

รูปที่ 5.5 แสดง sample regression line และเส้นแสดงความสัมพันธ์ในรูปเส้นโค้ง
ซึ่งสูexe้ากับจุดทั้ง 11 จุดมากกว่าเส้นตรง

หมายเหตุ

ในการทำ Regression Analysis เราควร plot scatter diagram เป็นอย่างแรกเพื่อเป็นแนวทางในการตั้งค่าแบบที่จะใช้แสดงความสัมพันธ์ระหว่าง X และ Y แต่ตัวอย่างนี้ต้องการใช้เป็นตัวอย่างในการทดสอบ lack of fit จึง plot scatter diagram ที่หลังเพื่อตรวจสอบกัน

ตัวอย่างที่ 5.9 ปฏิกริยาทางเคมีที่วัดได้ ณ อุณหภูมิต่าง ๆ กันมีดังนี้

x(°C)	y(%)	x(°C)	y(%)
150	77.4	250	88.9
150	76.7	250	89.2
150	78.2	250	89.7
200	84.1	300	94.8
200	84.5	300	94.7
200	83.7	300	95.9

$T_1 = 232.3$

$T_2 = 252.3$

$T_3 = 267.8$

$T_4 = 285.4$

羌ໍາ sample regression line ແລະ ຖດສອນ lack of fit

ຈາກຂໍ້ມູນຂ້າງຕົ້ນ $n = 12$

$$\sum x_i = 2700, \sum x_i^2 = 645000, \bar{x} = 225$$

$$\sum y_i = 1037.8, \sum y_i^2 = 90265.52, \bar{y} = 86.4833$$

$$\sum x_i y_i = 237875$$

$$S_{xx} = 645000 - (2700)^2/12 = 37500$$

$$S_{yy} = 90265.52 - (1037.8)^2/12 = 513.1167$$

$$S_{xy} = 237875 - (2700)(1037.8)/12 = 4370$$

$$b_1 = \hat{\beta}_1 = \frac{4370}{37500} = .1165$$

$$b_0 = \hat{\beta}_0 = 86.4833 - (.1165)(225) = 60.2708$$

ດັ່ງນີ້ sample regression line ຕີ່ມີ $\hat{y} = 60.2708 + .1165x$

ຈະຖດສອນ 1) $H_0: \beta_1 = 0$ ແລະ

2) $H_0:$ ໃນມື້ lack of fit

ຄໍານວດ $SST = S_{yy} = 513.1167$

$$SSR = b_1 S_{xy} = (.1165)(4370) = 509.1050$$

$$SSE = SST - SSR = 4.0117$$

ຄໍານວດ SS(p.e.) ໂດຍຫາ $T_1 = 232.3, T_2 = 252.3, T_3 = 267.8$ ແລະ $T_4 = 285.4$

$$SS(p.e.) = \sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij}^2 - \frac{\left(\sum_{j=1}^{n_i} y_{ij} \right)^2}{n_i}$$

$$= \sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij}^2 - \frac{\sum_{i=1}^k T_i^2}{n_i}$$

$$= 90265.52 - \frac{(232.3)^2 + (252.3)^2 + (267.8)^2 + (285.4)^2}{3}$$

$$= 2.66$$

$$\begin{aligned} \text{SS(los)} &= \text{SSE} - \text{SS(p.e.)} \\ &= 4.0117 - 2.66 = 1.3517 \end{aligned}$$

ANOVA

S.V.	d.f.	SS	MS	f_c
Regression	1	509.1050	509.1050	$f_1 = 1531.60^{**}$
Error	10	4.0117	.40117	
Lack of fit	2	1.3517	.6758	$f_2 = 2.03$ (n.s.)
Pure error	8	2.6600	.3324	
Total	11	513.1167		

- ทดสอบ 1) $H_0 : \beta_1 = 0$
 2) $H_1 : \beta_1 \neq 0$
 3) $\alpha = .05$
 4) CR : $F > f_{(1,10), .05} = 4.9646$
 5) $f_c = f_1 = 1531.60$
 6) $f_c > 4.9646$ เรากล่าวว่า H_0 ที่ $\alpha = .05$ (เรากล่าวว่า H_0 ที่ $\alpha = .01$ ด้วย)

นั่นคือ มีความสัมพันธ์ระหว่าง X และ Y

- ทดสอบ 1) H_0 : ไม่มี lack of fit
 2) H_1 : มี lack of fit
 3) $\alpha = .05$
 4) CR : $F > f_{(2,8), .05} = 4.459$
 5) $f_c = f_2 = 2.03$
 6) $f_c < 4.459$ เราไม่สามารถปฏิเสธ H_0 ที่ $\alpha = .05$ ได้ นั่นคือไม่มี lack of fit

$$\text{ถ้าค่า R^2} = \frac{\text{SSR}}{\text{SST}} = .9922$$

$100r^2\% = 99.22\%$ ของความแปรปรวนของ Y เนื่องมาจากความสัมพันธ์กับ X ในรูปเชิงเส้นตรง

สรุปได้ว่า ตัวแบบเชิงเส้นตรงที่ใช้อยู่เหมาะสมดีแล้ว

5.7 สมสัมพันธ์ (Correlation) : เครื่องมือวัดความสัมพันธ์เชิงเส้นตรง

เมื่อเรามีตัวแปรเชิงสุ่ม Y ซึ่งขึ้นอยู่กับ X ที่ถูกควบคุมโดยผู้ทำการทดลอง และทำการวิเคราะห์เพื่อคุณิตพารามิเตอร์ของ X หรือประโยชน์ของมันในการทำนายค่า Y เราจะใช้วิธีการของ regression แต่ผู้ทำการทดลองอาจต้องชุดมุ่งหมายที่จะทราบความสัมพันธ์ของตัวแปรเชิงสุ่ม 2 ตัว โดยที่เราไม่สนใจว่าตัวใดเป็นเหตุและตัวใดเป็นผล เราจะใช้วิธีการของ correlation

โครงสร้างของข้อมูลประกอบด้วย ตัวอย่างสุ่มจาก Bivariate Normal distribution ขนาด n : $(X_1, Y_1), \dots, (X_n, Y_n)$ และต่างๆกันจะเป็นอิสระต่อกัน

5.7.1 Bivariate Normal distribution

$$f(x, y) = \frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)}Q}, -\infty < x < \infty \text{ และ } -\infty < y < \infty$$

$$\text{โดยที่ } Q = \left(\frac{x - \mu_x}{\sigma_x}\right)^2 - 2\rho\left(\frac{x - \mu_x}{\sigma_x}\right)\left(\frac{y - \mu_y}{\sigma_y}\right) + \left(\frac{y - \mu_y}{\sigma_y}\right)^2$$

ρ คือสัมประสิทธิ์สมสัมพันธ์ของประชากร (Population correlation coefficient)

$$\begin{aligned} \rho &= \text{corr}(X, Y) \\ &= \frac{\text{cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}} \end{aligned}$$

5.7.2 พิสูจน์ว่า $-1 \leq \rho \leq 1$

ให้ $\rho = \text{corr}(X, Y)$

: ถ้า standardized random variable 2 ฟังก์ชัน

$$\left| \begin{array}{l} Z_1 = \frac{X - \mu_x}{\sigma_x} \\ = \frac{X}{\sigma_x} - \frac{\mu_x}{\sigma_x} \end{array} \right| \quad \left| \begin{array}{l} Z_2 = \frac{Y - \mu_y}{\sigma_y} \\ = \frac{Y}{\sigma_y} - \frac{\mu_y}{\sigma_y} \end{array} \right|$$

$\text{corr}(Z_1, Z_2) = \rho$ (เพราะว่าสัมประสิทธิ์ของ X และ Y ใน Z_1 และ Z_2 มีเครื่องหมายเหมือนกัน)

$\therefore \text{Variance ของตัวแปรเชิงสูงๆ } \geq 0$

$$\text{Var}(Z_1 + Z_2) = \text{Var}(Z_1) + \text{Var}(Z_2) - 2 \text{Cov}(Z_1, Z_2) \geq 0$$

$$\text{Var}(Z_1) + \text{Var}(Z_2) \geq 2 \text{Cov}(Z_1, Z_2)$$

$$1 + 1 \geq \frac{2 \text{Cov}(Z_1, Z_2)}{\sigma_{z_1} \sigma_{z_2}}$$

$$1 \geq \text{Corr}(Z_1, Z_2) = \rho \quad \dots\dots\dots(1)$$

$$\text{Var}(Z_1 - Z_2) = \text{Var}(Z_1) + \text{Var}(Z_2) + 2 \text{Cov}(Z_1, Z_2) \geq 0$$

$$\text{Var}(Z_1) + \text{Var}(Z_2) \geq -2 \text{Cov}(Z_1, Z_2)$$

$$1 + 1 \geq -2\rho$$

$$-1 \leq \rho \quad \dots\dots\dots(2)$$

จาก (1) และ (2) สรุปได้ว่า $-1 \leq \rho \leq 1$

5.7.3 สัมประสิทธิ์สหสัมพันธ์ของตัวอย่าง (Sample correlation coefficient หรือ Pearson's product moment)

$$r = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}} , \text{ โดยที่ } S_{xx}, S_{yy} \text{ และ } S_{xy} \text{ ได้กล่าวถึงในหัวข้อ 5.2}$$

R เป็น Unbiased estimator ของ ρ ซึ่งมี mean $E(R) = \rho$ และ variance

$$\text{Var}(R) = \left(\frac{1 - \rho^2}{n - 2} \right)$$

คุณสมบัติของ r (ค่าของตัวแปรเชิงสูง R)

$$1) -1 \leq r \leq 1$$

$r = 1$ เกิดขึ้นเมื่อจุดทั้ง n จุดอยู่บนเส้นตรงที่มี slope เป็น +

$r = -1$ เกิดขึ้นเมื่อจุดทั้ง n จุดอยู่บนเส้นตรงที่มี slope เป็น

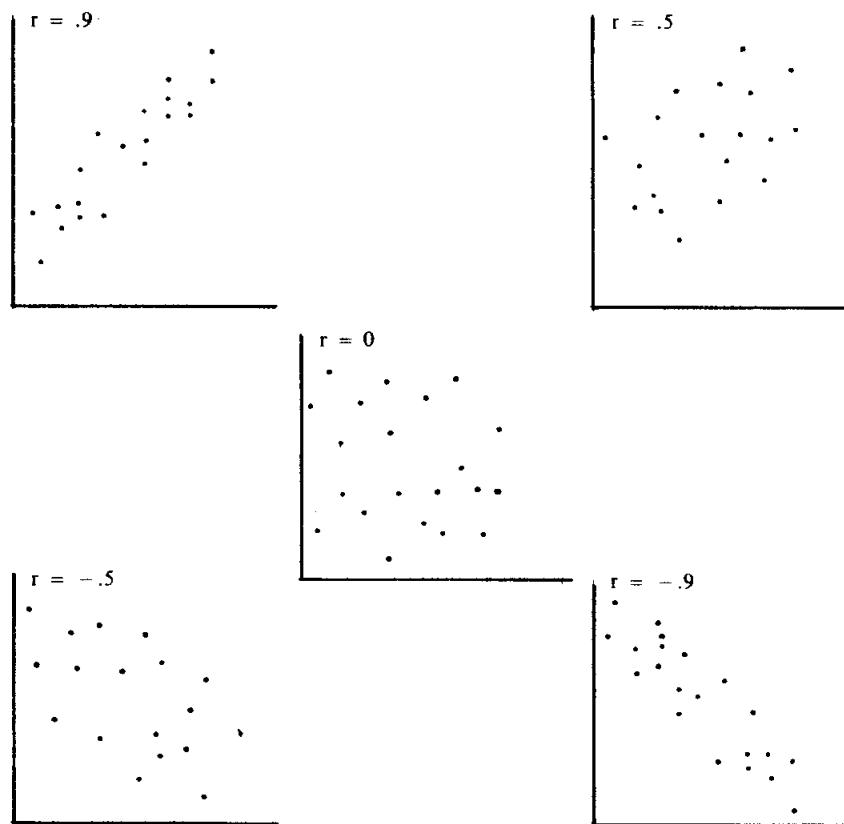
$r = \pm 1$ แสดงว่า X และ Y มีความสัมพันธ์เชิงเส้นตรงอย่างสมบูรณ์

