



ทฤษฎีที่ใช้ในการวิเคราะห์ปัญหา

ในการหาความสัมพันธ์ของข้อมูล ๒ พวก โดยที่พวกหนึ่งเป็นตัวแปรอิสระ (independent-variable) คือตัวแปรที่มีการเปลี่ยนแปลงค่าไม่ขึ้นกับตัวแปรอื่น กับอีกพวกหนึ่งเป็นตัวแปรตาม (dependent variable) ซึ่งหมายถึงตัวแปรที่มีการเปลี่ยนแปลงค่าไปตามตัวแปรอิสระ และอาจเป็นการเปลี่ยนแปลงไปในทางเดียว หรือในทางกลับกันก็ได้ นั้น เราใช้ทฤษฎีความถดถอย (Regression Theory) เป็นเครื่องมือในการพิจารณา ซึ่งนอกจากจะโคจรมาถึงความสัมพันธ์ระหว่างข้อมูลสองพวกนั้นแล้ว ยังสามารถใช้สมการถดถอยทำนายค่าของตัวแปรตาม โดยอาศัยค่าของตัวแปรอิสระเป็นตัวทำนายได้อีกด้วย สำหรับข้อมูลในปัญหาที่จะทำการวิจัยนี้ เมื่อโคลงนำไปเขียน แผนภาพการกระจาย (scatter diagram) ของตัวแปรอิสระกับตัวแปรตามแล้วพบว่า ความสัมพันธ์ของตัวแปรทั้งสองพวกจัดเป็นความสัมพันธ์ในเชิงเส้นตรง ดังนั้นทฤษฎีที่จะใช้ในการวิเคราะห์ปัญหานี้ จึงใช้ทฤษฎีความถดถอยในรูปของสมการเส้นตรง (Linear Regression Theory) ^(๒) ซึ่งมีแบบจำลองทางคณิตศาสตร์ (Mathematical Model) ดังนี้

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + e \quad (1)$$

เมื่อ Y เป็นตัวแปรตาม

X_i เป็นตัวแปรอิสระตัวที่ i , ($i = 1, 2, 3, \dots, p$)

β_0 เป็นจุดตัดบนแกน Y ของเส้นถดถอย (regression line) หรือเป็นค่าเฉลี่ยของ Y ซึ่งแยกเอาความเปลี่ยนแปลงของค่า Y เนื่องจากค่า X_i ต่าง ๆ ออกไปแล้ว

β_i เป็นค่าสัมประสิทธิ์ของความถดถอย (regression coefficient) ($i = 1, 2, 3, \dots, p$) ซึ่งแสดงถึง อัตราการเปลี่ยนแปลงของ Y เมื่อเทียบกับตัวแปรอิสระตัวที่ i ในขณะที่กำหนดให้ตัวแปรอิสระตัวอื่น ๆ คงที่

(๒) Draper & Smith, Applied Regression Analysis, John Wiley & Sons Inc.

- ๑ เป็นความคลาดเคลื่อนซึ่งเนื่องจากการที่ ค่าประมาณของตัวแปรตาม Y สำหรับตัวแปรอิสระแต่ละชุด ที่ได้จากสมการถดถอย แตกต่างไปจากค่าจริง มีลักษณะเป็นตัวแปรสุ่ม (random variable) อาจเขียนได้ว่า $e = Y - \hat{Y}$ เมื่อ \hat{Y} เป็นค่าประมาณของตัวแปรตาม Y ที่ได้จากสมการถดถอย โดยมีข้อสมมติว่า $E(e) = 0 ; V(e) = V(Y) = \sigma^2 ; COV(e_i, e_k) = 0 , i \neq k ; e \sim N(0, \sigma^2)$

จาก (1) ถ้าเราจะใช้ค่าตัวแปรอิสระต่าง ๆ ทำนายค่าตัวแปรตาม จะได้สมการเป็น

$$\hat{Y} = b_0 + b_1X_1 + b_2X_2 + \dots + b_pX_p \quad (2)$$

เมื่อ \hat{Y} เป็นค่าประมาณของ Y สำหรับ X_i แต่ละชุด

b_0 เป็นค่าประมาณของ β_0

b_i เป็นค่าประมาณของ $\beta_i , (i = 1, 2, 3, \dots, p)$

การคำนวณค่าจุดตัดบนแกน Y (b_0) และค่าสัมประสิทธิ์ของความถดถอย (b_i)

ในการคำนวณค่า b_0 และค่าสัมประสิทธิ์ของความถดถอย $b_i (i = 1, 2, \dots, p)$ เพื่อให้ค่าที่ได้ออกมา เมื่อนำไปใช้ในสมการถดถอยแล้วทำให้ค่าประมาณของ Y ที่ได้ มีความคลาดเคลื่อนน้อยที่สุด วิธีคำนวณจึงใช้วิธีกำลังสองน้อยที่สุด (Method of Least Squares) โดยให้ผลบวกของกำลังสองของความคลาดเคลื่อน มีค่าน้อยที่สุด

จาก (1) เมื่อมีข้อมูลจำนวน n ชุด ($j = 1, 2, \dots, n$) จะได้สมการเป็น

$$Y_j = \beta_0 + \beta_1X_{1j} + \beta_2X_{2j} + \dots + \beta_pX_{pj} + e_j \quad (3)$$

$$\therefore e_j = Y_j - \beta_0 - \beta_1X_{1j} - \beta_2X_{2j} - \dots - \beta_pX_{pj}$$

เมื่อรวมกำลังสองของความคลาดเคลื่อนทั้งข้อมูล จะได้เป็น

$$\sum_{j=1}^n e_j^2 = \sum_{j=1}^n (Y_j - \beta_0 - \beta_1X_{1j} - \beta_2X_{2j} - \dots - \beta_pX_{pj})^2 \quad (4)$$

จาก (4) ทำการแกสมการโดยอาศัยวิธี พหุคูณดิริเวทีฟ (partial derivative) เทียบกับ β_i ($i = 1, 2, \dots, p$) และ β_0 แล้วเทียบให้เท่ากับ 0 ซึ่งจะไดสมการทั้งหมด ($p+1$) สมการ จากสมการเหล่านี้ เมื่อแกสมการแล้วก็จะไดค่าประมาณของ β_0 และ β_i ออกมาคือ b_0, b_i ($i = 1, 2, 3, \dots, p$)

