

บทที่ 2

สถิติที่ใช้ในการวิจัย

แผนแบบการทดลองแบบแฟคทอเรียล (Factorial Experiments)

แผนแบบการทดลองแบบแฟคทอเรียลเป็นแผนแบบการทดลองที่น่าสนใจมาก เพราะเป็นแผนแบบการทดลองที่สามารถทดลองปัจจัย (factor) ได้ทีละหลาย ๆ ปัจจัยไปพร้อม ๆ กันได้ ซึ่งทำให้สามารถประหยัดเวลา และค่าใช้จ่ายได้เป็นอย่างมาก และแผนแบบการทดลองแบบแฟคทอเรียลที่สำคัญแบบหนึ่งคือแผนแบบการทดลองแบบแฟคทอเรียลขนาด 2^2 หรือแผนแบบการทดลองที่มีปัจจัย 2 ปัจจัย ซึ่งแต่ละปัจจัยมี 2 ระดับ หรืออาจเรียกว่าระดับต่ำ (low) และระดับสูง (high) ซึ่งมีรูปแบบผลบวกของสิ่งทดลอง (treatment combinations) ดังนี้

อิทธิพลแฟคทอเรียล (Factorial effect)	Treatment combinations			
	(1)	a	b	ab
A	-	+	-	+
B	-	-	+	+
AB	+	-	-	+

ซึ่งการหาค่าเฉลี่ยของอิทธิพลของปัจจัย (the average of factor) สามารถหาได้ดังนี้

$$A = 1/2r [ab + a - b - (1)]$$

$$B = 1/2r [ab + b - a - (1)]$$

$$AB = 1/2r [ab + (1) - a - b]$$

และการหาผลรวมกำลังสอง (Sum of Squares) สามารถหาได้ดังนี้

$$SSA = [ab + a - b - (1)]^2 / (r \cdot 4)$$

$$SSB = [ab + b - a - (1)]^2 / (r \cdot 4)$$

$$SSAB = [ab + (1) - a - b]^2 / (r \cdot 4)$$

$$SST = \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^n y_{ijk}^2 - \frac{y_{...}^2}{4r}$$

$$SSE = SST - SSA - SSB - SSAB$$

โดยองศาความเป็นอิสระของ SST เท่ากับ $(4r - 1)$

ของ SSA เท่ากับ 1

ของ SSB เท่ากับ 1

ของ SSAB เท่ากับ 1

ของ SSE เท่ากับ $4 \cdot (r - 1)$

เมื่อ r คือจำนวนการทำซ้ำ

ซึ่งสามารถเขียนตารางวิเคราะห์ความแปรปรวนได้ดังนี้

สาเหตุของความผันแปร	d.f.	ผลรวมกำลังสอง	ผลรวมกำลังสองเฉลี่ย	F
A	1	SSA	SSA/1	MSA/MSE
B	1	SSB	SSB/1	MSB/MSE
AB	1	SSAB	SSAB/1	MSAB/MSE
ความคลาดเคลื่อน	$4 \cdot (r - 1)$	SSE	$SSE / (4 \cdot (r - 1))$	
รวม	$(4r - 1)$	SST		

เพื่อให้สามารถหาความคลาดเคลื่อนกำลังสองเฉลี่ย (MSE) ได้ จึงต้องมีจำนวนการทำซ้ำอย่างน้อยที่สุดเท่ากับ 2 แต่เนื่องจากความจำเป็นบางประการอาจทำให้ไม่สามารถหาข้อมูลที่มีความคล้ายคลึงกันจำนวนมาก ๆ ได้ ซึ่งเป็นเหตุให้ไม่สามารถทำซ้ำในแผนแบบการทดลองได้ จึงได้มีการเสนอให้มีการทำซ้ำที่จุดศูนย์กลาง (center point) โดยการทำซ้ำที่จุดศูนย์กลางนี้คือการเก็บค่าสังเกตที่ตัวแปรในแผนแบบการทดลอง (x's) ทุกตัวมีค่าเป็น 0 ซึ่งการทำซ้ำที่จุดศูนย์กลางนี้ทำให้สามารถหาค่าความคลาดเคลื่อนกำลังสองเฉลี่ยได้ และจะไม่ทำให้ค่าสัมประสิทธิ์การถดถอยเปลี่ยนแปลง แต่จะส่งผลทำให้การประมาณค่า β_0 มีค่าเข้าใกล้ค่าเฉลี่ย (Grand Mean) มากยิ่งขึ้น ซึ่งการทำซ้ำที่จุดศูนย์กลางนี้ควรทำอย่างน้อยที่สุดเท่ากับ

