

$$= \frac{(\sum_{i=1}^n X_{2i}^2)(\sum_{i=1}^n X_{1i}EY_i) - (\sum_{i=1}^n X_{1i}X_{2i})(\sum_{i=1}^n X_{2i}EY_i)}{(\sum_{i=1}^n X_{1i}^2)(\sum_{i=1}^n X_{2i}^2) - (\sum_{i=1}^n X_{1i}X_{2i})^2}$$

$$\begin{aligned} & \frac{(\sum_{i=1}^n X_{2i}^2)(\beta_1 \sum_{i=1}^n X_{1i}^2 - \beta_2 \sum_{i=1}^n X_{1i}X_{2i}) - (\sum_{i=1}^n X_{1i}X_{2i})(\beta_1 \sum_{i=1}^n X_{1i}X_{2i} - \beta_2 \sum_{i=1}^n X_{2i}^2)}{(\sum_{i=1}^n X_{1i}^2)(\sum_{i=1}^n X_{2i}^2) - (\sum_{i=1}^n X_{1i}X_{2i})^2} \\ &= \frac{\beta_1 \left[(\sum_{i=1}^n X_{2i}^2)(\sum_{i=1}^n X_{1i}^2) - (\sum_{i=1}^n X_{1i}X_{2i})^2 \right]}{(\sum_{i=1}^n X_{1i}^2)(\sum_{i=1}^n X_{2i}^2) - (\sum_{i=1}^n X_{1i}X_{2i})^2} \\ &= \beta_1 \end{aligned}$$

$$\begin{aligned} (2) \quad V(\hat{\beta}_1) &= V \left[\frac{(\sum_{i=1}^n X_{2i}^2)(\sum_{i=1}^n X_{1i}Y_i) - (\sum_{i=1}^n X_{1i}X_{2i})(\sum_{i=1}^n X_{2i}Y_i)}{(\sum_{i=1}^n X_{1i}^2)(\sum_{i=1}^n X_{2i}^2) - (\sum_{i=1}^n X_{1i}X_{2i})^2} \right] \\ &= \frac{(\sum_{i=1}^n X_{2i}^2)^2(\sigma^2 \sum_{i=1}^n X_{1i}^2) - (\sum_{i=1}^n X_{1i}X_{2i})^2(\sigma^2 \sum_{i=1}^n X_{2i}^2)}{\left[(\sum_{i=1}^n X_{1i}^2)(\sum_{i=1}^n X_{2i}^2) - (\sum_{i=1}^n X_{1i}X_{2i})^2 \right]^2} \\ &= \frac{\sigma^2 \cdot \sum_{i=1}^n X_{2i}^2 \left[(\sum_{i=1}^n X_{2i}^2)(\sum_{i=1}^n X_{1i}^2) - (\sum_{i=1}^n X_{1i}X_{2i})^2 \right]}{\left[(\sum_{i=1}^n X_{1i}^2)(\sum_{i=1}^n X_{2i}^2) - (\sum_{i=1}^n X_{1i}X_{2i})^2 \right]^2} \\ &= \frac{\sum_{i=1}^n X_{2i}^2}{(\sum_{i=1}^n X_{1i}^2)(\sum_{i=1}^n X_{2i}^2) - (\sum_{i=1}^n X_{1i}X_{2i})^2} \sigma^2 \end{aligned}$$

三、某一研究想要知道房屋的價格Y與房屋的坪數 X_1 ，屋齡 X_2 ，房間數 X_3 與空屋率 X_4 的關係。今收集30間房屋的資料並對此資料配適一迴歸模型 $Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + \beta_4 X_{4i} + \varepsilon_i$ ，其中 ε_i 為 i.i.d. $N(0, \sigma^2)$ 。已知 $F_{2,27,0.05} = 3.35$ ， $F_{2,25,0.05} = 3.39$ 。

Source of Variation(變異來源)	SS	df	MS
SSR(X_1, X_2, X_3, X_4)	300	4	75
SSR($X_3, X_4 X_1, X_2$)	40	2	20
SSR(X_3, X_4)	150	2	75
SSE(X_1, X_2, X_3, X_4)	200	25	8

(一)假設迴歸模型中僅考慮房間數 X_3 與空屋率 X_4 。試就此模型在 $\alpha = 0.05$ 下檢定 $H_0: \beta_3 = \beta_4 = 0$ 。

(請務必將完整之檢定寫出，包括 $H_0, H_1, 檢定量, 拒絕區域, 結論等$) (10分)

(二)假設迴歸模型中已考慮坪數 X_1 與屋齡 X_2 。試就此模型在 $\alpha = 0.05$ 下檢定 $H_0: \beta_3 = \beta_4 = 0$ 。