จากสมการซึ่งทำ พหุคูณดิริเวทีฟ เทียบกับ β_0 และให้เท่ากับ 0 จะทำให้ไดค่า b_0 ดังต่อไปนี้

$$\begin{aligned} \frac{\partial \left(\sum_{j=1}^n e_j^2 \right)}{\partial b_0} &= -2 \sum_{j=1}^n (Y_j - b_0 - b_1 X_{1j} - b_2 X_{2j} - \dots - b_p X_{pj}) = 0 \\ \therefore \sum_{j=1}^n b_0 &= \sum_{j=1}^n (Y_j - b_1 X_{1j} - b_2 X_{2j} - \dots - b_p X_{pj}) \\ b_0 &= \bar{Y} - b_1 \bar{X}_1 - b_2 \bar{X}_2 - \dots - b_p \bar{X}_p \quad (5) \end{aligned}$$

จะเห็นไดวาค่า b_0 เป็นฟังก์ชันของ b_i , ($i = 1, 2, \dots, p$) ้วยเหตุนี้ถ้าเราแทนค่า b_0 ใน (4) เสียก่อนที่จะคำนวณ พหุคูณดิริเวทีฟ โดยเทียบกับ b_i , ($i = 1, 2, \dots, p$) แต่ละตัว ่ต่อไป ก็จะทำให้เหลือค่าที่จะต้องคำนวณเพียง P ค่า แทนที่จะเป็น $P+1$ ค่า ดังนี้

$$\sum_{j=1}^n e_j^2 = \sum_{j=1}^n \left[Y_j - \bar{Y} - b_1 (X_{1j} - \bar{X}_1) - b_2 (X_{2j} - \bar{X}_2) - \dots - b_p (X_{pj} - \bar{X}_p) \right]^2$$

เมื่อให้ $y = (Y_j - \bar{Y})$, $x_1 = (X_{1j} - \bar{X}_1)$, $x_2 = (X_{2j} - \bar{X}_2)$,
....., $x_p = (X_{pj} - \bar{X}_p)$

ก็จะสามารถเขียนสมการใหม่ได้ดังนี้

$$\sum_{j=1}^n e_j^2 = \sum_{j=1}^n (y - b_1 x_1 - b_2 x_2 - \dots - b_p x_p)^2 \quad (6)$$

จากนั้นก็ทำ พหุคูณดิริเวทีฟ ของสมการ (6) เทียบกับ b_i , ($i = 1, 2, \dots, p$) ทีละตัว จะได P สมการ ซึ่งเมื่อแกสมการเหล่านี้ก็จะไดค่า b_1, b_2, \dots, b_p ออกมา และนำไปหาค่า b_0 ได้

เมื่อจัดเป็นรูปเมทริกซ์จะได้สมการเป็น

$$Y = X\beta + e \quad (9)$$

เมื่อ

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} (n \times 1) \quad X = \begin{bmatrix} 1 & X_{11} & X_{21} & \dots & X_{p1} \\ 1 & X_{12} & X_{22} & \dots & X_{p2} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & X_{1n} & X_{2n} & \dots & X_{pn} \end{bmatrix} (n \times (p+1))$$

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} (p+1) \times 1 \quad e = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} (n \times 1)$$

โดยมีข้อสมมติเบื้องต้นดังนี้

$$E(e) = (0), \quad V(e) = V(Y) = I\sigma^2, \quad e \sim N(0, I\sigma^2)$$

และสำหรับค่าประมาณของ Y คือ \hat{Y} ก็จะมีสมการเป็น

$$\hat{Y} = Xb \quad (10)$$

เมื่อ

$$\hat{Y} = \begin{bmatrix} \hat{Y}_1 \\ \hat{Y}_2 \\ \vdots \\ \hat{Y}_n \end{bmatrix} (n \times 1) \quad b = \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_p \end{bmatrix} [(p+1) \times 1]$$

เมื่อพิจารณาจาก (10) จะได้ normal equations เพื่อนำไปหาค่า b ดังนี้

$$\begin{aligned} X'Xb &= X'Y \\ [X'X]^{-1}X'Xb &= [X'X]^{-1}X'Y \\ b [(p+1) \times 1] &= [X'X]^{-1}_{(p+1)(p+1)} \cdot X'Y_{(p+1) \times 1} \end{aligned} \quad (11)$$

จากสมการ (11) เราจะได้ออกค่า b_0, b_1, \dots, b_p ออกมา จะเห็นได้ว่าจาก normal equations ในสมการ (7) เราได้ออกค่า b_1, b_2, \dots, b_p เท่านั้น ค่า b_0 ต้องคำนวณใหม่

อีกครั้ง ดั่งนี้การหาค่าในรูปแมทริกซ์ จะให้ผลดีกว่าการแก้สมการ normal equations โดยตรง เพราะสามารถหาค่า b_0 ออกมาได้พร้อมกับค่า b_i แต่วิธีแมทริกซ์นี้จะใช้ได้เมื่อ iff (XX') เป็น nonsingular matrix

$$[XX']_{(p+1)(p+1)} = \begin{bmatrix} n & \sum X_{1j} & \sum X_{2j} & \dots & \sum X_{pj} \\ \sum X_{1j} & \sum X_{1j}^2 & \sum (X_{1j}X_{2j}) & \dots & \sum (X_{1j}X_{pj}) \\ \sum X_{2j} & \sum (X_{2j}X_{1j}) & \sum X_{2j}^2 & \dots & \sum (X_{2j}X_{pj}) \\ \dots & \dots & \dots & \dots & \dots \\ \sum X_{pj} & \sum (X_{pj}X_{1j}) & \sum (X_{pj}X_{2j}) & \dots & \sum X_{pj}^2 \end{bmatrix}_{(p+1)(p+1)}$$

$$[X'Y]_{(p+1) \times 1} = \begin{bmatrix} \sum Y_j \\ \sum (X_{1j}Y_j) \\ \sum (X_{2j}Y_j) \\ \dots \\ \sum (X_{pj}Y_j) \end{bmatrix}_{(p+1) \times 1}$$