$(2k - 1)$ และต้องมีจำนวนไม่มากเกินไป เนื่องจากการทำซ้ำที่จุดศูนย์กลางนี้ทำขึ้นเพื่อแก้ ปัญหาการไม่สามารถหาข้อมูลที่มีความคล้ายคลึงกันได้ โดยการหาค่าความคลาดเคลื่อนกำลัง สองเฉลี่ยจากการทำซ้ำที่จุดศูนย์กลางสามารถหาได้จากสูตรดังนี้

$$\begin{aligned} \text{MSE} &= \text{SSE}/(n_0 - 1) \\ &= \frac{\sum_{i=1}^{n_0} (y_i - \bar{y})^2}{n_0 - 1} \end{aligned}$$

เมื่อ n_0 คือจำนวนการทำซ้ำที่จุดศูนย์กลาง และการทำซ้ำที่จุดศูนย์กลางนี้ยังมีประโยชน์ในการตรวจสอบความโค้ง (curvature) ในรูปแบบ ได้ด้วย โดยจะพิจารณารูปแบบ

$$y = \beta_0 + \sum_{j=1}^p \beta_j x_j + \sum_{j=1}^p \beta_{jj} x_j^2 + \sum_{i < j} \sum_{i=1}^p \beta_{ij} x_i x_j + \varepsilon$$

ทดสอบสมมติฐาน

$$H_0 : \sum_{j=1}^p \beta_{jj} = 0$$

$$H_a : \sum_{j=1}^p \beta_{jj} \neq 0$$

ซึ่งการคำนวณผลรวมกำลังสองของความโค้ง (Sum of Squares of Curvature) สามารถหาได้จากสูตรดังนี้

$$\text{SSC} = \frac{n_f n_0 (\bar{y}_f - \bar{y}_0)^2}{n_f + n_0}$$

เมื่อ n_f คือจำนวนตัวอย่างในแผนแบบการทดลองแบบแฟคทอเรียล

n_0 คือจำนวนการทำซ้ำที่จุดศูนย์กลาง

\bar{y}_f คือค่าเฉลี่ยของ y ที่จุดในแผนแบบการทดลอง

\bar{y}_o คือค่าเฉลี่ยของ y ที่จุดศูนย์กลาง

และ SSC มีองศาความเป็นอิสระเท่ากับ 1

ซึ่งสามารถเขียนตารางวิเคราะห์ความแปรปรวนได้ดังนี้

สาเหตุของความผันแปร	d.f.	ผลรวมกำลังสอง	ผลรวมกำลังสองเฉลี่ย	F
A	1	SSA	SSA/1	MSA/MSE
B	1	SSB	SSB/1	MSB/MSE
AB	1	SSAB	SSAB/1	MSAB/MSE
Curvature	1	SSC	SSC/1	MSC/MSE
ความคลาดเคลื่อน	(n_0-1)	SSE	SSE/ (n_0-1)	
รวม	$n-1$	SST		

เราจะยอมรับสมมติฐานหลัก ก็ต่อเมื่อ F_{cur} มีค่าน้อยกว่า $F_{\alpha,1,(n_0-1)}$ ซึ่งแสดงว่าในรูปแบบมีส่วนโค้งรวมอยู่ด้วย ซึ่งในการประมาณค่าเราจะต้องนำส่วนโค้งนั้นมาพิจารณาด้วย

การวิเคราะห์พื้นผิว (Response Surface Methodology (RSM))

การวิเคราะห์พื้นผิว (Response Surface Methodology (RSM)) เป็นวิธีการวิเคราะห์ที่ใช้เทคนิคทางสถิติ โดยรูปแบบที่ได้จากวิธีนี้จะเป็นการผสมผสานระหว่างรูปแบบที่ได้จากการสังเกต (Empirical model) และรูปแบบที่จะเป็นประโยชน์ (Exploitation model) โดยมีหลักการคือจะพิจารณาบริเวณที่สนใจ (Interest Region) ของตัวแปรพยากรณ์ (predictor variable) โดยจะพยายามหาบริเวณของตัวแปรพยากรณ์ที่เหมาะสมที่สุดที่จะทำให้รูปแบบที่ได้ มีความเอนเอียง (bias) น้อยที่สุด และมีความเหมาะสมของรูปแบบมากที่สุด ซึ่งความเหมาะสมของรูปแบบนี้สามารถพิจารณาได้จากค่าพารามิเตอร์ไร้ศูนย์กลาง (noncentrality parameter)