(請務必將完整之檢定寫出，包括 $H_0, H_1, 檢定量, 拒絕區域, 結論等$) (10分)

(三)試解釋上面兩小題結果不盡相同之原因。(5分)

答：

(一)model : $Y_i = \beta_0 + \beta_3 X_{3i} + \beta_4 X_{4i} + \varepsilon_i$

$$\begin{aligned} SSTO &= SSR(X_1 X_2 X_3 X_4) + SSE(X_1 X_2 X_3 X_4) \\ &= 300 + 200 = 500 \end{aligned}$$

$$\begin{aligned} SSE(X_1 X_2) &= SSTO - SSR(X_1 X_2) \\ &= 500 - 150 = 350 \end{aligned}$$

1.檢定假設

$$\begin{cases} H_0: \beta_3 = \beta_4 = 0 \\ H_1: \text{不全為0} \end{cases}$$

2.拒絕區域

$$C\{F | F > F_{0.05}(2, 30-3) = 3.35\}$$

3.檢定統計值

$$F = \frac{\frac{SSR(X_1 X_2)}{2}}{\frac{SSE(X_1 X_2)}{(30-3)}} = \frac{\frac{75}{2}}{\frac{350}{27}} = 5.78 \in C$$

\therefore reject H_0 , 有充分證據顯示 β_3, β_4 不全為0

(二)model : $Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + \beta_4 X_{4i} + \varepsilon_i$

1.檢定假設

$$\begin{cases} H_0: \beta_3 = \beta_4 = 0 \\ H_1: \text{不全為0} \end{cases} \Leftrightarrow \begin{cases} H_0: \text{不值得引進 } X_3, X_4 \\ H_1: \text{值得引進 } X_3, X_4 \end{cases}$$

2.拒絕區域

$$C\{F | F > F_{0.05}(2, 30-5) = 3.39\}$$

3.檢定統計值

$$F = \frac{\frac{SSR(X_3 X_4 | X_1 X_2)}{2}}{\frac{SSE(X_1 X_2 X_3 X_4)}{(30-5)}} = \frac{\frac{20}{2}}{\frac{8}{8}} = 2.5 \notin C$$

\therefore Not reject H_0 , 無充分證據顯示 β_3, β_4 不全為0

(三)

1. 在(一)中，表示房屋的價格可以用房間數(X_3)與空屋率(X_4)來預測
2. 在(二)中，表示當model已含有 X_1 (坪數)， X_2 (屋齡)兩個自變數時，不值得引進 X_3 (房屋數)與 X_4 (空屋率)當Y之預測變數。
 \therefore (一)採用迴歸ANOVA F-test與(二)偏F檢定結果未必相同

四、在模型診斷時，我們常用DFFITS, Cook's Distance, DFBETAS方法辨認具有影響力的個案 (Influential cases)。

(一)試比較DFFITS, Cook's Distance, DFBETAS此三種方法之差異。(15分)

(二)試說明此三種方法辨認具有影響力的個案之判定原則。(10分)

答：

(一)

1. DFFITS：剔除ith個離群值後，對單一適配值(fitted value)的影響。

$$\text{其中}, (DFFITS)_i = \frac{\hat{Y}_i - \hat{Y}_{i(i)}}{\sqrt{MSE_{(i)} h_{ii}}}$$

2. Cook's Distance：剔除ith個離群值後，對所有適配值(fitted value)的影響。

$$\text{其中}, D_i = \frac{\sum_{j=1}^n (\hat{Y}_j - \hat{Y}_{j(i)})^2}{pMSE}$$

3. DFBETAS：剔除ith個離群值後，對於每一個迴歸係數(regression coefficients)的影響。

$$\text{其中}, (DFBETAS)_{k(i)} = \frac{b_k - b_{k(i)}}{\sqrt{MSE_{(i)} c_{kk}}} \quad k = 0, 1, \dots, p-1$$

(二)

1. DFFITS判定原則：

對於中、小規模的資料集而言，當 $|DFFITS| > 1$ ；或者是，大規模資料集而言， $|DFFITS| > 2\sqrt{\frac{p}{n}}$ ，

則可視為離群值具有影響力。

2. Cook's Distance判定原則：

Cook's Distance中， D_i 與 $F(p, n-p)$ 所對應的百分位數值來確認個別的影響。當百分位數值小於10%或是20%時，表示第*i*個個案對於適配值影響不大；當百分位數值大於50%時，表示第*i*個個案對於適配迴歸函數具有影響。

3. DFBETAS判定原則：

對於中、小規模的資料集而言，當 $DFBETAS > 1$ ；或者是，大規模資料集而言，

$$DFBETAS > \frac{2}{\sqrt{n}} \text{，則可視為離群值具有影響力。}$$