

# 《迴歸分析》

試題評析	<p>本年考題需要相當熟悉理論架構，計算量較去年少，但申論問答的部分又較去年多，考題難度約略中上，程度好的同學應可拿80分以上，但較難拿到90分以上。</p> <p>第一題：考判定係數的評述，需就簡單迴歸與複迴歸加以分述才可以完整拿分，完整表達絕對是是否能全拿分的關鍵。屬於須考量全部觀點的申論題。</p> <p>第二題：考複迴歸分析，模型設定及解釋係數的能力，且要熟悉殘差分析才有辦法設定正確模型。</p> <p>第三題：考迴歸模型診斷，是正課與總複習再三強調的重點，故熟讀就能拿分。</p> <p>第四題：考報表綜合題型，須熟悉額外平方和的概念輔以報表判讀解釋能力，才有辦法完整回答。</p>
考點命中	<p>1.《高點·高上迴歸分析講義》第一回，許勝雄編撰，頁36-37實例說明。 《高點·高上迴歸分析講義》第二回，許勝雄編撰，頁58-59定理說明。</p> <p>2.《高點·高上迴歸分析講義》第三回，許勝雄編撰，與頁52例題7相同。</p> <p>3.《高點·高上迴歸分析講義》第三回，許勝雄編撰，頁1-2、81；《第二回》頁33相同。</p> <p>4.《高點·高上迴歸分析講義》第三回，許勝雄編撰，與頁48例題5相同。</p>

一、迴歸分析中分別對下列陳述做一評述：

(一)如  $R^2$  (coefficient of determination 判定係數) 大 (譬如 0.95)，則此模型良好，應採用。(6分)

(二)如  $R^2$  (coefficient of determination 判定係數) 小 (譬如 0.35)，則此模型不佳，不應採用。(6分)

**答：**

(一)若在簡單迴歸分析中，當  $R^2 = \frac{SSR}{SSTO} = \frac{\text{可解釋變異}}{\text{總變異}}$ ，若  $R^2$  (判定係數) 大，

則 F 統計量  $= \frac{(n-2)R^2}{1-R^2}$ ，F 越顯著，則代表 ANOVA-F 檢定顯著，意味模型配適良好；而在複迴歸中，由於

自變數越多，會使得複判定係數  $R^2 = \frac{SSR}{SSTO} = \frac{\text{可解釋變異}}{\text{總變異}}$  越大，故自變數個數增加，複判定係數就上升，

但無法確認引進的變數是否皆為有效解釋變數，而在複迴歸中，由 F 統計量  $= \frac{(n-k)R^2}{(k-1)(1-R^2)}$  可知，當若要 F

檢定量夠大而拒絕虛無假設時，仍須考量樣本數、自變數個數與其他的因素方可判定，故複迴歸分析則須參考修正判定係數  $R_a^2 = \bar{R}^2 = 1 - \frac{n-1}{n-k}(1-R^2)$  為判定複迴歸的模型良好與否指標，避免引進過多自變數而令複判定係數變大造成無法辨別引進自變數是否正確。

(二)若在簡單迴歸分析中，當  $R^2 = \frac{SSR}{SSTO} = \frac{\text{可解釋變異}}{\text{總變異}}$ ，若  $R^2$  (判定係數) 小

