

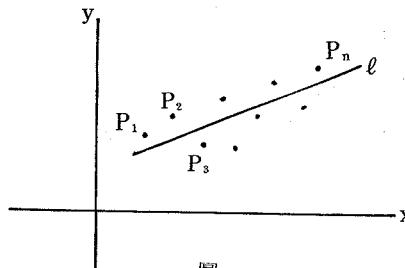
# 迴歸直線與相關係數

葉東進

國立科學園區實驗高中

科學活動的主要內容在通過某些基本模式的建立，以尋找變量之間的關係或規律，但是有些變量之間並不存在必然的關係或是明顯的規律，譬如人的身高與體重、學生的理化成績與數學成績、產品的廣告費與銷售量等；然而它們變化之間卻又隱約有某個趨勢存在。例如兩個變量  $x$  與  $y$  的數據：

X		$x_1, x_2, x_3, \dots, x_n$
y		$y_1, y_2, y_3, \dots, y_n$



圖一

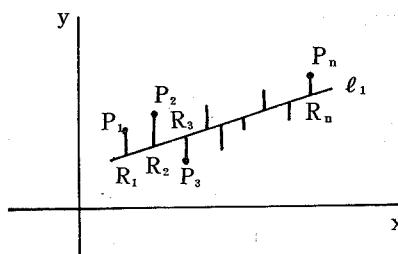
從它們的散佈情形來看（圖一），雖然沒能找到一個基本的模式  $f$  來直接表出  $x$  與  $y$  的關係  $y = f(x)$ ，但是隱約之間，我們看到由這些數據所建構的點  $P_i (x_i, y_i)$  是有向某一直線靠攏的趨勢。因此，退而求其次，我們設法找出這一直線，並對靠攏的趨勢程度加以量化。

## 一、迴歸直線

如（圖一）中的直線  $\ell$ ，其意義可以有許多不同的解釋，只要不離開這些點太離譖

，隨便哪一條直線都可以，問題是：什麼方法才能使找到的直線可以最佳地表達出數據間的變化趨勢？「最小方差法」便是其中的方法之一：

首先，考慮諸點  $P_i(x_i, y_i)$  ( $i = 1, 2, 3, \dots, n$ ) 沿縱軸方向在  $\ell_1$  上的投影  $R_i$  (圖二)，使  $\sum P_i R_i^2$  為最小。



圖二

爲方便處理，取  $\ell_1$  的方程式爲

則有  $\sum \overline{P_i R_i}^2 = \sum (mx_i + k - y_i)^2$  ( 記爲  $f(m, k)$  )

$$\text{由 } \frac{\partial f}{\partial m} = \frac{\partial f}{\partial k} = 0$$

$$\left\{ \begin{array}{l} \sum (mx_i + k - y_i) x_i = 0 \\ \sum (mx_i + k - y_i) = 0 \end{array} \right.$$

$$\Rightarrow \begin{cases} (\sum x_i^2) m + (\sum x_i) k = \sum x_i y_i \\ (\sum x_i) m + n \cdot k = \sum y_i \end{cases}$$

$$解得 \quad m = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2}, \text{ 其中 } \bar{x} = \frac{1}{n} \sum x_i, \bar{y} = \frac{1}{n} \sum y_i$$

$$由 \quad k = \frac{1}{n} \sum y_i - (\frac{1}{n} \sum x_i) m$$

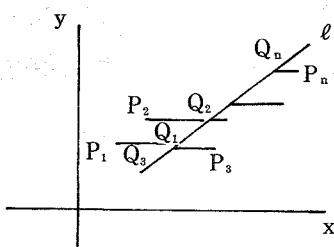
$$\Rightarrow k = \bar{y} - m\bar{x}$$

$$\therefore \bar{y} = m\bar{x} + k$$

此式表示直線  $\ell_1$ :  $y = mx + k$  通過點  $(\bar{x}, \bar{y})$ ，因此  $\ell_1$  的方程式亦可寫成：

我們把直線  $\ell_1$  稱為 y 對 x 的迴歸直線。

其次，我們也必須考慮諸點  $P_i(x_i, y_i)$  ( $i = 1, 2, 3, \dots, n$ ) 沿橫軸方向



圖三

在  $\ell_2$  上的投影  $Q_i$  (圖三), 使  $\sum \overline{P_i Q_i}^2$  為最小。

取  $\ell_2$  的方程式爲

$$x = my + k$$

則有  $\sum \overline{P_i Q_i}^2 = \sum (my_i + k - x_i)^2$  ( 記為  $g(m, k)$  )

$$\text{由 } \frac{\partial g}{\partial m} = \frac{\partial g}{\partial k} = 0$$

得

$$\begin{cases} \sum (m y_i + k - x_i) y_i = 0 \\ \sum (m y_i + k - x_i) = 0 \end{cases} \Rightarrow \begin{cases} (\sum y_i)^2 m + (\sum y_i) k = \sum x_i y_i \\ (\sum y_i) m + n \cdot k = \sum x_i \end{cases}$$

$$解得 \quad m = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum y_i^2 - n \bar{y}^2}$$

$$由 \quad k = \frac{1}{n} \sum x_i - (\frac{1}{n} \sum y_i) m$$

$$\Rightarrow k = \bar{x} - mv$$

$$\therefore \bar{x} = m\bar{y} + k$$

此式表示直線  $\ell_2$  :  $x = my + k$  通過點  $(\bar{x}, \bar{y})$ ，因此  $\ell_2$  的方程式亦可寫成：

$$\ell_2: \quad x - \bar{x} = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum y_i^2 - n \bar{y}^2} (y - \bar{y})$$