ในปี ค.ศ. 1963 จอร์จ อี พี บอกซ์ และนอร์แมน อาร์ แดรปเปอร์ (George E.P. Box and Norman R. Draper) ได้พิจารณารูปแบบที่จัดกลุ่มไม่ได้ (model misspecification) โดยทั่วไปมักจะประมาณค่าสัมประสิทธิ์การถดถอยด้วยรูปแบบดังนี้

$$y(\mathbf{x}) = \mathbf{X}_1'\beta_1$$

โดยในที่นี้เราพิจารณารูปแบบกำลังหนึ่ง (First order model) ซึ่งมีรูปแบบดังนี้

$$y(x) = \beta_0 + \sum_{j=1}^p \beta_j x_j + \varepsilon$$

แต่ในความเป็นจริงแล้วรูปแบบความสัมพันธ์ที่แท้จริงอาจจะเป็นรูปแบบอื่น ๆ ที่เราละเลยตัวแปรบางตัวไป ก็จะมีรูปแบบดังนี้

$$\eta(x) = \mathbf{X}_1' \beta_1 + \mathbf{X}_2' \beta_2$$

โดยในที่นี้เราจะพิจารณารูปแบบกำลังสอง (Second order model) ซึ่งมีรูปแบบดังนี้

$$\eta(x) = \beta_0 + \sum_{j=1}^p \beta_j x_j + \sum_{j=1}^p \beta_{jj} x_j^2 + \sum_{i < j} \beta_{ij} x_i x_j + \varepsilon$$

และสามารถเขียนในรูปแบบของข้อมูลจำนวน n ชุดได้ดังนี้

$$y(x) = \mathbf{X}_1 \beta_1 \quad (1)$$

$$\eta(x) = \mathbf{X}_1 \beta_1 + \mathbf{X}_2 \beta_2 \quad (2)$$

ซึ่งถ้ารูปแบบ (2) เป็นจริง แต่เราใช้รูปแบบ (1) ในการประมาณค่า ดังนั้นตัวประมาณที่ได้จาก (1) ซึ่งเท่ากับ $\hat{\beta} = (\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1' y$ จึงเป็นตัวประมาณที่เอนเอียงของ β ซึ่งค่าคาดหวังของ $\hat{\beta}$ คือ

$$E(\hat{\beta}) = \beta + A \beta_2$$

โดยที่ A คือเมทริกซ์ของความเอนเอียง (bias matrix) และ A สามารถเขียนได้ดังนี้

$$A = (\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1' \mathbf{X}_2$$

ซึ่งสามารถคำนวณค่าความเอนเอียงของการประมาณ (bias in prediction) ได้จากสูตรดังนี้

$$B_1 = (nK/\sigma^2) \int_R B^2(y)$$

เมื่อ

$$K^{-1} = \int_R dx$$

R คือบริเวณที่สนใจ

และ B(y) คือความเอนเอียงของ y

โดยเราสามารถหาค่าของ $\int_R B^2(y)$ ได้ในเทอมของโมเมนต์เมทริกซ์ (moment matrices) ดังนี้

$$\int_R B^2(y) = n\beta_2' T \beta_2 / \sigma^2$$

เนื่องจาก

$$B_1 = (n/\sigma^2) \int_R w(x) \{Ey(x) - \eta(x)\}^2 dx$$

เมื่อ

$$y(x) = X_1 \beta_1$$

$$\eta(x) = X_1 \beta_1 + X_2 \beta_2$$

และ

$$M_{11} = n^{-1} X_1' X_1$$

$$M_{12} = n^{-1} X_1' X_2$$

$$M_{22} = n^{-1} X_2' X_2$$

ดังนั้น

$$\begin{aligned} B_1 &= (n/\sigma^2) \int_R \beta_2' (x_1 A - x_1)' (x_1 A - x_2) \beta_2 dx \\ &= (n/\sigma^2) \beta_2' T \beta_2 / \sigma^2 \end{aligned}$$

โดยที่

$$T = \mu_{22} + M_{12}' M_{11}^{-1} \mu_{11} M_{11}^{-1} M_{12} - 2\mu_{12}' M_{11}^{-1} M_{12}$$

n คือขนาดตัวอย่าง

และ μ_{11} , μ_{12} และ μ_{22} คือโมเมนต์เมทริกซ์

เมื่อ

$$\mu_{11} = K \int_R x_1 x_1' dx$$

$$\mu_{12} = K \int_R x_1 x_2' dx$$

และ $\mu_{22} = K \int_R x_2 x_2' dx$

ซึ่งสามารถหาเมทริกซ์ต่างๆ ดังกล่าวที่กล่าวมาข้างต้นได้โดยจะยกตัวอย่างในกรณีแผนแบบ การทดลองแบบแฟกทอเรียลขนาด 2^2 จะได้ว่า