$r = 0$ หมายความว่า X และ Y ไม่มีความสัมพันธ์กันในเชิงเส้นตรง แต่อาจมีความสัมพันธ์กันในรูปอื่น หรือไม่มีความสัมพันธ์กันเลย

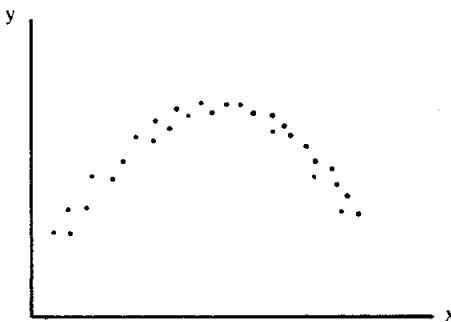
2) ค่าของ r ใช้วัดอัตราความมากน้อยของความสัมพันธ์เชิงเส้นตรง เครื่องหมายของ r บอกทิศทางของความสัมพันธ์ นั่นคือ ถ้า r เป็น + ค่าของตัวหนึ่งเพิ่ม ค่าของอีกตัวหนึ่งจะเพิ่มขึ้นด้วย และถ้า r เป็น - ค่าของทั้ง 2 ตัว จะมีทิศทางตรงกันข้าม

3) r^2 เป็นอัตราส่วนของความแปรปรวนของ Y ที่อธิบายได้ด้วย regression line ที่คำนวนได้ด้วยวิธีการของ least square method หรืออาจกล่าวว่าอัตราส่วนของความแปรปรวนของ Y ที่เนื่องมาจากการอิทธิพลของความสัมพันธ์เชิงเส้นตรงกับ X

4) $\text{Corr}(aX + b, cY + d) = \text{Corr}(X, Y)$ ถ้า a และ c มีเครื่องหมายเดียวกัน



รูปที่ 5.6 แสดงค่าของ r และความมากน้อยของการกระจายของจุด



รูปที่ 5.7 แสดงความสัมพันธ์ของ X และ Y ในรูปเส้นโค้ง ถ้าคำนวณค่า r , ค่าเกือบเป็นศูนย์

5.8 ความสัมพันธ์ระหว่าง Sample regression coefficient (b_1) และ Sample correlation coefficient (r)

$$b_1 = \frac{S_{xy}}{S_{xx}}$$

$r = \frac{b_1 \sqrt{S_{xx}}}{\sqrt{S_{yy}}}$ นั้นคือ r และ b_1 จะมีเครื่องหมายเหมือนกัน

$$r^2 = \frac{b_1^2 S_{xx}}{S_{yy}}$$

$$\therefore SSE = SST - SSR$$

$$= S_{yy} - b_1 S_{xy}$$

$$= S_{yy} - b_1^2 S_{xx}$$

$$\frac{SSE}{S_{yy}} = 1 - \frac{b_1^2 S_{xx}}{S_{yy}}$$

$$b_1^2 \frac{S_{xx}}{S_{yy}} = 1 - \frac{SSE}{S_{yy}}$$

$$r^2 = 1 - \frac{SSE}{S_{yy}}$$

$$= \frac{S_{yy} - SSE}{S_{yy}}$$

$$= \frac{SSR}{SST}$$

ถ้า $SSE = 0, r^2 = 1$ นั่นคือ $r = -1$ หรือ $r = +1$ ซึ่งกรณีนี้เกิดเมื่อจุดทุกจุดอยู่บนเส้นตรง

นั่นคือ ถ้า $SSE \rightarrow 0, r^2 \rightarrow 1$ นั่นคือ ถ้าค่า r^2 ยิ่งสูง ค่า SSE จะยิ่งน้อย เราจึงใช้ r^2 เป็นมาตรการช่วยชี้ว่าความแปรปรวนของ γ ที่เนื่องมาจากการอิทธิพลของความสัมพันธ์เชิงเส้น ตรงกับ X มีมากน้อยเท่าไร

5.9 การทดสอบสมมติฐานและการหาช่วงความเชื่อมั่นเกี่ยวกับ ρ

5.9.1 ทดสอบ $H_0 : \rho = 0, H_1 : \rho \neq 0$

$$\text{ตัวสถิติที่ใช้ในการทดสอบ คือ } T = \frac{R - \rho}{\sqrt{\frac{1 - R^2}{n - 2}}} \sim t_{n-2}$$

ถ้า H_0 จริง $t_c = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}$ เป็นค่าของตัวสถิติ T ซึ่งมีการกระจายเป็น t -distribution

ที่มี d.f. = $n - 2$

5.9.2 ทดสอบ $H_0 : \rho = \rho_0 (\rho_0 \neq 0)$

$\frac{1}{2} \ln \left(\frac{1+r}{1-r} \right)$ เป็นค่าของตัวแปรเชิงสุ่ม ซึ่งมีการกระจายเป็น Normal distribution

ที่มี mean $\frac{1}{2} \ln \left(\frac{1+\rho}{1-\rho} \right)$ และ variance $\frac{1}{n-3}$

($\ln = \log$ ฐาน $e = 2.3026 \log_{10}$)

$$\begin{aligned} \text{ถ้า } H_0 \text{ จริง } z_c &= \sqrt{\frac{n-3}{2}} \left[\ln \left(\frac{1+r}{1-r} \right) - \ln \left(\frac{1+\rho_0}{1-\rho_0} \right) \right] \\ &= \sqrt{\frac{n-3}{2}} \ln \left[\frac{(1+r)(1-\rho_0)}{(1-r)(1+\rho_0)} \right] \text{ จะเป็นค่าของ } Z \sim N(0, 1) \end{aligned}$$

ตาราง A10 เราสามารถอ่านค่าของ $z = \frac{1}{2} \ln \left(\frac{1+r}{1-r} \right)$ เราจึงควรใช้สูตร

$$z_c = \frac{\frac{1}{2} \ln \left(\frac{1+r}{1-r} \right) - \frac{1}{2} \ln \left(\frac{1+\rho_0}{1-\rho_0} \right)}{\sqrt{\frac{1}{n-3}}}$$

$$= \frac{z_r - z_{\rho}}{\sqrt{\frac{1}{n-3}}}$$

$$\text{โดยที่ } z_r = \frac{1}{2} \ln \left(\frac{1+r}{1-r} \right) \text{ และ } z_{\rho} = \frac{1}{2} \ln \left(\frac{1+\rho}{1-\rho} \right)$$

5.9.3 การหา $100(1 - \alpha)\%$ C.I. ของ ρ

$$\text{จาก } Z = \frac{z_r - z_{\rho}}{\sqrt{\frac{1}{n-3}}} \sim N(0, 1)$$

$$\Pr \left[-z_{\frac{\alpha}{2}} < Z < z_{\frac{\alpha}{2}} \right] = 1 - \alpha$$

$$\Pr \left[-z_{\frac{\alpha}{2}} < \frac{z_r - z_{\rho}}{\sqrt{\frac{1}{n-3}}} < z_{\frac{\alpha}{2}} \right] = 1 - \alpha$$

$$\Pr \left[z_r - z_{\rho} \frac{1}{\sqrt{n-3}} < z_{\rho} \leq z_r + z_{\rho} \frac{1}{\sqrt{n-3}} \right] = 1 - \alpha$$

$$100(1 - \alpha)\% \text{ C.I. ของ } z_{\rho} \text{ คือ } z_{\rho} \pm z_{\frac{\alpha}{2}} \frac{1}{\sqrt{n-3}}$$

ในการหา C.I. ของ ρ อาจทำได้โดยใช้ตาราง A10 เพื่อเปลี่ยนค่า $z_r - z_{\rho} \frac{1}{\sqrt{n-3}}$

และ $z_r + z_{\rho} \frac{1}{\sqrt{n-3}}$ ไปเป็น correlation coefficients หรืออาจใช้ Hyperbolic tangent function

จะได้ $100(1 - \alpha)\% \text{ C.I. ของ } \rho \text{ คือ}$

$$\left[\tanh \left(z_r - z_{\rho} \frac{1}{\sqrt{n-3}} \right), \tanh \left(z_r + z_{\rho} \frac{1}{\sqrt{n-3}} \right) \right]$$

ตัวอย่างที่ 5.10 [การศึกษาถึงความสัมพันธ์ของปริมาณน้ำฝนที่ตกลงมากับปริมาณอากาศเสีย (air pollution) ที่จางหายไปหลังจากฝนตก ได้ผลดังนี้]

x : ปริมาณฝนตกต่อวัน (.01 นิว)	4.3 4.5 5.9 5.6 6.1 5.2 3.8 2.1 7.5
y : ปริมาณของอากาศสกปรกที่หายไป (ในไมครอรัม/ตารางเมตร)	126 121 116 118 114 118 132 141 108

$$\begin{aligned}
 n &= 9, \quad \Sigma x_i = 45, \quad \Sigma x_i^2 = 244.26, \quad \bar{x} = 5 \\
 \Sigma y_i &= 1094, \quad \Sigma y_i^2 = 133786, \quad \bar{y} = 121.56 \\
 \Sigma x_i y_i &= 5348.2 \\
 S_{xx} &= 19.26 \\
 S_{yy} &= 804.22223 \\
 S_{xy} &= -121.8 \\
 r &= \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}} = \frac{-121.8}{\sqrt{(19.26)(804.22223)}} \\
 &= -.9786
 \end{aligned}$$

(ถ้าหาค่า $b_1, b_0 = -6.3240$)