การหา Inverse Matrix $[XX']^{-1}$ ใช้วิธี Abbreviated Doolittle

(ดูจากภาคผนวก)

$[XX']^{-1}_{(p+1)(p+1)}$ จะเป็น symmetric matrix ด้วย

$$[XX']^{-1} = C = \begin{bmatrix} c_{00} & c_{01} & c_{02} & \dots & c_{0p} \\ c_{10} & c_{11} & c_{12} & \dots & c_{1p} \\ \dots & \dots & \dots & \dots & \dots \\ c_{p0} & c_{p1} & c_{p2} & \dots & c_{pp} \end{bmatrix}$$

ดังนั้นจะได้

$$\begin{aligned} b_0 &= c_{00} \sum Y_j + c_{01} \sum (X_{1j}Y_j) + c_{02} \sum (X_{2j}Y_j) + \dots + c_{0p} \sum (X_{pj}Y_j) \\ b_1 &= c_{10} \sum Y_j + c_{11} \sum (X_{1j}Y_j) + c_{12} \sum (X_{2j}Y_j) + \dots + c_{1p} \sum (X_{pj}Y_j) \\ b_2 &= c_{20} \sum Y_j + c_{21} \sum (X_{1j}Y_j) + c_{22} \sum (X_{2j}Y_j) + \dots + c_{2p} \sum (X_{pj}Y_j) \\ &\dots \\ b_p &= c_{p0} \sum Y_j + c_{p1} \sum (X_{1j}Y_j) + c_{p2} \sum (X_{2j}Y_j) + \dots + c_{pp} \sum (X_{pj}Y_j) \end{aligned}$$

และ $V(\hat{\beta}) = [X'X]^{-1} \sigma^2$ เมื่อ $V(\hat{\beta})$ เป็น Variance Covariance Matrix ของ $\hat{\beta}$
นั่นคือ $s_{b_i}^2 = c_{ii} s^2$ เมื่อ s^2 คือค่าประมาณของ σ^2 , $V(Y) = I \sigma^2$

การวิเคราะห์ข้อมูลโดยใช้ทฤษฎีความถดถอยเชิงซ้อนหนึ่งเป็นเส้นตรง

ในการทดสอบสมมติฐานเพื่อหาความสัมพันธ์ของข้อมูล ๒ พวก คือระหว่างตัวแปรตามกับตัวแปรอิสระ โดยใช้ทฤษฎีความถดถอยเชิงซ้อนหนึ่งเป็นเส้นตรงนี้ วิธีหนึ่งที่มีประโยชน์ก็คือ การวิเคราะห์ความแปรปรวน (Analysis of Variance) ซึ่งใช้ค่า F เป็นเครื่องพิจารณาตามทฤษฎีนี้ เรากำหนดข้อสมมติว่า ค่าความคลาดเคลื่อน e มีการกระจายแบบโค้งปกติ (normal curve) มีค่าเฉลี่ย (mean) = 0 และความแปรปรวน (variance) = σ^2 นั่นคือ $e \sim N(0, \sigma^2)$

ค่าประมาณที่ไม่มีค่าเอนกของ σ^2 คือ $s^2 = \frac{\sum_{j=1}^n (Y_j - \hat{Y}_j)^2}{[n - (p+1)]}$

เมื่อ n เป็นจำนวนข้อมูล

p + 1 เป็นจำนวน พารามิเตอร์ (parameters) ที่ใช้ในแบบจำลอง

เนื่องจากวิธีที่ใช้ในการหาค่า $b_0, b_1, b_2, \dots, b_p$ ของสมการถดถอยเราใช้วิธีผลบวกของกำลังสองของความคลาดเคลื่อนน้อยที่สุด ดังนั้นในการพิจารณาวิเคราะห์เพื่อทดสอบสมมติฐานความสัมพันธ์ระหว่างตัวแปรตามและตัวแปรอิสระ เราจึงต้องพิจารณาจากค่า ผลบวกของกำลังสองของความคลาดเคลื่อนดังนี้

$$\sum_{j=1}^n e_j^2 = \sum_{j=1}^n (Y_j - \hat{Y}_j)^2 = \sum_{j=1}^n (Y_j - \bar{Y})^2 + \sum_{j=1}^n (\hat{Y}_j - \bar{Y})^2 - 2 \sum_{j=1}^n (Y_j - \bar{Y})(\hat{Y}_j - \bar{Y}) \quad (12)$$

เพื่อให้เห็นได้ชัดจะพิจารณากรณีที่มีตัวแปรอิสระ X_I เพียงตัวเดียวในแบบจำลอง

จาก (12) จะเห็นได้ว่า

$$\begin{aligned} \sum_{j=1}^n (Y_j - \bar{Y})(\hat{Y}_j - \bar{Y}) &= \sum_{j=1}^n (Y_j - \bar{Y}) [\bar{Y} + b_1 (X_{Ij} - \bar{X}_I) - \bar{Y}] \\ &= \sum_{j=1}^n (Y_j - \bar{Y}) [b_1 (X_{Ij} - \bar{X}_I)] \\ &= b_1 \sum_{j=1}^n (Y_j - \bar{Y})(X_{Ij} - \bar{X}_I) \\ &= b_1 (b_1 \sum_{j=1}^n (X_{Ij} - \bar{X}_I)^2) \end{aligned} \quad b_1 = \frac{\sum_{j=1}^n (X_{Ij} - \bar{X}_I)(Y_j - \bar{Y})}{\sum_{j=1}^n (X_{Ij} - \bar{X}_I)^2}$$

$$\begin{aligned}
 &= b_1^2 \sum_{j=1}^n (X_j - \bar{X})^2 = \sum_{j=1}^n [b_1 (X_j - \bar{X})]^2 \\
 &= \sum_{j=1}^n (\hat{Y}_j - \bar{Y})^2 \quad (\text{เพราะ } b_1 (X_j - \bar{X}) = \hat{Y}_j - \bar{Y})
 \end{aligned}$$