我們把直線  $\ell_2$  稱為 x 對 y 的迴歸直線。

顯然， $\ell_1$  與  $\ell_2$  的差異如果愈小，諸點  $P_i$  靠攏  $\ell_1$  與  $\ell_2$  的趨向便愈明顯。

## 二、相關係數

把前述  $y$  對  $x$  的迴歸直線(1)及  $x$  對  $y$  的迴歸直線(2)的斜率分別記為  $m_1$  與  $m_2$ ：

$$m_1 = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2}$$

$$m_2 = \frac{\sum y_i^2 - n \bar{y}^2}{\sum x_i y_i - n \bar{x} \bar{y}}$$

$$\begin{aligned} \text{由 } \sum (x_i - \bar{x})(y_i - \bar{y}) &= \sum x_i y_i - (\sum x_i) \bar{y} - (\sum y_i) \bar{x} + n \bar{x} \bar{y} \\ &= \sum x_i y_i - n \bar{x} \bar{y} - n \bar{y} \bar{x} + n \bar{x} \bar{y} \\ &= \sum x_i y_i - n \bar{x} \bar{y} \end{aligned}$$

$$\begin{aligned} \text{及 } \sum (x_i - \bar{x})^2 &= \sum (x_i^2 - 2x_i \bar{x} + \bar{x}^2) \\ &= \sum x_i^2 - 2(\sum x_i) \bar{x} + n \bar{x}^2 \\ &= \sum x_i^2 - 2n \bar{x} \cdot \bar{x} + n \bar{x}^2 \\ &= \sum x_i^2 - n \bar{x}^2 \end{aligned}$$

$$\text{所以 } m_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

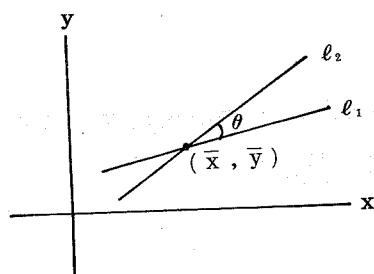
$$m_2 = \frac{\sum (y_i - \bar{y})^2}{\sum (x_i - \bar{x})(y_i - \bar{y})}$$

$$\text{考慮 } \frac{m_1}{m_2} = \frac{[\sum (x_i - \bar{x})(y_i - \bar{y})]^2}{\sum (x_i - \bar{x})^2 \cdot \sum (y_i - \bar{y})^2} (> 0), \text{ 不失一般性, 底下的討論均}$$

假定  $m_1 > 0$ ,  $m_2 > 0$ , 而且  $m_2 \geq m_1$ 。

$$\text{取 } K = \frac{m_1}{m_2}$$

可以看出，在  $m_1$  的值固定的情況下， $m_2$  與  $K$  成反比。就是說  $m_2$  愈小， $K$  值隨之

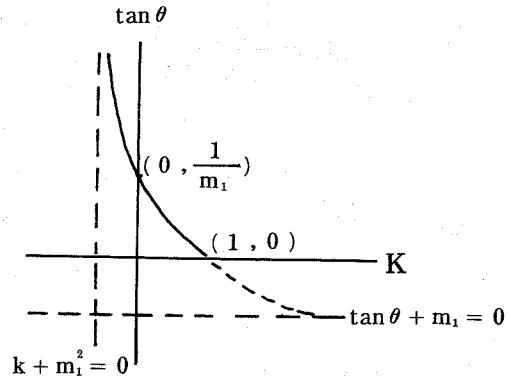


圖四

愈大； $m_2$  愈大， $K$  值隨之愈小。幾何意義看來，便是  $\ell_1$  與  $\ell_2$  的銳夾角  $\theta$  愈小， $K$  值隨之愈大； $\theta$  愈大， $K$  值隨之愈小（圖四）。因此， $K$  值的大小能夠適度地反映出  $\ell_2$  與  $\ell_1$  之間的差異。

事實上，由

$$\tan \theta = \frac{m_2 - m_1}{1 + m_2 m_1} = \frac{1 - \frac{m_1}{m_2}}{\frac{1}{m_2} + m_1} = \frac{1 - K}{\frac{K}{m_1} + m_1}$$



圖五

因為  $\theta$  為銳角， $\tan \theta \geq 0$ ，即  $0 < K \leq 1$ 。

另外，從  $(\tan \theta + m_1)(K + m_1^2) = m_1 + m_1^3$  可以繪出  $\tan \theta$  與  $K$  之間的變化關係圖形（圖五），從圖形中我們看出  $\tan \theta$  是隨  $K$  值的增大而遞減；隨  $K$  值的減小而遞增。也可以說  $\theta$  與  $K$  值成遞減關係。

既然  $K$  值的大小可以反映出  $\ell_2$  與  $\ell_1$  間的差異，而  $\ell_2$  與  $\ell_1$  間的差異又反映出諸點  $P_i$  靠攏  $\ell_1$  與  $\ell_2$  的趨向，因此  $K$  值的大小便反映出諸點  $P_i$  靠攏  $\ell_1$  與  $\ell_2$  的趨向程度。

由  $K = \frac{m_1}{m_2} = \frac{[\sum (x_i - \bar{x})(y_i - \bar{y})]^2}{\sum (x_i - \bar{x})^2 \cdot \sum (y_i - \bar{y})^2} > 0$

取  $\sigma = \sqrt{K}$

我們把  $\sigma$  稱作是兩個變量  $x$  與  $y$  之間的相關係數。

因此，兩個變量  $x$  與  $y$  之間的相關係數便是反應諸點  $P_i(x_i, y_i)$  靠攏迴歸直線的程度的一個數量。

## 參考資料

1. 高中基礎數學第四冊，師大科教中心，國立編譯館。
2. 普通數學教程，楊維哲、蔡聰明，文仁。
3. 統計學初階，賴建業譯，中央。