- 1) $H_0: \rho = 0$
- 2) $H_1: \rho \neq 0$
- 3) $\alpha = .05$
- 4) CR : $T < -t_{1-\alpha/2} = -2.365$ และ $T > 2.365$

$$5) t_c = \frac{-0.9786\sqrt{7}}{\sqrt{1 - (-0.9786)^2}} = -12.9032$$

6) $t_c < -2.365$ (ตกใน CR) เราปฏิเสธ H_0 ที่ $\alpha = .05$ นั้นคือ X และ Y มีความสัมพันธ์กันเชิงเส้นตรง

ตัวอย่างที่ 5.11 จากการศึกษาถึงความสัมพันธ์ของคะแนน Midterm กับคะแนนสอบໄลในวิชาสถิติเบื้องต้นว่ามีความสัมพันธ์กันหรือไม่ ได้สุ่มนักศึกษามา 20 คนได้คะแนนดังนี้

x : คะแนน Midterm	35	60	55	35	35	50	30	60	50	20
y : คะแนนสอบໄล	60	80	60	80	75	90	60	105	60	30

x : คะแนน Midterm	55	45	40	60	40	60	50	55	50	35
y : คะแนนสอบໄล	90	75	80	80	45	80	80	95	100	75

$$\begin{aligned}
 n &= 20, \Sigma x_i = 920 & \Sigma x_i^2 &= 44900 \\
 \Sigma y_i &= 1500 & \Sigma y_i^2 &= 118850 \\
 \Sigma x_i y_i &= 71600 \\
 S_{xx} &= 2580 \\
 S_{yy} &= 6350 \\
 S_{xy} &= 2600
 \end{aligned}$$

- ก. 1) $H_0 : \rho = .9$
 2) $H_1 : \rho \neq .9$
 3) $\alpha = .05$
 4) CR : $Z < -1.96$ หรือ $Z > 1.96$

$$5) r = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}} = \frac{2600}{\sqrt{(2580)(6350)}} = .642$$

จากตาราง A10 : $z_r = z_{.642} = .7616$

$$z_{\rho_0} = z_{.9} = 1.4722$$

$$z_c = \frac{.7616 - 1.4722}{\sqrt{1/17}} = -2.93$$

$$\begin{aligned}
 \text{หรือ } z_c &= \frac{\sqrt{n-3}}{2} \ln \left[\frac{(1+r)(1-\rho_0)}{(1-r)(1+\rho_0)} \right] \\
 &= \frac{\sqrt{17}}{2} \ln \left[\frac{(1.642)(.1)}{(.358)(1.9)} \right] \\
 &= \frac{\sqrt{17}}{2} \ln (.2413996) \\
 &= (2.0615528)(-1.4213017) \\
 &= -2.930085 \approx -2.93
 \end{aligned}$$

$$6) z_c = -2.93 < -1.96 \text{ เราปฏิเสธ } H_0 \text{ ที่ } \alpha = .05$$

ข. หา 90% C.I. ของ ρ

$$\text{ก่อนอื่นหา 90% C.I. ของ } z_{.9} = z_r \pm z_{\frac{\alpha}{2}} \sqrt{\frac{1}{n-3}}$$

$$\begin{aligned}
 &= .7616 \pm \frac{1.645}{\sqrt{17}} \\
 &= .7616 \pm (1.645) (.2425) \\
 &= (.348, .821)
 \end{aligned}$$

1) ใช้ตาราง A10 เปลี่ยนค่า .3627 และ 1.1605 เป็น correlation coefficients
 $\therefore 90\% \text{ C.I. ของ } \rho \text{ คือ } (.348, .821)$

2) ใช้ Hyperbolic tangent function

90% C.I. ของ ρ คือ ($\tanh (.3627)$, $\tanh (1.1605)$)

$$= (.34759, .8212028) = (.348, .821)$$

5.10 การทดสอบการเท่ากันของ correlation coefficient 2 ตัว ($H_0 : \rho_1 = \rho_2$)

โครงสร้างของข้อมูล

x_{1i}	x_{11}	x_{12}	...	x_{1n_1}
y_{1i}	y_{11}	y_{12}	...	y_{1n_1}

คำนวณ $r_1 = \frac{S_{xy}^{(1)}}{\sqrt{S_{xx}^{(1)} S_{yy}^{(1)}}}$

x_{2i}	x_{21}	x_{22}	...	x_{2n_2}
y_{2i}	y_{21}	y_{22}	...	y_{2n_2}

คำนวณ $r_2 = \frac{S_{xy}^{(2)}}{\sqrt{S_{xx}^{(2)} S_{yy}^{(2)}}}$

1) $H_0 : \rho_1 = \rho_2$

2) $H_1 : \rho_1 \neq \rho_2$

3) กำหนด α

4) CR : $Z > z_{\alpha/2}$ และ $Z < -z_{\alpha/2}$

5) $z_c = \frac{z_1 - z_2}{\sqrt{\sigma_1^2 + \sigma_2^2}}$

โดยที่ $z_1 = \frac{1}{2} \ln \left(\frac{1+r_1}{1-r_1} \right)$,

$z_2 = \frac{1}{2} \ln \left(\frac{1+r_2}{1-r_2} \right)$,

$\sigma_1^2 = \frac{1}{n_1 - 3}$ และ $\sigma_2^2 = \frac{1}{n_2 - 3}$

6) สรุปผล

ตัวอย่างที่ 5.12 ในการศึกษาถึงความแตกต่างของสัมประสิทธิ์สหสัมพันธ์ของคะแนนวิชาค่านวน (X) และคะแนนวิชาภาษาอังกฤษ (Y) ของนักเรียนจาก 2 โรงเรียน ได้ทำการสุ่มนักเรียนชั้น ม.1 12 คนจากโรงเรียนที่ 1 และ 10 คนจากโรงเรียนที่ 2 ผลปรากฏดังนี้

โรงเรียนที่ 1	x_{1i}	41 39 53 67 61 67 46 50 55 72 63 59
	y_{1i}	29 19 30 27 28 27 22 29 24 33 25 20

โรงเรียนที่ 2	x_{2i}	56 38 52 40 65 61 64 64 53 51
	y_{2i}	34 21 25 24 32 29 27 26 24 25

จงทดสอบที่ $\alpha = .05$ ว่าสัมประสิทธิ์สหสัมพันธ์ (ของคะแนนทั้ง 2 วิชา) ของ 2 โรงเรียนนี้ไม่แตกต่างกัน

$$n_1 = 12, \Sigma x_{1i} = 673, \Sigma x_{1i}^2 = 38985$$

$$\Sigma y_{1i} = 313, \Sigma y_{1i}^2 = 8359$$

$$\Sigma x_{1i} y_{1i} = 17759$$

$$S_{xx}^{(1)} = 1240.9167$$

$$S_{yy}^{(1)} = 194.9167$$

$$S_{xy}^{(1)} = 204.9167$$

$$r_1 = \frac{S_{xy}^{(1)}}{\sqrt{S_{xx}^{(1)} S_{yy}^{(1)}}} = .417$$

จากตาราง A10 : $z_1 = z_{.417} = .4441$

$$n_2 = 10, \Sigma x_{2i} = 544, \Sigma x_{2i}^2 = 30432$$

$$\Sigma y_{2i} = 267, \Sigma y_{2i}^2 = 7269$$

$$\Sigma x_{2i} y_{2i} = 14750$$

$$S_{xx}^{(2)} = 838.4$$

$$S_{yy}^{(2)} = 140.1$$

$$S_{xy}^{(2)} = 225.2$$

$$r_2 = \frac{S_{xy}^{(2)}}{\sqrt{S_{xx}^{(2)} S_{yy}^{(2)}}} = .657$$

จากตาราง A10 : $z_1 = z_{.68} = .7875$

1) $H_0 : \rho_1 = \rho_2$

2) $H_1 : \rho_1 \neq \rho_2$

3) $\alpha = .05$

4) CR : $Z < -z_{.025} = -1.96$ และ $Z > 1.96$

$$5) z_c = \frac{z_1 - z_2}{\sqrt{\frac{1}{n_1-3} + \frac{1}{n_2-3}}}$$

$$= \frac{.4441 - .7875}{\sqrt{\frac{1}{9} + \frac{1}{7}}}$$

$$= \frac{- .3434}{.5039526} = - .6814$$

6) $\because z_c = - .6814$ ไม่ตกอยู่ใน CR เราไม่สามารถปฏิเสธ H_0 ที่ $\alpha = .05$ ได้นั้นคือ สัมประสิทธิ์สหสัมพันธ์ (ของคะแนนทั้ง 2 วิชา) ของ 2 โรงเรียนนี้ไม่แตกต่างกันอย่างมีนัยสำคัญ

5.11 Serial correlation

เมื่อเราสังเกตค่าของตัวแปรตัวหนึ่ง ณ เวลาที่ต่าง ๆ กัน สิ่งสำคัญที่เราควรพิจารณา คือโอกาสที่จะเกิดความสัมพันธ์ระหว่างค่าสังเกตที่ติดตามกันมา หรืออาจกล่าวได้ว่าหากคุณ-สมบัติของการสุ่ม (randomness) คำว่า “serial correlation” นี้หมายถึงความสัมพันธ์ระหว่าง ค่าสังเกตที่ติดต่อกันตามลำดับเวลา หรือตามลำดับอื่น ๆ ของการเก็บรวบรวมข้อมูล

กำหนดให้ x_1, \dots, x_n แทนค่าสังเกตของตัวแปร X ณ เวลา t จุดที่ต่อเนื่องกัน เพื่อ ที่จะตรวจสอบว่าสมาชิกที่อยู่ต่อเนื่องกันนั้นมีความสัมพันธ์กันหรือไม่ เราจับคู่ค่าสังเกตที่อยู่ติดกันคือ (x_{i-1}, x_i) โดยที่ $i = 2, \dots, n$ ดังนั้นเราจะได้ $(n-1)$ คู่ ดังนี้คือ

$$(x_1, x_2), (x_2, x_3), \dots, (x_{i-1}, x_i), \dots, (x_{n-1}, x_n)$$

เพื่อช่วยการสำรวจในขั้นต้น เราอาจสร้าง Scatter diagram ซึ่งประกอบด้วย $(n-1)$ จุดข้างต้น โดยให้สมาชิกตัวแรกของ (x_{i-1}, x_i) เป็นค่านแกน x และสมาชิกตัวที่สองเป็นค่า บันแกน y

สัมประสิทธิ์สหสัมพันธ์ซึ่งคำนวณจากค่าสังเกต $(n-1)$ คู่ข้างบนเรารอเรียกว่า first serial correlation หรือ first-order autocorrelation หรือ lag 1 correlation