เมื่อนำไปแทนค่าใน (12) จะได้ดังนี้

$$\begin{aligned}
 \sum_{j=1}^n e_j^2 &= \sum_{j=1}^n (Y_j - \hat{Y}_j)^2 = \sum_{j=1}^n (Y_j - \bar{Y})^2 + \left(\sum_{j=1}^n (\hat{Y}_j - \bar{Y})^2 - 2 \sum_{j=1}^n (Y_j - \bar{Y})(\hat{Y}_j - \bar{Y}) \right) \\
 &= \sum_{j=1}^n (Y_j - \bar{Y})^2 - \sum_{j=1}^n (\hat{Y}_j - \bar{Y})^2 \\
 \therefore \sum_{j=1}^n (Y_j - \bar{Y})^2 &= \sum_{j=1}^n (\hat{Y}_j - \bar{Y})^2 + \sum_{j=1}^n (Y_j - \hat{Y}_j)^2
 \end{aligned}$$

นั่นคือ ผลบวกกำลังสองของความคลาดเคลื่อนเนื่องจากค่าจริงของตัวแปรตาม แตกต่างไปจากค่าเฉลี่ย (Total sum of Squares = $\sum_{j=1}^n (Y_j - \bar{Y})^2$) เท่ากับ ผลบวกของกำลังสองของความคลาดเคลื่อนเนื่องจากค่าประมาณของตัวแปรตามแตกต่างไปจากค่าเฉลี่ย (คือความเบี่ยงเบนคิดจากเส้นถดถอยกับค่าเฉลี่ย = $\sum_{j=1}^n (\hat{Y}_j - \bar{Y})^2$) รวมกับ ผลบวกของกำลังสองของความคลาดเคลื่อนเนื่องจากค่าจริงของตัวแปรตามแตกต่างไปจากค่าประมาณของตัวเอง (sum of squares - about regression หรือ Residual sum of Squares = $\sum_{j=1}^n (Y_j - \hat{Y}_j)^2$)

ค่าต่าง ๆ ดังกล่าวนั้นเราพิจารณาข้อมูลในรูปเมทริกซ์ จะได้เป็นรูปดังนี้

$$\begin{aligned}
 \sum_{j=1}^n (Y_j - \bar{Y})^2 &= (YY)_{1 \times 1} \\
 \sum_{j=1}^n (\hat{Y}_j - \bar{Y})^2 &= (b'XY)_{1 \times 1} \\
 \sum_{j=1}^n (Y_j - \hat{Y}_j)^2 &= YY - b'XY
 \end{aligned}$$

จากค่าผลบวกกำลังสองทั้งสามค่านี้ เรานำไปวิเคราะห์ความแปรปรวนเพื่อทดสอบสมมติฐานของความสัมพันธ์ต่อไป

เนื่องจากค่าเฉลี่ยที่เรานำมาใช้มีค่าเฉลี่ยที่แท้จริงของข้อมูลในประชากร แต่เป็นค่าที่เราประมาณออกมาจากข้อมูลที่เป็นตัวอย่งในการพิจารณา เพื่อมิให้เกิดความลำเอียงในการคำนวณค่าต่าง ๆ จึงต้องมีการแก้ไข ดังแสดงไว้ในตารางวิเคราะห์ความแปรปรวนต่อไปนี้



ตารางวิเคราะห์ความแปรปรวน (ANOVA)

003537

แหล่งความแปรปรวน S.v.	D.F.	SUM SQUARES	MEAN SQUARES	F	R^2 R %
Due to Regression	p	$\sum \hat{Y}^2 - n\bar{Y}^2 = SSR$	$\frac{SSR}{p} = A$	A/B	$\frac{SSR}{SST} \times 100$
Residual(error)	n-p-1	$\sum Y^2 - \sum \hat{Y}^2 = SSE$	$\frac{SSE}{n-p-1} = B$		
Total(corrected)	n-1	$\sum Y^2 - n\bar{Y}^2$			

จะเห็นได้ว่าในการวิเคราะห์ความแปรปรวน เราอาจพิจารณาจากค่า F ประกอบกับค่าสัมประสิทธิ์แห่งการตกลงใจ (coefficient of determination) สำหรับค่า F เป็นค่าที่ใช้ในการทดสอบสมมติฐาน (null hypothesis) $H_0 : \beta_i = 0$ หรือ $H_0 : \beta_0 = \beta_1 = \beta_2 = 0$ การหาค่า F พิจารณาจากอัตราส่วนของ ผลบวกกำลังสองของความแตกต่างระหว่างค่าจริงของตัวแปรตามที่แตกต่างกันไปจากค่าเฉลี่ยของมัน กับ ผลบวกกำลังสองของความแตกต่างระหว่างค่าจริงกับค่าประมาณของตัวแปรตาม โดยที่ได้นำค่า degrees of freedom ของค่า นั้น ๆ ไปหารแล้ว สำหรับค่า degrees of freedom ของค่า F กำหนดด้วยค่า degrees of freedom ของค่าผลบวกกำลังสองของความแตกต่างทั้งสองค่าที่นำมาใช้คำนวณหาค่า F นั้น

$$F = \frac{\text{Mean Squares due to Regression (D.F. = p)}}{\text{Residual Mean Squares (D.F. = n - p - 1)}}$$

ค่า degrees of freedom ของ F คือ (p , n - p - 1)

จากค่า F ที่คำนวณได้ เรานำไปเปรียบเทียบกับค่า F จากตาราง ซึ่งมี degrees of freedom เดียวกัน ณ ระดับนัยสำคัญทางสถิติ ที่ต้องการพิจารณา เช่น 0.01 หรือ 0.05 ซึ่งไร้สัญลักษณ์ α สำหรับในการวิเคราะห์ของปัญหานี้ จะใช้ระดับนัยสำคัญ $\alpha = 0.01$

การเปรียบเทียบ ถ้าค่า F ที่คำนวณได้ มากกว่า ค่า $F_{(p, n-p-1)}$ จากตาราง เราจะต้อง ปฏิเสธสมมติฐาน null Hypothesis

ถ้าค่า $F_{(p, n-p-1)}$ ที่คำนวณได้ น้อยกว่า ค่า $F_{(p, n-p-1)}$ จากตาราง เราจะ ไม่ปฏิเสธสมมติฐาน null Hypothesis