代表此一模型 SSE 較高，配適較差，則 F 統計量  $= \frac{(n-2)R^2}{1-R^2}$  所計算的檢定統計量會越不顯著(越小)，模型相

對比較不佳(容易無法拒絕虛無假設)，但實務上仍需端視所研究的主題而來判讀而定；就複迴歸的觀點而言，

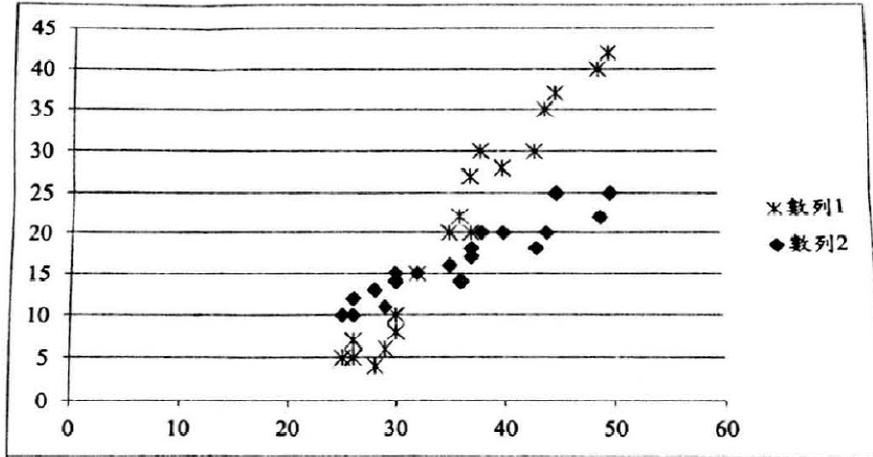
由於 F 統計量為由 F 統計量  $= \frac{(n-k)R^2}{(k-1)(1-R^2)}$  可知，則  $R^2$  (判定係數) 小使 F 檢定量夠而拒絕虛無假設時，仍須考

量樣本數、自變數個數其他的因素方可判定，故複迴歸分析仍需參考修正判定係數  $R_a^2 = \bar{R}^2 = 1 - \frac{n-1}{n-k}(1-R^2)$

【版權所有，重製必究！】

$R^2$ )來判定複迴歸模型優劣為宜。

二、某公司兩業務單位 (A, B) 之員工業績 (Y) 及年齡 (X) 散佈圖如下，數列1為單位A，數列2為單位B。欲觀察年齡與所屬業務單位如何影響員工業績，請以員工業績 (Y) 為應變數設一個複迴歸模式，並解釋模式中每個迴歸係數之涵意。(24分)



答：

由題意，建立以下複迴歸模式：

$$Y_i = \beta_0 + \beta_1 x_i + \beta_2 D_i + \beta_3 X_i D_i + \varepsilon_i, \quad \text{其中 } D_i = \begin{cases} 1, & \text{單位A} \\ 0, & \text{單位B} \end{cases}$$

當  $D_i = 0$  時代表單位B，模式變為  $Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ ， $\beta_0$  為員工基本業績， $\beta_1$  為單位B每單位年齡變動造成員工業績的影響數。

當  $D_i = 1$  時代表單位A，模式變為  $Y_i = \beta_0 + \beta_1 x_i + \beta_2 D_i + \beta_3 X_i D_i + \varepsilon_i$ ， $\beta_0$  為員工基本業績， $\beta_1$  為單位A每單位年齡變動造成員工業績的影響數。

$\beta_2$  代表每增加一單位年齡，所能增加單位A員工業績， $\beta_3$  代表單位A與單位B平均員工基本業績差。

三、在複迴歸模型診斷中，

- (一)「某人對應變數 (Y) 做常態假設之檢定，發現非常態，故採Y之變數變換處理」，評論之。(3分)
- (二)如何觀察各解釋變數與應變數之線性假設是否成立？(3分)
- (三)VIF (變異數膨脹係數) 值過大，表示為何？(3分)
- (四)某  $h_{ii} = 0.7$ ，表示為何？(3分)

答：

(一)若應變數不符合常態分配，則實務上會使用以下轉換：

**Box - Cox Transformation :**

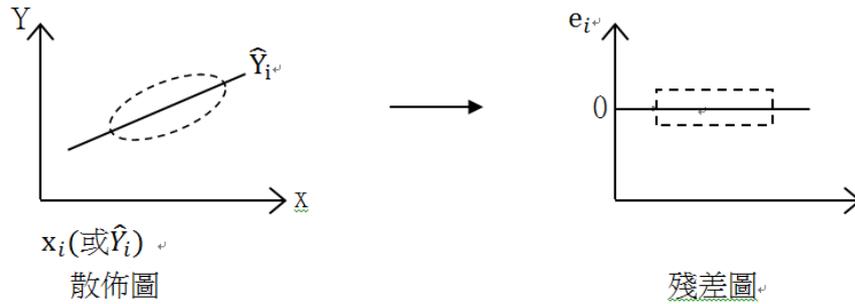
$$\text{令 } T(\varepsilon) = \begin{cases} \frac{\varepsilon^\lambda - 1}{\lambda}, & \text{if } \lambda \neq 0 \\ \ln \varepsilon, & \text{if } \lambda = 0 \end{cases} \quad \text{where: } \lim_{\lambda \rightarrow 0} \frac{\varepsilon^\lambda - 1}{\lambda} = \ln \varepsilon, \text{ 求 } \lambda \text{ 之 MLE}$$