Lag 1 correlation หรือ first-order autocorrelation

$$r_1 = \frac{\sum_{i=2}^n (x_{i-1} - \bar{x})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\begin{aligned} \text{พิจารณา } \sum_{i=2}^n (x_{i-1} - \bar{x})(x_i - \bar{x}) &= \sum_{i=2}^n (x_{i-1}x_i - x_{i-1}\bar{x} - x_i\bar{x} + \bar{x}^2) \\ &= \sum_{i=2}^n x_{i-1}x_i - \bar{x} \left[\sum_{i=2}^n x_{i-1} + \sum_{i=2}^n x_i \right] + (n-1)\bar{x}^2 \\ &= \sum_{i=2}^n x_{i-1}x_i - \bar{x} \left[\sum_{i=1}^{n-1} x_i + \sum_{i=2}^n x_i \right] + (n-1)\bar{x}^2 \\ \sum_{i=1}^n (x_i - \bar{x})^2 &= \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2/n \end{aligned}$$

ดังนั้นในการคำนวณโดยใช้เครื่องคิดเลขแนะนำให้ใช้สูตร

$$r_1 = A/B$$

$$\text{โดยที่ } A = \sum_{i=2}^n x_{i-1}x_i - \bar{x} \left[\sum_{i=1}^{n-1} x_i + \sum_{i=2}^n x_i \right] + (n-1)\bar{x}^2$$

$$\text{และ } B = \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2/n$$

r_1 เป็นค่าค่าหนึ่งของตัวแปรเชิงสุ่ม R_1 ซึ่งเป็นตัวประมาณค่า (Estimator) ของ ρ_1 (population first-order autocorrelation) สำหรับอย่างมีขันดั่งใหญ่ การทดสอบสมมติฐาน H_0 : $\rho_1 = 0$ (ข้อมูลไม่มีความสัมพันธ์กันในตัวข้อมูลเอง) เราจะใช้ความจริงที่ว่าเราอาจประมาณการแยกของ $\sqrt{n}R_1$ เป็น $N(0, 1)$

สรุปขั้นในการทดสอบสมมติฐานเกี่ยวกับ ρ_1 (เมื่อตัวอย่างมีขนาดใหญ่)

- 1) $H_0 : \rho_1 = 0$
- 2) $H_1 : 1) \rho_1 \neq 0$
2) $\rho_1 > 0$
3) $\rho_1 < 0$
- 3) กำหนด α
- 4) CR : 1) $Z > z_{\alpha/2}$ และ $Z < -z_{\alpha/2}$
2) $Z > z_\alpha$
3) $Z < -z_\alpha$
- 5) $z_c = \sqrt{n} r_1$

6) สรุป

ก) ถ้า z_c ตกใน CR. เราปฏิเสธ H_0 ที่ระดับนัยสำคัญ $= \alpha$ นั้นคือข้อมูลชุดนี้มีความสัมพันธ์กันในตัวข้อมูลเอง

ข) ถ้า z_c ไม่ตกอยู่ใน CR. เราไม่ปฏิเสธ H_0 ที่ระดับนัยสำคัญ $= \alpha$ นั้นคือข้อมูลชุดนี้ไม่มีความสัมพันธ์กันในตัวข้อมูลเอง หรือข้อมูลมีคุณสมบัติของการสุ่ม (randomness) ตัวอย่างที่ 5.13 โรงงานผลิตแผ่นพลาสติกพยาຍາมที่จะควบคุมน้ำหนักของแผ่นพลาสติกให้ใกล้เคียงกับ 125 ปอนด์ต่อหนึ่งหน่วยพื้นที่ ข้อมูลทั้ง 28 ค่าต่อไปนี้ซึ่งเป็นค่าที่ต่อเนื่องกันนั้นได้มาจากการหาค่าน้ำหนักที่ต่างไปจาก 125

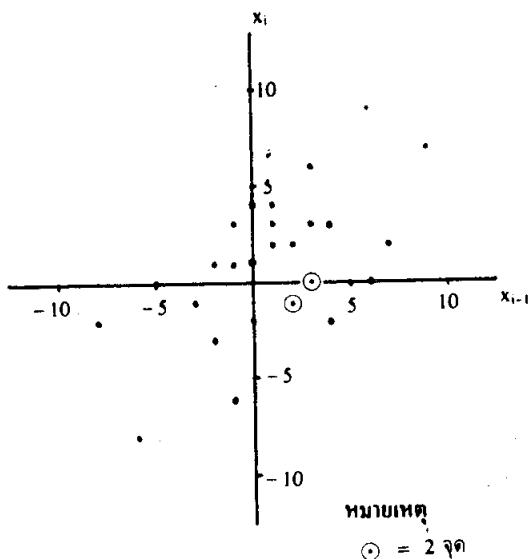
ค่าสังเกตตัวที่	ค่าสังเกต									
1 – 10	6	9	7	2	-1	-6	-8	-2	1	3
11 – 20	0	1	4	3	3	0	4	-2	-3	-1
21 – 28	1	2	2	-1	3	6	0	-2		

ก) จงสร้าง Scatter diagram

ข) จงทดสอบ $H_0 : \rho_1 = 0$, $H_1 : \rho_1 \neq 0$ ที่ $\alpha = .05$

ก) 27 จุดที่จะใช้สร้าง Scatter diagram คือ

(6, 9), (9, 7), (7, 2), ..., (6, 0), (0, -2)



$$\text{ก) } n = 28, \sum_{i=1}^n x_i = 31, \sum_{i=1}^n x_i^2 = 409, \bar{x} = 1.1071$$

$$\sum_{i=1}^{n-1} x_i = \sum_{i=1}^n x_i - x_n = 31 - (-2) = 33$$

$$\sum_{i=2}^n x_i = \sum_{i=1}^n x_i - x_1 = 31 - 6 = 25$$

$$\sum_{i=2}^n x_{i-1}x_i = (6)(9) + (9)(7) + \dots + (6)(0) + (0)(-2) = 244$$

$$\sum_{i=2}^n x_{i-1}x_i - \bar{x} \left[\sum_{i=2}^n x_i + \sum_{i=1}^{n-1} x_i \right] + (n-1) \bar{x}^2 = 212.8813$$

$$\sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2/n = 409 - \frac{(31)^2}{28} = 374.6786$$

$$\therefore r_1 = \frac{212.8813}{374.6786} = .5682$$

- 1) $H_0 : \rho_1 = 0$
- 2) $H_1 : \rho_1 \neq 0$
- 3) $\alpha = .05$
- 4) CR : $Z > z_{.025} = 1.96$ และ $Z < -1.96$
- 5) $z_c = \sqrt{n} r_1 = \sqrt{28} (.5682) = 3.0066$
- 6) $\because z_c > 1.96$, เราปฏิเสธ H_0 ที่ $\alpha = .05$ นั้นคือข้อมูลชุดนี้มีความสัมพันธ์กันในตัวข้อมูลเอง

5.12 Model checking และ Multiple linear regression

5.12.1 Examining the residuals

วิธีการนี้ไม่ได้ใช้สำหรับตัวแบบเชิงเส้นตรงเท่านั้น แต่อาจใช้ได้กับ model อื่น ๆ ด้วย เมื่อเราได้ทำการ fit ตัวแบบด้วยวิธีกำลังสองน้อยที่สุดแล้ว ส่วนของความแปรปรวนของ y ซึ่งไม่สามารถอธิบายได้ด้วยตัวแบบจะรวมอยู่ใน residuals

$$e_i = y_i - \hat{y}_i, i = 1, \dots, n$$

โดยที่ y_i คือค่า observed y และ $\hat{y}_i = b_0 + b_1 x_i$

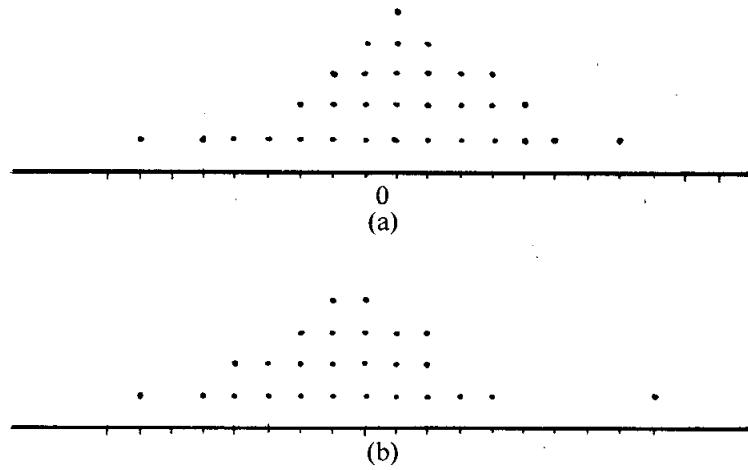
ในการวิเคราะห์การถดถอยที่ผ่านมาเรามีข้อสมมติเกี่ยวกับ random error (E_i) ว่า E_i เป็นอิสระต่อกัน และต่างก็มีการแจกแจงเป็น normal distribution ที่มี mean = 0 และ constant variance σ^2

นั่นคือ $E_i \sim NID(0, \sigma^2)$

ถ้าข้อสมมติเป็นจริงดังกล่าว เราจึงจะเชื่อถือผลการวิเคราะห์ได้ เมื่อตัวแบบถูกต้องเราใช้ $e_i = \hat{E}_i$ ดังนั้นในการตรวจสอบสมมติฐานข้างต้น เราจะสำรวจ residuals โดยวิธีการทางกราฟ

1. Histogram หรือ Dot diagram

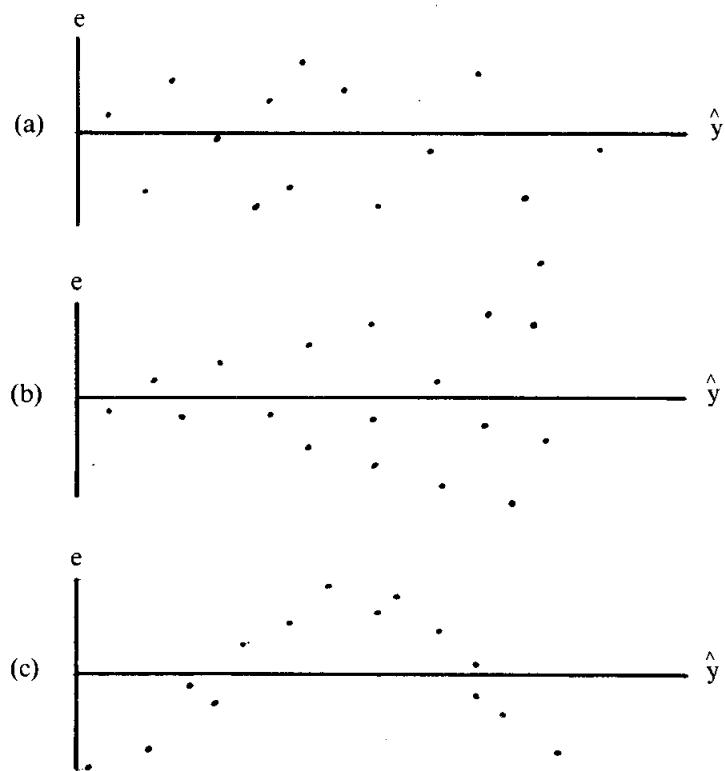
เพื่อดูลักษณะรวม ๆ ของ residuals เราจัด residuals เป็นแบบ grouped data สำหรับข้อมูลจำนวนมาก และ plot histogram สำหรับข้อมูลจำนวนน้อยเราใช้วิธี plot dot diagram บนเส้นตรง ตัวอย่างการ plot ในรูป (a) และ (b) รูป (a) แสดงว่าข้อมูลมีการแจกแจงคล้ายแบบปกติที่มี mean ใกล้ 0 ส่วนรูป (b) อาจพอสรุปได้ว่าข้อมูลมีการแจกแจงแบบปกติ แต่มีค่าสังเกตตัวหนึ่งอยู่ห่างไปจากกลุ่ม (wild observation) ในสถานการณ์เช่นนี้ควรจะได้พิจารณาข้อมูลตัวนั้นอีกครั้งเพื่อให้แน่ใจว่าข้อมูลนี้จะมีอิทธิพลต่อการ fit model หรือไม่อย่างไร (ซึ่งอาจพบว่าเป็นข้อมูลผิดปกติที่อาจต้องยกไปได้)