ค่า Residual Mean Square ในตารางวิเคราะห์ความแปรปรวน (ANOVA) นั้น
ก็คือค่า กำลังสองของค่าประมาณของความคลาดเคลื่อนมาตรฐาน (standard error) ที่
เกิดขึ้นในการประมาณค่า ตัวแปรตาม Y ด้วยค่าตัวแปรอิสระ X_i ต่าง ๆ ในแบบจำลองนั้น

จาก
$$s_{Y.1 \dots p} = \sqrt{\frac{\sum_{j=1}^n (Y_j - \hat{Y}_j)^2}{n - p - 1}}$$

เมื่อ $s_{Y.1 \dots p}$ คือค่าประมาณของความคลาดเคลื่อนมาตรฐานในการประมาณค่าตัวแปร
ตาม Y ดังนั้นถ้าพิจารณาจากตารางวิเคราะห์ความแปรปรวนจะได้

$$s_{Y.1 \dots p} = \sqrt{\text{residual mean square}}$$

สำหรับค่าประมาณของความคลาดเคลื่อนมาตรฐานในการประมาณค่า b_1, b_2, \dots, b_p
พิจารณาจาก Variance Covariance Matrix ของ \hat{b} , $V(\hat{b}) = V = [XX']^{-1} \sigma^2 = C \sigma^2$

$$C = \begin{bmatrix} c_{00} & c_{01} & c_{02} & \dots & c_{0p} \\ c_{10} & c_{11} & c_{12} & \dots & c_{1p} \\ c_{20} & c_{21} & c_{22} & \dots & c_{2p} \\ \dots & \dots & \dots & \dots & \dots \\ c_{p0} & c_{p1} & c_{p2} & \dots & c_{pp} \end{bmatrix} (P+1)(P+1)$$

$$s_{b_i} = \sqrt{c_{ii} \sigma^2} = s_{Y.1 \dots p} \sqrt{c_{ii}}$$

เมื่อ s_{b_i} เป็นค่าประมาณของความคลาดเคลื่อนมาตรฐานในการประมาณค่า b_i
 c_{ii} เป็นค่าบนแนวทแยงมุม (diagonal) ของ แมทริกซ์ C ตัวที่ ii
($i = 0, 1, 2, \dots, p$)

σ^2 เป็นความแปรปรวนของข้อมูลที่นำมาใช้เป็นตัวแปรตาม

นอกจากนี้ยังอาจหาค่าความแปรปรวนรวม (covariance) ระหว่าง b_i กับ b_k

ได้ดังนี้

$$COV(b_i, b_k) = \frac{c_{ik} \sigma^2}{\sqrt{c_{ii} \sigma^2} \sqrt{c_{kk} \sigma^2}}$$

c_{ik} เป็นค่าในแมทริกซ์ C ตัวที่ ik

($k = 0, 1, 2, \dots, p$)

ค่าสหสัมพันธ์เชิงเดียว (Simple Correlation)

ในกรณีที่เราต้องการทราบความสัมพันธ์ระหว่างข้อมูล ๒ พวก ว่ามีความสัมพันธ์กันมากน้อย เป็นไปในทางเดียวกันหรือกลับกันอย่างไร เราสามารถพิจารณาได้จากค่า สหสัมพันธ์ ระหว่างข้อมูลทั้ง ๒ ค่าสหสัมพันธ์ระหว่างข้อมูลที่ตีตัวแปรข้าม และตัวแปรอิสระ เพียงตัวเดียว เราเรียกว่าค่า สหสัมพันธ์เชิงเดียว หรือ สหสัมพันธ์อย่างง่าย ซึ่งมีสัญลักษณ์ r ถ้าต้องการหาค่า สหสัมพันธ์เชิงเดียวระหว่างข้อมูล X และ Y เราจะใช้สัญลักษณ์ r_{XY} ซึ่งหาได้ดังนี้

$$r = \frac{\text{COV}(X, Y)}{\sqrt{V(X) V(Y)}}$$

$$= \frac{\sum_{j=1}^n (X_j - \bar{X})(Y_j - \bar{Y})}{\sqrt{\sum_{j=1}^n (X_j - \bar{X})^2 \sum_{j=1}^n (Y_j - \bar{Y})^2}}$$

ค่า สหสัมพันธ์ จะมีค่าอยู่ระหว่าง -1 ถึง $+1$

ถาค่าสหสัมพันธ์เข้าใกล้ -1 แสดงว่า Y มีความสัมพันธ์กับ X มาก แต่เป็นไปในทางกลับกัน คือถ้า X มีค่าเพิ่มขึ้น Y จะมีค่าลดลง

ถาค่าสหสัมพันธ์เข้าใกล้ $+1$ แสดงว่า Y มีความสัมพันธ์กับ X มาก และเป็นไปในทางเดียวกัน คือถ้า X มีค่าเพิ่มขึ้น Y จะมีค่าเพิ่มขึ้นด้วย

กล่าวโดยสรุปได้ว่า ถาค่าสัมบูรณ์ของ r , $|r|$ มีค่าใกล้ 1 แสดงว่า Y กับ X มีความสัมพันธ์กันสูง

ถาค่าสัมบูรณ์ของ r , $|r|$ มีค่าใกล้ 0 แสดงว่า Y กับ X มีความสัมพันธ์กันน้อยมาก

ค่าสหสัมพันธ์เชิงซ้อน (Multiple Correlation)