當這些假設成立時，才開始進行參數的估計與檢定。

(二)若要「觀察」解釋變數與應變數之線性假設是否成立，則通常會以殘差圖進行分析

直線假設: modal:  $Y = \beta_0 + \beta_1 x + \varepsilon$ ，若以0為中心線水平帶，正負都沒有規則性傾向，則線性假設應成立。

【版權所有，重製必究！】



(三)變異數膨脹係數 =  $VIF_j = \frac{1}{1 - R_j^2}$ ,  $1 \leq VIF < \infty$ ,  $R_j^2$  代表為共線性模型之複判定係數,  $VIF$  越大代表模

型存在共線性問題, 會造成  $OLSE$  估計值的變異數膨脹, 進而使估計產生誤差。

(四)  $h_{ii}$ : 稱為槓桿值 (亦為帽子矩陣對角線元素), 係用來衡量第  $i$  個觀測值與預測值的距離。

槓桿值有兩個主要的性質:  $0 \leq h_{ii} \leq 1$ ,  $\sum h_{ii} = k - 1$  where  $k =$  參數個數

又  $cov(\mathbf{e}) = \sigma^2(\mathbf{I} - \mathbf{H})$ ,  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$ ,  $h_{ii} = \mathbf{X}_i(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}_i^T$ , 若  $h_{ii} = 0.7$  則

依照  $t$  化殘差圖定義:  $t_i = \frac{e_i - \bar{e}}{S_{e_i}} = \frac{e_i}{S_{e_i}}$  ( $\because \sum e_i = 0 \Rightarrow \bar{e} = 0$ )

Where  $S_{e_i} = \sqrt{\left[1 - \frac{1}{n} - \frac{(x_i - \bar{x})^2}{S_{XX}}\right] \cdot MSE} = \sqrt{(1 - h_{ii}) \cdot MSE}$

也可知, 當  $h_{ii} = 0.7$  使得殘差的變異數變小, 使得  $t$  化殘差變大, 較容易造成離群值的狀況。

四、針對某公司員工, 得一迴歸分析如下:

模型  $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_M M + \beta_R R + \varepsilon$ 。

$Y =$  量化工作績效,  $X_1 =$  年資 (以年為單位),  $X_2 =$  薪資 (以萬元為單位),  $X_3 =$  所屬小組之人數,  $(M, R) = (1, 0)$  為行銷部門,  $(M, R) = (0, 1)$  為研發部門,  $(M, R) = (0, 0)$  為行政部門。

應變數:  $Y$

使用的觀測值數目: 70

$SSTO = 36.304$

$SSR(X_1) = 14.288$ ,  $SSR(X_2 | X_1) = 0.676$ ,  $SSR(X_3 | X_1, X_2) = 5.766$ ,

$SSR(M | X_1, X_2, X_3) = 1.212$ ,  $SSR(R | X_1, X_2, X_3, M) = 0.357$

參數估計值

變數	參數估計值	標準誤差	t 值	Pr >  t
Intercept	0.535	0.491	1.09	0.2793
X1	0.029	0.006	4.84	<.0001
X2	-0.180	0.107	-1.68	0.0951
X3	0.044	0.011	4.00	0.0001
M	0.433	0.163	2.66	0.0099
R	0.175	0.137	1.28	0.2061

(一) 如薪資 ( $X_2$ ) 之單位由萬元改為千元, 上述提供之資料有那些不會改變? 有那些會改變? 變化為何? (8分)