รูปที่ 5.8 Dot diagram of residuals

2. Residual vs. predicted value (\hat{y})

การ plot e_i กับ \hat{y}_i นั้น จะช่วยตรวจจับว่าข้อมูลเกี่ยวกับ constant error variance ว่าจริงหรือไม่



รูปที่ 5.9 Plot of residuals vs. predicted values

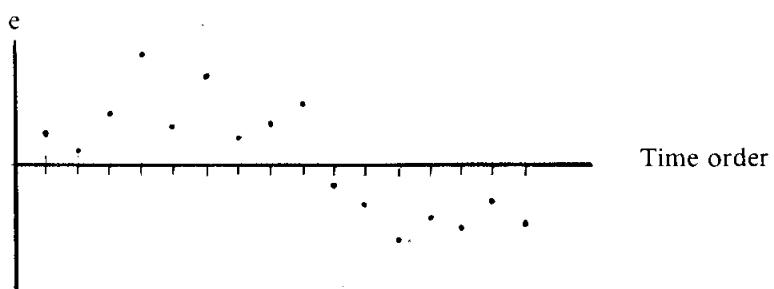
จากรูป (a) จะเห็นว่าจุดกระจายเป็นรูป horizontal band แสดงว่าไม่มีการบ่งชี้ว่า ขาดคุณสมบัติของ constant error variance

รูป (b) จะเห็นว่าจุดกระจายในลักษณะที่เมื่อ y มีค่ามากขึ้น e ; จะกระจายมากขึ้น แสดงว่า error variance มีแนวโน้มจะมากขึ้นเมื่อค่าของ y_i มากขึ้น ควรตรวจสอบคุณสมบัติของ constant error variance ทั้งนี้รวมถึงการแก้โดยการแบ่งข้อมูลเพื่อทำให้ variance คงที่

รูป (c) จะเห็นว่าจุดกระจายอย่างมีแบบแผน (systematic pattern) แทนที่จะกระจายอยู่อย่างสุ่มรอบ ๆ แกน y เริ่มด้วยการเพิ่มขึ้นอย่างคงที่แล้วลดลง ลักษณะนี้จะบ่งชี้ให้เรา สงสัยว่าตัวแบบไม่พอเพียงที่จะอธิบายข้อมูล เราอาจเพิ่มเทอมกำลังสอง หรือพิจารณาลองใช้ตัวแบบที่ไม่ใช่ตัวแบบเชิงเส้นตรง (nonlinear model)

3. Residual vs. Time order

ในการวิเคราะห์การทดสอบ การขาดข้อสมมติที่ถือว่าร้ายแรงคือการที่ error E_i ไม่เป็นอิสระต่อกัน การขาดความเป็นอิสระต่อกันมักจะเกิดขึ้นเมื่อประยุกต์ใช้ทางธุรกิจ และเศรษฐศาสตร์ โดยเกิดขึ้นเมื่อมีการรวมข้อมูลตามลำดับเวลาเพื่อทำการวิเคราะห์ การทดสอบในการทำนายแนวโน้มได้ ๆ ก็ตามเมื่อทำการทดลองในเวลาที่แบ่งไปอย่างมาก การ plot residuals กับลำดับเวลา (time order) มักจะทำให้เราตรวจจับการขาดคุณสมบัติ ของความเป็นอิสระได้ตัวอย่างในรูปที่ 5.10 เป็นรูปที่มีแบบแผน แสดงว่ากลุ่มของค่าที่สูง จะถูกตามมาด้วยกลุ่มของค่าที่ต่ำ



รูปที่ 5.10 Plot of residuals vs. time order

สิ่งนี้บ่งชี้ว่า residual ที่อยู่ต่อเนื่องกันจะมีความสัมพันธ์กัน ทำให้เราสงสัยได้ว่าเกิด การขาดคุณสมบัติของความเป็นอิสระ เราอาจตรวจสอบความเป็นอิสระของ error โดย การ plot residual ($n - 1$) คู่ในรูปคู่ลำดับ (e_{i-1}, e_i) , $i = 2, \dots, n$ ถ้าการกระจายของจุดไม่มี

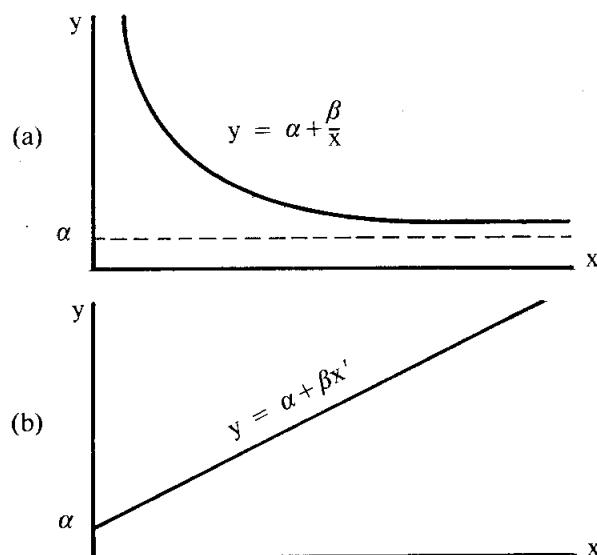
รูปแบบที่ชัดเจน ก็ถือได้ว่าข้อมูลเป็นอิสระกัน ถ้ามีการจับกลุ่มกันอยู่ในแนวเส้นตรง เราจะถือว่าขาดความเป็นอิสระระหว่างข้อมูลที่อยู่ติดกัน วิธีการนี้ได้กล่าวมาแล้วในหัวข้อ 5.11.

5.12.2 ความสัมพันธ์ที่ไม่เป็นเส้นตรงและการแปลงให้เป็นเส้นตรง

(Nonlinear relations and Linearizing Transformations)

เท่าที่ผ่านมาเราได้ศึกษาถึงรูปความสัมพันธ์ระหว่าง X กับ Y ที่เป็นสมการเส้นตรง หรือใช้ตัวแบบที่เป็นสมการเส้นตรง อย่างไรก็ได้ในชีวิตจริงรูปความสัมพันธ์ที่เราสนใจนั้น ไม่เป็นเส้นตรงเสมอไป จากการ plot scatter diagram จะพบว่ารูปความสัมพันธ์อาจใกล้ไปจากเส้นตรงมาก และเมื่อทำการคำนวณค่า r^2 จะได้ค่าเล็ก ๆ เมื่อเราใช้สมการเส้นตรง หรือเมื่อทดสอบ lack of fit ก็จะได้ผลสรุปว่ามี lack of fit วิธีการทางสถิติที่จะต้องปฏิบัติกับความสัมพันธ์ที่ไม่เป็นเส้นตรงนั้นยุ่งยากกว่าวิธีการจัดการกับความสัมพันธ์เชิงเส้นตรง ทั้งนี้ยกเว้นกับการจัดการกับตัวแบบที่เรียกว่า polynomial regression model ซึ่งจะกล่าวในหัวข้อ 5.12.4

ในบางสถานการณ์ เราอาจทำการแปลงตัวแปร X และ/หรือ Y เพื่อทำให้รูปความสัมพันธ์ใหม่ใกล้เคียงกับรูปเชิงเส้นตรง เราอาจเขียนตัวแบบเชิงเส้นตรงให้อยู่ในรูปของตัวแปรที่ถูกแปลงแล้ว และทำการวิเคราะห์กับข้อมูลที่ถูกแปลงแล้ว ในการวิเคราะห์นี้ ควรจะมีการตรวจสอบคุณสมบัติของ residuals ในตัวแบบที่ถูกแปลงแล้วด้วย เพราะว่า ข้อสมมุติที่ว่าตัว error เป็นอิสระต่อกัน และมีการแจกแจงเป็นแบบปกติที่มี variance คงที่นั้น จะต้องใช้กับตัวแบบที่ถูกแปลงแล้วเช่นกัน



รูปที่ 5.11

จากูปที่ 5.11

ถ้า $x =$ เวลาที่ฝึกฝน

$y =$ เวลาที่ทำงานเสร็จ

จะเห็นได้ว่าเวลาที่ทำงานเสร็จจะลดลงเมื่อใช้เวลาที่ฝึกฝนเพิ่มขึ้น แต่เมื่อถึงจุดหนึ่ง จะเห็นได้ว่า y ไม่ลดลงอีกแม้ว่าจะเพิ่มระยะเวลาที่ฝึกฝน

ดังนั้นเราอาจพิจารณาความสัมพันธ์ในรูป

$$y = \alpha + \frac{\beta}{x} \text{ ซึ่งมีรูปความสัมพันธ์ดังในรูป (a)}$$

รูป (a) แสดงว่าความสัมพันธ์ไม่เป็นเส้นตรง

แทนที่เราจะใช้ x เราจะทำการแปลงโดยใช้ $x' = \frac{1}{x}$ ดังนั้นตัวแบบจะกลายเป็น

$$y = \alpha + \beta x' \text{ ซึ่งมีรูปความสัมพันธ์ดังในรูป (b)}$$

การวิเคราะห์ข้อมูลจะทำกับ (y_i, x'_i) , $i = 1, \dots, n$ โดยที่ $x'_i = \frac{1}{x_i}$