เป็นค่าสหสัมพันธ์ระหว่างตัวแปรตาม Y กับตัวแปรอิสระ X_1, X_2, \dots, X_p ทั้งหมดที่นำมาใช้ในแบบจำลองความถดถอย ใช้สัญลักษณ์ R เป็นค่าที่บอกให้ทราบว่า ตัวแปรตาม Y มีความสัมพันธ์ กับ ตัวแปรอิสระทุกตัวที่นำมาใช้ในแบบจำลองความถดถอย มากหรือน้อยเพียงใด (ไม่ใช่ค่าสหสัมพันธ์ เฉพาะกับตัวแปรอิสระตัวใดตัวหนึ่ง แต่เป็นค่าสหสัมพันธ์กับตัวแปรอิสระทั้งหมดทุกตัวในแบบจำลอง) หรือบอกให้ทราบถึงความสัมพันธ์ ระหว่าง ค่าตัวแปรตามจริง กับ ค่าประมาณของตัวแปรตาม ซึ่งคำนวณได้จากแบบจำลองความถดถอยนั้น จึงอาจเขียนได้ดังนี้

$$R = r_{\hat{Y}Y}$$

$$= \frac{\text{COV}(Y, \hat{Y})}{\sqrt{V(Y) V(\hat{Y})}}$$

ค่า R มีค่าอยู่ระหว่าง -1 ถึง +1 ความหมายของค่า R ที่แสดงถึงความสัมพันธ์ของตัวแปรตามกับตัวแปรอิสระ พิจารณาได้เช่นเดียวกับค่าสหสัมพันธ์เชิงเดียว

ค่า R² เราเรียกว่า สัมประสิทธิ์แห่งการตกลงใจ (coefficient of determination) ซึ่งจะบอกให้ทราบว่า สัดส่วนของความแปรปรวนของตัวแปรตาม Y ที่ขึ้นกับความแปรปรวนของค่าตัวแปรอิสระ X₁ ต่าง ๆ ในแบบจำลอง มีมากน้อยอย่างไร เมื่อเทียบกับค่าความแปรปรวนทั้งหมดของตัวแปรตาม Y นั้น

พิจารณาจากตารางวิเคราะห์ความแปรปรวน

$$R^2 = \frac{\text{regression sum of squares}}{\text{total sum of squares}}$$

$$= \frac{\sum_{j=1}^n (\hat{Y}_j - \bar{Y})^2}{\sum_{j=1}^n (Y_j - \bar{Y})^2}$$

หรือ

$$= \frac{b'XY - n\bar{Y}^2}{n\bar{Y}^2 - n\bar{Y}^2}$$

(ในรูปแมทริกซ์)

เพราะว่า

$$\sum_{j=1}^n (Y_j - \hat{Y}_j)^2 = \sum_{j=1}^n Y_j^2 - R^2 \sum_{j=1}^n Y_j^2$$

$$\therefore \sum_{j=1}^n Y_j^2 = \sum_{j=1}^n (Y_j - \hat{Y}_j)^2 + R^2 \sum_{j=1}^n Y_j^2$$

$$\sum_{j=1}^n Y_j^2 = \sum_{j=1}^n Y_j^2 (1 - R^2) + R^2 \sum_{j=1}^n Y_j^2$$

ซึ่งจะได้

$$1 = (1 - R^2) + R^2$$

ค่า (1 - R²) เป็นสัดส่วนของความแปรปรวนของค่าตัวแปรตาม Y ซึ่งไม่ขึ้นกับความแปรปรวนของตัวแปรอิสระทุกตัวในแบบจำลอง (coefficient of nondetermination)

ค่า R² คือค่าสัมประสิทธิ์ของการตกลงใจ (coefficient of determination) แสดงสัดส่วนของ Y² ที่มีความแปรปรวนขึ้นกับความแปรปรวนของตัวแปรอิสระทุกตัวในแบบจำลอง

ค่า R² มีค่าอยู่ระหว่าง ๐ ถึง ๑

ถ้าค่า R^2 มีค่าเข้าใกล้ ๑ แสดงว่า ความแปรปรวนของค่าตัวแปรตาม Y ที่เกิดขึ้น เนื่องจากความแปรปรวนของค่า X_i ต่าง ๆ ในแบบจำลอง มีมาก หรือก็คือ แบบจำลองความถดถอยที่ใช้ในการประมาณค่าตัวแปรตาม Y นั้น ให้ความประมาณของ Y ได้ใกล้เคียงกับค่าจริงมาก ซึ่งแสดงว่าแบบจำลองที่ใช้ในการประมาณค่าตัวแปรตาม Y นั้น มีความเหมาะสมคตินั่นเอง

ถ้าค่า R^2 มีค่าน้อยมากจนเข้าใกล้ ๐ ก็แสดงว่า สัดส่วนความแปรปรวนของค่าตัวแปรตาม Y ที่ขึ้นอยู่กับความแปรปรวนของค่าตัวแปรอิสระ X_i ต่าง ๆ ในแบบจำลองความถดถอยที่ใช้ มีน้อยมาก ซึ่งยอมแสดงว่า แบบจำลองความถดถอยที่ใช้ในการประมาณค่าตัวแปรตาม Y นั้น ยังไม่เหมาะสมที่จะนำมาใช้ จะต้องแก้ไข เปลี่ยนแปลงใหม่

ค่าสหสัมพันธ์พหุเชิงเส้น (Partial Correlation) หรือสหสัมพันธ์เชิงส่วน

ในแบบจำลองความถดถอยที่มีตัวแปรหลายตัว นอกจากเราจะหาค่าสหสัมพันธ์เชิงเดี่ยวระหว่างตัวแปรต่าง ๆ เป็นคู่ ๆ ได้ เช่น r_{ik} ซึ่งก็คือค่าสหสัมพันธ์เชิงเดี่ยว ระหว่างตัวแปรที่ i กับตัวแปรที่ k โดยอาศัยทฤษฎีที่ได้อธิบายมาข้างต้นแล้ว ในกรณีที่ตัวแปรต่าง ๆ ทุกตัว ผันแปรพร้อมกัน แบบ การกระจายปกติ (normal distribution) เราก็อาจจะหาค่าสหสัมพันธ์เชิงส่วนของตัวแปรใด ลองพิจารณาในแบบจำลองความถดถอยที่มีตัวแปรตาม ๑ ตัว และตัวแปรอิสระ ๓ ตัว ซึ่งมีแบบจำลองความถดถอยดังนี้

$$\hat{Y} = b_0 + b_1X_1 + b_2X_2 + b_3X_3$$

นอกจากเราจะได้อ่านค่าสหสัมพันธ์เชิงเดี่ยว เช่น r_{YX_1} , r_{YX_2} , $r_{X_1X_2}$ แล้ว เรายังสามารถหาค่า สหสัมพันธ์เชิงส่วนได้ เช่น