(二) 在其他變數固定下, 檢定薪資 ( $X_2$ ) 對  $Y$  之效果是否顯著為負? ( $\alpha = 0.05$ ) (8分)

(三)  $H_0: \beta_0 = 0$  vs.  $H_1: \beta_0 \neq 0$  之檢定 ( $\alpha = 0.05$ ), 有何結論? 又, 是否要去除截距項? (敘

【版權所有, 重製必究!】

明理由) (8分)

(四)就檢定  $H_0: \beta_M = \beta_R = 0$  vs.  $H_1: \text{not } H_0$ ，請算出F檢定統計量之值及寫出決策法則。

( $\alpha = 0.05$ ) (12分)

(五)求  $R^2$  (coefficient of determination判定係數)，並解釋其涵意。(8分)

(六)經變數選取後，得估計迴歸式為  $\hat{Y} = 0.01 + 0.025X_1 + 0.045X_3 + 0.342M$ ，請分別解釋  $X_3$  及  $M$  的係數估計值之涵意。(8分)

**答：**

(一)由於統計量不具單位性，故 T 值不變進而 P-Value 也不變。

關於參數估計值與標準誤差，試先算出標準迴歸係數(Beta 係數)

$$\text{Beta}(\hat{\beta}_2) = \frac{S_2}{S_y} \hat{\beta}_2 = \frac{0.1070}{0.7254} \times -0.180 = -0.026651$$

$$\text{Where } S_y = \sqrt{S_y^2} = \sqrt{\frac{\sum(y_i - \bar{y})^2}{n-1}} = \frac{36.304}{70-1} = \sqrt{0.526145} = 0.725358$$

$$\hat{\beta}_2^* = -0.026651 * \frac{0.725358}{1.07} = -0.018$$

$$\text{Where } \hat{\sigma}_2^* = \sqrt{\frac{\sum(10X_2 - 10\bar{X}_2)^2}{70-1}} = \sqrt{100} * S_y = 10 * S_y$$

故單位變換後，參數估計值為原本 0.1 倍，而標準誤差為原本之 10 倍。

總結單位變換後，參數估計值與標準誤差會改變，檢定 T 值與 P-Value 不變。

(二)

$$\text{Step1: } \begin{cases} H_0: \beta_2 = 0 \\ H_1: \beta_2 \neq 0 \end{cases}$$

Step2: p-value = 0.0951 > 0.05 → 在顯著水準  $\alpha = 0.05$  之下，Do not reject  $H_0$

(三)

$$\text{Step1: } \begin{cases} H_0: \beta_0 = 0 \\ H_1: \beta_0 \neq 0 \end{cases}$$

Step2: p-value = 0.2793 > 0.05 → 在  $\alpha = 0.05$  之下，Do not reject  $H_0$

故截距項在檢定水準  $\alpha = 0.05$  不具統計顯著性，建議刪除。

(四) 採偏 F-Test

$$\text{Step: } \begin{cases} H_0: \beta_M = \beta_R = 0 \\ H_1: \beta_M = \beta_R \neq 0 \end{cases}$$

Step2: 決策法則: Reject  $F > F_{\alpha=0.05}(2, 70-6) \leftrightarrow \text{Reject } H_0$

Step3: F 檢定統計量 =

$$\frac{SSE(\text{Reduced}) - SSE(\text{Full})/2}{SSE(\text{Full})/(70-6)} = \frac{(15.574 - 14.005)/2}{14.005/64} = 3.5850$$

(五)

$$R^2 = \frac{SSR}{SSTO} = \frac{\text{可解釋變異}}{\text{總變異}} = \frac{22.299}{36.304} = 0.6142 \text{ 代表此模型可以解釋總變異的 } 61.42\%$$

(六)

$\hat{\beta}_3 = 0.045$  代表在固定其他因素下，所屬小組人數每增加一人，量化工作績效可以增加 0.045 個單位。

$\hat{\beta}_M = 0.342$  代表在固定其他因素下，行銷部門與非行銷部門量化工作績效平均值差異。

【版權所有，重製必究！】