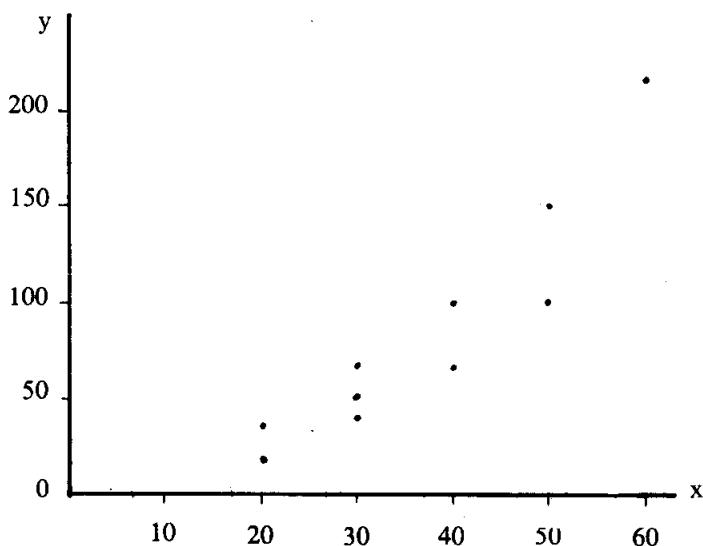
ในการ plot scatter diagram เพื่อดูว่าควรจะแปลงตัวแปรอย่างไรเพื่อให้ได้ความสัมพันธ์เชิงเส้นตรงนั้น อาจ plot (x_i, y_i) , $i = 1, \dots, n$ ลงบนกระดาษกราฟเช่น semilog paper (แกนใดแกนหนึ่งอยู่ใน log scale) หรือ double-log paper (ทั้ง 2 แกนอยู่ใน log scale) เพื่อดูว่าควรแปลงตัวใดตัวหนึ่งหรือทั้ง 2 ตัวจึงจะเหมาะสม ตัวอย่างเช่นความสัมพันธ์ของข้อ (a) ในตารางที่ 5.3 นั้น เมื่อ plot ลงบน semilog paper กราฟจะเป็นเส้นตรง บางครั้งถ้า diagram แสดงว่า ถ้า y เพิ่มเร็วมากเมื่อเทียบกับ x เราอาจใช้ \sqrt{y} หรือกำลังเศษส่วนอื่น ๆ ของ y เพื่อทำให้รูปความสัมพันธ์เป็นเส้นตรง

ตารางที่ 5.3
แสดงการแปลงรูปตัวแปรเพื่อให้ได้ความสัมพันธ์ในรูปเส้นตรง

Nonlinear model	Transformation	Transformed model
a) $y = ae^{bx}$	$y' = \ln y, x' = x$	$y' = \alpha + \beta x', \alpha = \ln a, \beta = b$
b) $y = ax^b$	$y' = \log y, x' = \log x$	$y' = \alpha + \beta x', \alpha = \log a, \beta = b$
c) $y = \frac{1}{a+bx}$	$y' = \frac{1}{y}, x' = x$	$y' = \alpha + \beta x', \alpha = a, \beta = b$
d) $y = \frac{1}{(a+bx)^2}$	$y' = \frac{1}{\sqrt{y}}, x' = x$	$y' = \alpha + \beta x', \alpha = a, \beta = b$
e) $\frac{1}{y} = a + \frac{b}{1+x}$	$y' = \frac{1}{y}, x' = \frac{1}{1+x}$	$y' = \alpha + \beta x', \alpha = a, \beta = b$
f) $y = a + b\sqrt{x}$	$y' = y, x' = \sqrt{x}$	$y' = \alpha + \beta x', \alpha = a, \beta = b$

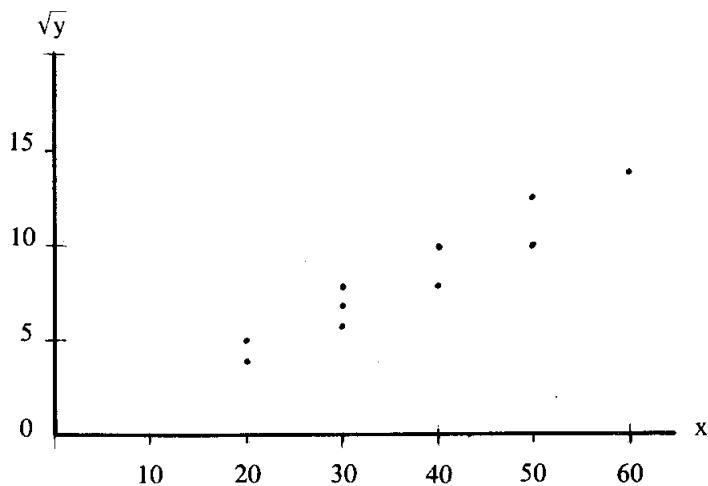
ตัวอย่างที่ 5.14 เพื่อทดสอบความสามารถสูงสุดของการหยุดของรถยนต์เมื่อเหยียบเบรค เดิมที่ รถ 10 คันได้ถูกขับด้วยอัตราความเร็วที่ระบุไว้สำหรับแต่ละคัน และวัดระยะทางที่ต้องใช้ก่อนหยุดรถได้สนใจ ความเร็วที่ขับสำหรับรถแต่ละคันและระยะทางที่ใช้หยุด เป็นดังนี้

Initial speed :	x	20 20 30 30 30 40 40 40 50 50 60
Stopping distance :	y	16.3 26.7 39.2 63.5 51.3 98.4 65.7 104.1 155.6 217.2



จากรูปจะเห็นว่าความสัมพันธ์ไม่เป็นเส้นตรง y เพิ่มขึ้นด้วยอัตราสูงที่ค่าสูง ๆ ของ x มากกว่าที่ค่าต่ำกว่าของ x เป็นการแนะนำว่าควรจะแปลงโดยใช้ $y' = \sqrt{y}$ หรือ y กำลังเศษส่วนอื่น ๆ การ plot (x_i, \sqrt{y}_i) ทำในรูปด้านไป จะเห็นได้ว่ารูปความสัมพันธ์เข้าใกล้เส้นตรงมากกว่ารูปแรก

x	20	20	30	30	30	40	40	50	50	60
$y' = \sqrt{y}$	4.037	5.167	6.261	7.969	7.162	9.920	8.106	10.203	12.474	14.738



เมื่อใช้ model $y' = \alpha + \beta x + E$

จากข้อมูลข้างต้น $\bar{x} = 37$, $\bar{y}' = 8.604$

$$S_{xx} = 1610, S_{y'y'} = 97.773, S_{xy'} = 381.621$$

$$\therefore \hat{\alpha} = -.166, \hat{\beta} = .237$$

ดังนั้น prediction equation คือ $\hat{y}' = -.166 + .237x$

$$\text{คำนวน } r^2 = \frac{\text{SSR}}{\text{SST}} = \frac{\hat{\beta} S_{xy'}}{S_{y'y'}} = \frac{.237(381.621)}{97.773} = .925$$

5.12.3 Multiple linear regression

Model : $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{3i} + E_i, i = 1, \dots, n$

โดยที่

x_{1i}, x_{2i}, x_{3i} เป็น fixed values ของตัวแปรอิสระ (independent variable) 3 ตัวในหน่วยทดลองที่ i

y_i คือ response

$E_i \sim NID(0, \sigma^2)$

$\beta_0, \beta_1, \beta_2, \beta_3$ เป็นพารามิเตอร์ที่เราจะประมาณค่า

Observations : $(x_{1i}, x_{2i}, x_{3i}, y_i), i = 1, \dots, n$

เราจะประมาณค่า $\beta_0, \beta_1, \dots, \beta_3$ ด้วยวิธีการกำลังสองน้อยที่สุด คือหาค่า b_0, b_1, \dots, b_3

เพื่อทำให้ $\sum_{i=1}^n e_i^2$ มีค่าน้อยที่สุด

$$S = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - b_0 - b_1 x_{1i} - \dots - b_3 x_{3i})^2$$

จากวิธีทาง calculus

$$\frac{\partial S}{\partial b_0} = 0$$

$$\frac{\partial S}{\partial b_1} = 0$$

⋮

$$\frac{\partial S}{\partial b_3} = 0$$

จะได้ normal equations ดังนี้

$$nb_0 + \sum x_{1i}b_1 + \sum x_{2i}b_2 + \sum x_{3i}b_3 = \sum y_i$$

$$\sum x_{1i}b_0 + \sum x_{1i}^2 b_1 + \sum x_{1i}x_{2i}b_2 + \sum x_{1i}x_{3i}b_3 = \sum x_{1i}y_i$$

$$\sum x_{2i}b_0 + \sum x_{1i}x_{2i}b_1 + \sum x_{2i}^2 b_2 + \sum x_{2i}x_{3i}b_3 = \sum x_{2i}y_i$$

$$\sum x_{3i}b_0 + \sum x_{1i}x_{3i}b_1 + \sum x_{2i}x_{3i}b_2 + \sum x_{3i}^2 b_3 = \sum x_{3i}y_i$$

หารืออยู่ในรูป

$$b_1 S_{x_1}^2 + b_2 S_{x_1 x_2} + b_3 S_{x_1 x_3} = S_{x_1 y}$$

$$b_1 S_{x_1 x_2} + b_2 S_{x_2}^2 + b_3 S_{x_2 x_3} = S_{x_2 y}$$

$$b_1 S_{x_1 x_3} + b_2 S_{x_2 x_3} + b_3 S_{x_3}^2 = S_{x_3 y}$$

$$b_0 = \bar{y} - b_1 \bar{x}_1 - b_2 \bar{x}_2 - b_3 \bar{x}_3$$

รูปทั่วไปของ multiple linear regression model ในรูปเมตริกซ์

$$\text{Model : } Y_{nx1} = X_{n \times (p+1)} \beta_{(p+1) \times 1} + E_{nx1}$$

โดยที่ Y คือ observed vector

$$Y' = [y_1 \ y_2 \ \dots \ y_n]$$

$$E' = [E_1 \ E_2 \ \dots \ E_n]$$

$$\beta' = [\beta_0 \ \beta_1 \ \dots \ \beta_p]$$

X = design matrix

$$X = \begin{bmatrix} 1 & x_{11} & \dots & x_{p1} \\ 1 & x_{12} & \dots & x_{p2} \\ \vdots & \vdots & & \vdots \\ 1 & x_{1n} & & x_{pn} \end{bmatrix}_{n \times (p+1)}$$

$$\text{Normal equations : } (X'X)\hat{\beta} = X'Y$$

โดยที่

$$X'X = \begin{bmatrix} n & \sum x_{1i} & \sum x_{2i} & \dots & \sum x_{pi} \\ \sum x_{1i} & \sum x_{1i}^2 & \sum x_{1i}x_{2i} & \dots & \sum x_{1i}x_{pi} \\ \vdots & \vdots & \vdots & & \vdots \\ \sum x_{pi} & \sum x_{1i}x_{pi} & \sum x_{2i}x_{pi} & \dots & \sum x_{pi}^2 \end{bmatrix}_{(p+1) \times (p+1)}, \quad X'Y = \begin{bmatrix} \sum y_i \\ \sum x_{1i}y_i \\ \vdots \\ \sum x_{pi}y_i \end{bmatrix}_{(p+1) \times 1}$$

$$\hat{\beta}' = [b_0 \ b_1 \ \dots \ b_p]$$

$$\text{ดังนั้น } \hat{\beta} = (X'X)^{-1}X'Y$$

ตัวอย่างที่ 5.15 เราสนใจว่าความดันโลหิต y มีความสัมพันธ์กับน้ำหนัก (x_1) และอายุ (x_2) ในกลุ่มผู้ชายที่มีความสูงเท่า ๆ กันหรือไม่ จากชาย 13 คนซึ่งมีน้ำหนักและอายุที่เราเลือกไว้ก่อน วัดความดันโลหิตของชายทั้ง 13 คนได้ข้อมูลในตารางที่ 5.4 ให้ fit multiple regression model