- $r_{Y1.2}$ คือค่าสหสัมพันธ์เชิงส่วนระหว่างตัวแปร Y กับตัวแปร X_1 โดยถือว่า X_2 คงที่
- $r_{Y2.1}$ คือค่าสหสัมพันธ์เชิงส่วนระหว่างตัวแปร Y กับตัวแปร X_2 โดยถือว่า X_1 คงที่
- $r_{12.3}$ คือค่าสหสัมพันธ์เชิงส่วนระหว่างตัวแปร X_1 กับตัวแปร X_2 โดยถือว่า X_3 คงที่
- $r_{13.2}$ คือค่าสหสัมพันธ์เชิงส่วนระหว่างตัวแปร X_1 กับตัวแปร X_3 โดยถือว่า X_2 คงที่

.....
ค่าสหสัมพันธ์เชิงส่วนต่าง ๆ นั้น จะหาได้โดยอาศัยสูตรต่อไปนี้

$$r_{Y3.12} = \frac{r_{Y3.1} - r_{Y2.1}r_{32.1}}{[(1-r_{Y2.1}^2)(1-r_{32.1}^2)]^{1/2}} = \frac{r_{Y3.2} - r_{Y1.2}r_{31.2}}{[(1-r_{Y1.2}^2)(1-r_{31.2}^2)]^{1/2}}$$

$$r_{Y1.2} = \frac{r_{Y1} - r_{Y2}r_{12}}{[(1-r_{Y2}^2)(1-r_{12}^2)]^{1/2}}$$

$$r_{I2.3} = \frac{r_{I2} - r_{I3}r_{23}}{[(1-r_{I3}^2)(1-r_{23}^2)]^{1/2}}$$

สำหรับค่าอื่น ๆ ก็คงพิจารณาจากรูปสูตรเดียวกัน

ค่าสหสัมพันธ์เชิงส่วน มีค่าอยู่ระหว่าง -1 ถึง +1 และการพิจารณาความหมายของความสัมพันธ์ของตัวแปร จากค่าสหสัมพันธ์เชิงส่วนที่ได้นี้ ก็คงอาศัยหลักการเกี่ยวกับการพิจารณาความหมายของค่าสหสัมพันธ์เชิงเดี่ยว ดังได้กล่าวมาแล้ว

วิธีต่าง ๆ ในการสร้างสมการถดถอย

ในกรณีที่มีตัวแปรตาม และตัวแปรอิสระเพียงตัวเดียว การสร้างแบบจำลองความถดถอยย่อมไม่มีปัญหายุ่งยาก แต่ในการสร้างแบบจำลองความถดถอยในกรณีที่มีตัวแปรอิสระ หลาย ๆ ตัวนั้น ย่อมมีปัญหาก่เกิดขึ้นว่า เราจะมึวิธีการอย่างไรที่จะสร้างแบบจำลองความถดถอยให้โดนดในทางประมาณค่าตัวแปรตามได้ดีที่สุด และเป็นวิธีการที่ประหยัด เวลา และ แรงงาน ได้มากที่สุด

วิธีต่าง ๆ ที่นิยมใช้กันในการสร้างสมการถดถอย สำหรับกรณีที่มีตัวแปรอิสระ หลายตัวนั้นมีดังนี้

๑. วิธีใช้ทดสอบแบบจำลองทั้งหมดที่จะสามารถสร้างขึ้นได้ เพื่อหาแบบจำลองที่ดีที่สุด

(All Possible Regression)

๒. วิธีใช้ตัวแปรอิสระทุกตัวสร้างแบบจำลองความถดถอยขึ้นมาก่อน แล้วจึงทดสอบนัยสำคัญ เพื่อจะตัดตัวแปรอิสระออกไปทีละตัว (ตัดเฉพาะตัวที่ทดสอบแล้วไม่มีนัยสำคัญ)

(Backward Selection)

๓. วิธีทดลองเลือกตัวแปรอิสระที่มีค่าสหสัมพันธ์กับตัวแปรตามสูงสุดสร้างแบบจำลองความถดถอยอย่างง่ายขึ้นก่อน แล้วค่อยเพิ่มตัวแปรอิสระตัวที่มีความสำคัญรองลงไปทีละตัว ทดสอบนัยสำคัญ (Forward Selection) ซึ่งจะนำมาใช้ในกรณี (อธิบายหน้า ๒๑)

๔. วิธีทดสอบทีละตัว ลองเพิ่มเขาไปและสลับตำแหน่ง (Stepwise Regression)

จากวิธีการต่าง ๆ ดังกล่าวผู้วิจัยพิจารณาเห็นว่า วิธีเลือกตัวแปรอิสระเพิ่มเข้าไปที่ละตัว มีความเหมาะสมกับปัญหามากกว่าวิธีอื่น กล่าวคือวิธีการแบบจำลองทั้งหมดที่จะเป็นไปได้นั้น ในขั้นแรกก็ต้องพิจารณาแบบจำลองที่มีตัวแปรอิสระเพียงตัวเดียว ว่าแบบจำลองที่ใช้ตัวแปรอิสระใด จะให้ผลดีที่สุด ซึ่งเมื่อพิจารณาจะเห็นว่าก็คือวิธีการขั้นแรกของวิธีที่จะนำมาใช้นั่นเอง จากที่ได้ ทดลองแบบจำลองที่มีตัวแปรอิสระเพียงตัวเดียวแล้ว ขั้นตอนต่อไปก็ต้องพิจารณาแบบจำลองทุกแบบที่จะ มีไว้ว่าแบบใดที่ดีที่สุด ผู้วิจัยเห็นว่าในขั้นนี้ผลที่ได้ก็คงจะคล้ายคลึงกับวิธีเลือกทีละตัวเพราะเรา ก็เลือกตัวที่มีความสัมพันธ์กับตัวแปรตามสูงที่สุดอยู่แล้ว ดังนั้นวิธีแรกนี้จึงออกจะยุ่งยาก และผล ก็จะไม่แตกต่างกับผลที่ได้จากวิธีเลือกทีละตัวโดยเฉพาะอย่างยิ่งในรูปแบบที่ว่าตัวแปรอิสระตัวใด จะเรียงลำดับก่อนหลังกันอย่างไรในแบบจำลองความถดถอยที่จะนำมาใช้ขั้นนี้ก็คงจะไม่ต่างกัน