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + E_i, i = 1, \dots, 13$$

ตารางที่ 5.4

แสดงน้ำหนัก อายุและความดันโลหิตของชาย 13 คน

คนที่	x_1	x_2	y
1	152	50	120
2	183	20	141
3	171	20	124
4	165	30	126
5	158	30	117
6	161	50	129
7	149	60	123
8	158	50	125
9	170	40	132
10	153	55	123
11	164	40	132
12	190	40	155
13	185	20	147

คำนวณค่าต่าง ๆ จากข้อมูลในตารางที่ 5.4 ดังนี้ $(\sum_{i=1}^{13})$

$$\Sigma x_{1i} = 2159 \quad \Sigma x_{1i}^2 = 360639 \quad \Sigma x_{1i}x_{2i} = 82335$$

$$\Sigma x_{2i} = 505 \quad \Sigma x_{2i}^2 = 21925 \quad \Sigma x_{1i}y_i = 282921$$

$$\Sigma y_i = 1694 \quad \Sigma x_{2i}y_i = 65135$$

$$S_{x_1}^2 = \sum_{i=1}^{13} x_{1i}^2 - (\sum_{i=1}^{13} x_{1i})^2 / 13 = 360639 - \frac{(2159)^2}{13} = 2078.9$$

$$S_{x_2}^2 = \sum_{i=1}^{13} x_{2i}^2 - (\sum_{i=1}^{13} x_{2i})^2 / 13 = 21925 - \frac{(505)^2}{13} = 2307.7$$

$$S_{x_1 x_2} = \sum_{i=1}^{13} x_{1i} x_{2i} - \frac{(\sum x_{1i})(\sum x_{2i})}{13} = 82335 - \frac{(2159)(505)}{13} = -1533.846$$

$$S_{x_1 y} = \sum_{i=1}^{13} x_{1i} y_i - \frac{(\sum x_{1i})(\sum y_i)}{13} = 282921 - \frac{(2159)(1694)}{13} = 1586.7$$

$$S_{x_2 y} = \sum_{i=1}^{13} x_{2i} y_i - \frac{(\sum x_{2i})(\sum y_i)}{13} = 65135 - \frac{(505)(1694)}{13} = -670.38$$

$$\bar{y} = 1694/13 = 130.308$$

$$\bar{x}_1 = 2159/13 = 166.077$$

$$\bar{x}_2 = 505/13 = 38.846$$

ตั้งนี้ Normal equations คือ

$$2078.9b_1 - 1533.8b_2 = 1586.7 \quad \dots \dots (1)$$

$$-1533.8b_1 + 2307.7b_2 = -670.38 \quad \dots \dots (2)$$

$$b_0 = 130.308 - 166.077b_1 - 38.846b_2 \quad \dots \dots (3)$$

จาก (1) และ (2) ได้ $b_1 = 1.077, b_2 = .425$

$$\therefore b_0 = -65.10$$

prediction equation คือ

$$\hat{y} = -65.10 + 1.077x_1 + .425x_2$$

หมายความว่า ความดันโลหิตเฉลี่ยจะเพิ่มขึ้น 1.077 หน่วย ถ้าอายุหนักเพิ่มขึ้น 1 หน่วย โดยที่อายุคงที่ และถ้าอายุคงที่แล้วอายุเพิ่มขึ้น 1 ปี จะทำให้ความดันเฉลี่ยเพิ่มขึ้น .425 หน่วย

5.12.4 Polynomial regression model : Quadratic function (Second degree polynomial)

เมื่อเรา plot scatter diagram แล้วความสัมพันธ์ไม่อยู่ในรูปเส้นตรง (nonlinear) เราอาจทำการแปลงตัวแปรตามวิธีที่กล่าวมาแล้วในหัวข้อ 5.12.2 แต่บางกรณีเราไม่อาจทำดังกล่าวได้ วิธีหนึ่งที่จะจัดการกับข้อมูลได้คือใช้ตัวแปร x ที่มีกำลังสูงกว่ากำลังหนึ่ง เช่น

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + E_i, i=1, \dots, n$$

ถ้าตัดเทอม E_i 去 กะจะเห็นว่าเป็น quadratic function ของตัวแปรอิสระ x เราเรียกตัวแบบรูปนี้ว่าตัวแบบการถดถอยรูป polynomial ของ y กับ x กำลังสูงสุดของ x จะเป็นตัวบวก *degree* หรือ *order* ของการถดถอยรูป polynomial และในการวิเคราะห์การถดถอยรูปนี้เมื่อได้ใช้วิธีพิเศษนอกเหนือไปจากวิธีของ multiple regression เลย

ตัวแบบที่ใช้คือ

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + E_i, i=1, \dots, n$$

โดยที่ $x_{1i} = x_i$ และ $x_{2i} = x_i^2$ นั้นเอง

ตัวอย่างที่ 5.16 จากข้อมูลที่กำหนดให้ จง fit quadratic model $Y = \beta_0 + \beta_1 x + \beta_2 x^2 + E$
โดยวิธีกำลังสองน้อยที่สุด

x	-2	-1	0	1	2
y	.4	1.3	2.2	2.5	3.0

	$x_1 = x$	$x_2 = x^2$	y	x_1^2	x_2^2	$x_1 x_2$	$x_1 y$	$x_2 y$
-2	4	.4	4	16	1	-8	-.8	1.6
-1	1	1.3	1	1	1	-1	-1.3	1.3
0	0	2.2	0	0	0	0	0	0
1	1	2.5	1	1	1	1	2.5	2.5
2	4	3.0	4	16	4	8	6.0	12.0
Total	0	10	9.4	10	3.4	0	6.4	17.4

คำนวณค่าต่าง ๆ ดังนี้ $\bar{x}_1 = 0$, $S_{x_1}^2 = 10$, $S_{x_1 y} = 6.4$

$\bar{x}_2 = 2$, $S_{x_2}^2 = 14$, $S_{x_2 y} = -1.4$

$\bar{y} = 1.88$, $S_{x_1 x_2} = 0$

Normal equations គឺ

$$10b_1 = 6.4 \quad \dots \dots (1)$$

$$14b_2 = -1.4 \quad \dots \dots (2)$$

$$b_0 = 1.88 - 2b_2 \quad \dots \dots (3)$$

តួនាទី $b_1 = .64, b_2 = -.1, b_0 = 2.08$

prediction equation គឺ

$$\hat{y} = 2.08 + .64x - .1x^2$$

តាមវិធីមេត្រិកចំ

$$\text{Design matrix } \underline{X} = \begin{bmatrix} 1 & -2 & 4 \\ 1 & -1 & 1 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 2 & 4 \end{bmatrix}, \underline{Y} = \begin{bmatrix} .4 \\ 1.3 \\ 2.2 \\ 2.5 \\ 3.0 \end{bmatrix}$$

$$\underline{X}'\underline{X} = \begin{bmatrix} 5 & 0 & 10 \\ 0 & 10 & 0 \\ 10 & 0 & 34 \end{bmatrix}, \underline{X}'\underline{Y} = \begin{bmatrix} 9.4 \\ 6.4 \\ 17.4 \end{bmatrix}$$

$$\hat{\beta} = (\underline{X}'\underline{X})^{-1}\underline{X}'\underline{Y}$$

$$= \begin{bmatrix} \frac{34}{70} & 0 & -\frac{1}{7} \\ 0 & \frac{1}{10} & 0 \\ -\frac{1}{7} & 0 & \frac{5}{70} \end{bmatrix} \begin{bmatrix} 9.4 \\ 6.4 \\ 17.4 \end{bmatrix}$$

$$= \begin{bmatrix} 2.08 \\ .64 \\ -.1 \end{bmatrix} = \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix}$$

$$\therefore \hat{y} = 2.08 + .64x - .1x^2$$

แบบฝึกหัด

- 5.1 พนักงานป่าไม้ต้องการทราบว่าเข่าจะใช้เส้นผ่าศูนย์กลางในการคำนวณความสูงของต้นไม้ได้หรือไม่ เขายังทำการวัดเส้นผ่าศูนย์กลางของต้นไม้ต่างส่วนสูง 4.5 ฟุตจากพื้นดิน และความสูงของต้นเมเปลที่ให้น้ำตาล 12 ต้น ผลปรากฏ ดังนี้

x : เส้นผ่าศูนย์กลาง (นิ้ว)	.9	1.2	2.3	3.1	3.3	3.9	4.3	6.2	9.6	12.6	16.1	25.8
y : ความสูง (ฟุต)	18	26	32	36	44.5	35.6	40.5	57.5	67.3	84	67	87.5

- ก) จงสร้าง Scatter diagram เพื่อถูกว่าควรใช้ Straight line model หรือไม่
ข) ถ้ากำหนดให้ $x' = \log x$ และ $y' = \log y$ จงสร้าง Scatter diagram ของ (x', y')

x'	-.05	.08	.36	.49	.52	.59	.63	.79	.98	1.10	1.21	1.41
y'	1.26	1.41	1.51	1.56	1.65	1.55	1.61	1.76	1.83	1.92	1.83	1.94

- ค) จาก model $y' = \beta_0 + \beta_1 x' + E$, จงหา Sample regression equation

- 5.2 ใน การเตรียมการที่จะคัดเลือกนักศึกษาปีที่ 1 เข้าเรียนวิชาเอกคณิตศาสตร์ โดยจะใช้วิธีการสอบวิชาคณิตศาสตร์พื้นฐาน (คะแนนเต็ม 10 คะแนน) เพื่อจะศึกษาถูกว่าการใช้คะแนนวิชานี้จะคำนวณผลสำเร็จในการเรียนวิชาเอกได้หรือไม่ วิทยาลัยแห่งหนึ่งจึงได้เอ่านักศึกษาที่จบแล้ว และได้คะแนนวิชาพื้นฐานต่าง ๆ กันมา 8 คน แล้วจดคะแนนตอนจบปริญญาตรี (Grade point average) ของแต่ละคนดังนี้

x : คะแนนคณิตศาสตร์พื้นฐาน	4	6	8	6.5	3.5	8	7	5
y : คะแนนจบปริญญาตรี	2.1	3.1	3.4	2.9	2.1	3.3	2.8	2.3

กำหนดค่าที่คำนวณจากข้อมูลดังนี้

$$\sum x_i = 46, \sum x_i^2 = 308.5, \bar{x} = 6, \sum x_i y_i = 137.9$$

$$\sum y_i = 22, \sum y_i^2 = 62.42, \bar{y} = 2.75$$

$$\sum (x_i - \bar{x})^2 = 20.5, \sum (y_i - \bar{y})^2 = 1.92, \sum (x_i - \bar{x})(y_i - \bar{y}) = 5.9$$

สมมุติว่าความสัมพันธ์ของ X และ Y อยู่ในรูป $\mu_{Y|x} = \beta_0 + \beta_1 x$

- ก) จงหา Least-squared estimated ของ β_0 และ β_1

- ข) จงเขียน prediction equation

- ก) จงหา point estimate ของ $\mu_{Y|X=3}$ และอธิบายความหมาย และถ้านายก.ได้คะแนน
วิชาพื้นฐาน 5.5 คะแนน จงประมาณว่านายก.จะเรียนจบด้วยคะแนนเฉลี่ยเท่าใด
- ก) จงทดสอบ $H_0 : \beta_1 = 0$ โดยใช้ F-test กำหนด $\alpha = .01$
- ก) จงหา $100(1-\alpha)\%$ และอธิบายความหมาย
- ก) ถ้าความสัมพันธ์ของ X และ Y อยู่ในรูปนี้จริง ค่าประมาณของ $\sigma^2 = ?$

5.3 ในปี 2519 บริษัทผู้ผลิตผงซักฟอกมีห้องต่าง ๆ มีรายจ่ายในการโฆษณาและยอดขาย
ดังนี้

ปีห้อง	x : รายจ่ายในการโฆษณา (ล้านบาท)	y : ยอดขาย (ล้านบาท)
บริษัท	7	9
รุน	9	12
นิส	6	8
ควิก	5	7
แฟบ	3	4

สมมุติรูปความสัมพันธ์ คือ $\mu_{Y|x} = \beta_0 + \beta_1 x$

- ก) จงเขียน prediction equation
- ข) จงอธิบายค่า $\hat{\beta}_0 = b_0$ และ $\hat{\beta}_1 = b_1$ ที่หาได้
- ก) ถ้าบริษัทผลิตกวิกมีรายจ่ายในการโฆษณา 4 ล้านบาท จงประมาณว่าจะได้ยอดขาย
กี่ล้านบาท
- ก) จงทดสอบที่ $\alpha = .05$ $H_0 : \beta_1 = 0$ โดยใช้ F-test
- ก) ถ้าหา $95\% \text{ C.I.}$ ของ $\mu_{Y|x}$ ได้ $(7.631, 8.369)$ จงอธิบายความหมาย

5.4 จากข้อมูลที่กำหนดให้

- ก) จงสร้าง Scatter diagram
- ข) จงเขียน prediction equation ของ model : $\mu_{Y|x} = \beta_0 + \beta_1 x$
- ก) จงทดสอบ lack of fit ที่ $\alpha = .05$

x	1	1	1	2	3	3	4	5	5
y	9	7	8	10	15	12	19	24	21

5.5 ปริมาณสารเคมี y (กรัม) ซึ่งละลายในน้ำ 100 กรัม ณ อุณหภูมิต่าง ๆ กันเป็นดังนี้

x ($^{\circ}\text{C}$)	y (กรัม)
0	8, 6, 8
15	12, 10, 14
30	25, 21, 24
45	31, 33, 28
60	44, 39, 42
75	48, 51, 44

ก) จงสร้าง Scatter diagram

ข) จงเขียน prediction equation ของ model : $\mu_{Y/x} = \beta_0 + \beta_1 x$

ค) จงประมาณปริมาณของสารเคมีที่จะละลายในน้ำ 100 กรัมที่อุณหภูมิ 50°C

ง) จงทดสอบที่ $\alpha = .01$ $H_0 : \beta_0 = 6$, $H_1 : \beta_0 \neq 6$

จ) จงทดสอบว่า Straight line model ที่ใช้ยังพอดีเพียงหรือไม่ นั่นคือทดสอบ $H_0 :$ ไม่มี lack of fit

5.6 กำหนดข้อมูลให้ต่อไปนี้

x	2.1	2.6	2.7	2.7	2.7	3.0	3.0	3.0	3.1	3.1	3.2	3.3	3.6	3.8
y	2.6	2.8	2.9	3.6	3.6	3.1	3.4	3.1	3.9	3.1	4.1	3.6	4.0	3.6

$$\sum x_i = 41.9, \sum y_i = 47.4, \sum x_i y_i = 143.59$$

$$\sum x_i^2 = 127.79, \sum y_i^2 = 163.26, n\bar{y}^2 = 160.4829$$

$$\sum (x_i - \bar{x})^2 = 2.3893, \sum (y_i - \bar{y})^2 = 2.7771, \sum (x_i - \bar{x})(y_i - \bar{y}) = 1.7286$$

Simple regression model : $\mu_{Y/x} = \beta_0 + \beta_1 x$ หรือ $Y = \beta_0 + \beta_1 x + E$

กำหนด $\hat{\beta}_0 = b_0 = 1.2203$ และ

$$\hat{\beta}_1 = b_1 = 0.7235$$

ก) จงเขียน prediction equation หรือ Fitted regression equation

ข) จงสร้างตาราง ANOVA เพื่อทดสอบที่ $\alpha = .05$

1) $H_0 : \beta_1 = 0$

2) H_0 : ไม่มี lack of fit , เรายังต้องปรับปรุง model ใหม่หรือไม่

3) คำนวณ $100 r^2\%$

- 5.7 สุ่มคนไข้ที่เป็นโรคชนิดหนึ่งมา 12 คน หลังจากที่ได้ทำการรักษาด้วยยาชนิดหนึ่งแล้ว 3 วัน ได้วัดอุณหภูมิ (X) และจำนวนการเต้นของชีพจรต่อ 1 นาทีของคนไข้แต่ละคน (Y) ผลปรากฏดังนี้

x : อุณหภูมิ	99.4	100.1	98.6	97.9	101.4	102.6
y : จำนวนการเต้นของชีพจร/นาที	78	83	70	76	75	89

x : อุณหภูมิ	99.0	98.6	101.1	98.0	100.5	99.3
y : จำนวนการเต้นของชีพจร/นาที	80	77	86	79	84	87

$$\sum x_i = 1196.5, \quad \sum y_i = 964.0$$

$$\sum x_i^2 = 119324.37, \quad \sum y_i^2 = 77786$$

$$\sum x_i y_i = 96170.2$$

ก) คำนวณค่า r (sample correlation coefficient)

ข) ทดสอบที่ $\alpha = .05 H_0 : \rho = 0, H_1 : \rho \neq 0$

ค) ทดสอบที่ $\alpha = .05 H_0 : \rho = .7, H_1 : \rho \neq .7$

ง) จงหา 95% C.I. ของ ρ

- 5.8 สุ่มนักเรียนมา 8 คน แต่ละคนวัด I.Q. (X) และวัดความเร็วในการอ่าน (Y) ผลปรากฏดังนี้

x : I.Q.	102	95	115	112	86	102	121	108
y : ความเร็ว	14.7	12.0	17.1	15.8	10.1	11.6	19.8	17.5

ก) คำนวณค่า r

ข) ทดสอบ $H_0 : \rho = 0, H_1 : \rho \neq 0$ ที่ $\alpha = .05$

ค) จงหา 95% C.I. ของ ρ

5.9 จากคะแนนวิชาคณิตศาสตร์ และคะแนนภาษาอังกฤษ ของนักศึกษา 6 คนที่สุ่มมาเป็นตัวอย่าง

x : คะแนนคณิตศาสตร์	70	92	80	74	65	83
y : คะแนนภาษาอังกฤษ	74	84	63	87	78	90

- ก) จงคำนวณค่า r
- ข) ทดสอบ $H_0: \rho = 0$, $H_1: \rho \neq 0$ ที่ $\alpha = .10$
- ค) จากข้อ ข) ถ้าเราไม่สามารถปฏิเสธ H_0 เราจะสรุปผลการทดสอบว่าอย่างไร

5.10 ต่อไปนี้เป็นความสูงและน้ำหนักของมิสอเมริกา 17 คน

x : ความสูง (นิ้ว)	65	67	66	65.5	65	66.5	66	67	66
y : น้ำหนัก (ปอนด์)	114	120	116	118	115	124	124	115	116
x : ความสูง	69	67	65.5	68	67	68	69	68	
y : น้ำหนัก (ปอนด์)	135	125	110	121	118	120	125	119	

- ก) จงสร้าง Scatter diagram
- ข) คำนวณค่า r และอธิบายความหมาย
- ค) ทดสอบที่ $\alpha = .10$ $H_0: \rho = 0$

5.11 จากข้อมูลที่กำหนดให้ 2 ชุด จงทดสอบ $H_0: \rho_1 = \rho_2$ ที่ $\alpha = .05$

x_{1i}	.20	.25	.25	.30	.40	.50	.50
y_{1i}	30	26	40	35	54	56	65

x_{2i}	.20	.25	.30	.40	.40	.50
y_{2i}	23	24	42	49	55	70

5.12 จากการวัดความถี่ของเสียงที่สีแยกหนึ่งในเวลาที่ต่อเนื่องกันได้ข้อมูลดังนี้

65, 64, 63, 61, 60, 58, 63, 64, 62, 64, 63, 63, 62, 60, 62, 64, 66, 68, 68, 69

จงทำการสำรวจเพื่อตรวจดูว่าอาจมี lag 1 independence หรือไม่โดย

- ก) สร้าง Scatter diagram
- ข) คำนวณ r_1 (first serial correlation coefficient)
- ค) ทดสอบ $H_0 : \rho_1 = 0$, $H_1 : \rho_1 \neq 0$ ที่ $\alpha = .05$

$$\text{กำหนด } \sum_{i=1}^n x_i^2 = 80667, \sum_{i=2}^n x_{i-1}x_i = 76140, \sum_{i=1}^n x_i = 1269$$

5.13 จากข้อมูลที่วัดต่อเนื่องกันของปริมาณชัลเพอร์ของโลหะชนิดหนึ่งจากเครื่องทำความร้อนในบ้าน เป็นดังนี้

2.0, 2.8, 2.0, 1.9, 1.1, 1.4, 1.2, 2.4, 2.5, 2.3, 3.1, 5.2, 2.8, 2.7, 2.6, 2.4, 3.3, 3.9, 3.0, 3.8, 2.9

- ก) จงสร้าง Scatter diagram
- ข) จงคำนวณ r_1
- ค) จงทดสอบ $H_0 : \rho_1 = 0$, $H_1 : \rho_1 \neq 0$ ที่ $\alpha = .05$

$$\text{กำหนด } \sum_{i=1}^n x_i^2 = 163.61, \sum_{i=2}^n x_{i-1}x_i = 148.48, \sum_{i=1}^n x_i = 55.3$$