ในวิธีสร้างแบบจำลองที่มีทุกตัวแปรอิสระขึ้นก่อนแล้วตัดออกทีละตัวนั้นนับว่าเป็นวิธีที่ไม่เหมาะสม อย่างยิ่ง เพราะในการสร้างแบบจำลองรวมนั้นเราจะไม่ได้พิจารณาถึงความสำคัญของตัวแปรอิสระ แต่ละตัว ว่าควรจะเรียงลำดับกันอย่างไร แบบจำลองที่ได้จากวิธีนี้จึงไม่ใช่แบบจำลองที่เหมาะสม ที่สุด วิธีนี้น่าจะใช้ในกรณีที่มีตัวแปรอิสระมาก ๆ

สำหรับวิธีเลือกเข้าไปที่ละตัวและทดลองสลับตำแหน่งนั้น จะเห็นได้ว่าตัวแปรอิสระของ เรามีน้อยมาก ในขั้นแรกเราก็ได้เลือกตัวแปรอิสระที่มีความสัมพันธ์กับตัวแปรตามสูงที่สุดมาใส่ก่อน ต่อไปที่พิจารณาจากตัวแปรอิสระที่เหลืออยู่ที่มีค่าสหสัมพันธ์เชิงสวนกับตัวแปรตามสูงที่สุดมาใส่ต่อไป และจากที่ได้พิจารณาเห็นว่า ตัวแปรอิสระที่สัมพันธ์กับตัวแปรตามอยู่ในเกณฑ์สูงนั้น มีเพียงสองตัว ซึ่งก็ให้ทดลองสร้างแบบจำลองของทั้งสองตัวแปรนี้เปรียบเทียบกันก่อนจะนำไปใช้เป็นแบบจำลอง จริงแล้ว โดยเลือกตัวแปรที่ให้ผลดีที่สุด ดังนั้นปัญหาในการที่จะสลับที่ของตัวแปรทั้งสองน่าจะไม่มี และสำหรับตัวแปรที่ ๓ และที่ ๔ ก็มีความสัมพันธ์กับตัวแปรตามแตกต่างกันมาก โดยเฉพาะตัว แปรอิสระตัวที่ ๔ นั้น มีความสัมพันธ์กับตัวแปรตามอยู่ในเกณฑ์ต่ำมาก และเมื่อนำไปสร้างแบบจำลอง ความถดถอยทำการทดสอบแล้วก็ไม่มีความสำคัญทางสถิติ ดังนั้นปัญหาในการสลับตำแหน่งจึงหมดไป

จากการพิจารณาถึงผลของการเลือกใช้วิธีต่าง ๆ ในการสร้างสมการถดถอยแล้ว ผู้วิจัย จึงได้ตกลงใจใช้วิธีเลือกตัวแปรเพิ่มเข้าไปที่ละตัว เป็นวิธีวิเคราะห์ความสัมพันธ์ในปัญหานี้ ดัง จะได้อธิบายวิธีที่จะนำมาใช้ในหัวข้อต่อไป

วิธี Forward Selection Procedure

เป็นวิธีในการสร้างสมการถดถอยสำหรับปัญหาที่มีตัวแปรอิสระหลาย ๆ ตัว โดยเลือกสร้างสมการถดถอยเชิงเดียว (simple regression) ระหว่างตัวแปรตาม กับ ตัวแปรอิสระตัวที่มีความสัมพันธ์เชิงเดียวกับตัวแปรตามนั้นสูงสุด หรือก็คือเลือกตัวแปรอิสระที่เห็นว่ามีความสัมพันธ์กับตัวแปรตามมากที่สุดนั่นเองมาพิจารณาก่อน แล้วทำการทดสอบแบบจำลองที่ได้ ถ้าพบว่ามีความสำคัญก็ดำเนินการต่อไป โดยพิจารณาเลือกตัวแปรอิสระที่เหลืออยู่ ซึ่งมีความสัมพันธ์เชิงส่วนกับตัวแปรตามนั้นสูงสุด (โดยไม่พิจารณาถึงตัวแปรอิสระตัวที่ได้ใช้ในแบบจำลองไปแล้ว) เพิ่มเข้าไปเป็นตัวแปรอิสระอีกตัวหนึ่ง และสร้างแบบจำลองความถดถอยใหม่ขึ้น ทำการทดสอบความสำคัญต่อไปว่าแบบจำลองความถดถอยใหม่ที่สร้างขึ้นโดยเพิ่มตัวแปรอิสระตัวใหม่เข้าไปนี้ จะมีความสำคัญเป็นนัยสำคัญในการประมาณค่าตัวแปรตาม Y หรือไม่ และแบบจำลองใหม่นี้ใหม่ลึกลับกว่าแบบจำลองเดิมหรือไม่ จากนั้นก็พิจารณาเลือกตัวแปรอิสระตัวต่อไป สร้างแบบจำลองความถดถอยขึ้นใหม่ และทดสอบ ตามวิธีที่กล่าวมาแล้ว ตัวแปรอิสระตัวใดที่ทดสอบในแบบจำลองแล้วไม่มีความสำคัญก็จะตัดทิ้งออกไปไม่นำมาใส่ในแบบจำลอง

วิธีนี้จะจบลงเมื่อพบว่า ตัวแปรอิสระที่เหลืออยู่เมื่อนำไปเพิ่มสร้างแบบจำลองใหม่สำหรับใช้ในการประมาณค่าตัวแปรตาม Y แล้ว แบบจำลองที่ได้ใหม่ไม่มีความสำคัญ หรือก็คือตัวแปรอิสระตัวที่เพิ่มเข้าไปใหม่นั้น ไม่มีความสัมพันธ์เป็นนัยสำคัญกับตัวแปรตาม Y หรืออาจจะพิจารณาจากค่า สัมประสิทธิ์แห่งการตกลงใจ R^2 ของแบบจำลองที่ทดสอบได้มีความสำคัญนั้นว่ามีค่าสูงพอหรือไม่

ศูนย์จตุรพิทยาการ
จุฬาลงกรณ์มหาวิทยาลัย