$$\mathbf{X} = [\mathbf{X}_1 \ \mathbf{X}_2]$$

โดยที่

$$\mathbf{X}_1 = [1 \quad x_1 \quad x_2]$$

$$\mathbf{X}_2 = [x_1 x_2 \quad x_1^2 \quad x_2^2]$$

$$x_1 = [1 \quad x_1 \quad x_2]$$

$$x_2 = [x_1 x_2 \quad x_1^2 \quad x_2^2]$$

ดังนั้น

$$\mathbf{X}'_1 \mathbf{X}_1 = \begin{bmatrix} n & x_1 & x_2 \\ x_1 & x_1^2 & x_1 x_2 \\ x_2 & x_1 x_2 & x_2^2 \end{bmatrix} \quad \mathbf{X}'_2 \mathbf{X}_2 = \begin{bmatrix} x_1^2 x_2^2 & x_1^3 x_2^2 & x_1^2 x_2^3 \\ x_1^3 x_2^2 & x_1^4 & x_1^2 x_2^2 \\ x_1^2 x_2^3 & x_1^2 x_2^2 & x_2^4 \end{bmatrix}$$

$$\mathbf{X}'_1 \mathbf{X}_2 = \begin{bmatrix} x_1 x_2 & x_1^2 & x_2^2 \\ x_1^2 x_2 & x_1^3 & x_1 x_2^2 \\ x_1 x_2^2 & x_1^2 x_2 & x_2^3 \end{bmatrix}$$

และ

$$\mathbf{M}_{11} = \frac{1}{n} \begin{bmatrix} n & x_1 & x_2 \\ x_1 & x_1^2 & x_1 x_2 \\ x_2 & x_1 x_2 & x_2^2 \end{bmatrix} \quad \mathbf{M}_{22} = \frac{1}{n} \begin{bmatrix} x_1^2 x_2^2 & x_1^3 x_2^2 & x_1^2 x_2^3 \\ x_1^3 x_2^2 & x_1^4 & x_1^2 x_2^2 \\ x_1^2 x_2^3 & x_1^2 x_2^2 & x_2^4 \end{bmatrix}$$

$$\mathbf{M}_{12} = \frac{1}{n} \begin{bmatrix} x_1 x_2 & x_1^2 & x_2^2 \\ x_1^2 x_2 & x_1^3 & x_1 x_2^2 \\ x_1 x_2^2 & x_1^2 x_2 & x_2^3 \end{bmatrix}$$

และ

$$\mathbf{x}'_1 \mathbf{x}_1 = \begin{bmatrix} 1 & x_1 & x_2 \\ x_1 & x_1^2 & x_1 x_2 \\ x_2 & x_1 x_2 & x_2^2 \end{bmatrix} \quad \mathbf{x}'_2 \mathbf{x}_2 = \begin{bmatrix} x_1^2 x_2^2 & x_1^3 x_2^2 & x_1^2 x_2^3 \\ x_1^3 x_2^2 & x_1^4 & x_1^2 x_2^2 \\ x_1^2 x_2^3 & x_1^2 x_2^2 & x_2^4 \end{bmatrix}$$

$$\mathbf{x}'_1 \mathbf{x}_2 = \begin{bmatrix} x_1 x_2 & x_1^2 & x_2^2 \\ x_1^2 x_2 & x_1^3 & x_1 x_2^2 \\ x_1 x_2^2 & x_1^2 x_2 & x_2^3 \end{bmatrix}$$

และ

$$\mu_{11} = \iint_R \begin{bmatrix} 1 & x_1 & x_2 \\ x_1 & x_1^2 & x_1 x_2 \\ x_2 & x_1 x_2 & x_2^2 \end{bmatrix} dx_1 dx_2$$

$$\mu_{22} = \iint_R \begin{bmatrix} x_1^2 x_2^2 & x_1^3 x_2^2 & x_1^2 x_2^3 \\ x_1^3 x_2^2 & x_1^4 & x_1^2 x_2^2 \\ x_1^2 x_2^3 & x_1^2 x_2^2 & x_2^4 \end{bmatrix} dx_1 dx_2$$

$$\mu_{12} = \iint_R \begin{bmatrix} x_1 x_2 & x_1^2 & x_2^2 \\ x_1^2 x_2 & x_1^3 & x_1 x_2^2 \\ x_1 x_2^2 & x_1^2 x_2 & x_2^3 \end{bmatrix} dx_1 dx_2$$

ในปี 1959 บอซซ์และแคร์ปเปอร์ ได้เสนอบริเวณที่สนใจในแผนแบบการทดลองซึ่งบริเวณนี้จะเป็นบริเวณที่มีค่าความเอนเอียงต่ำที่สุด โดยมีเงื่อนไขคือ

$$\mathbf{M}'_{11} \mathbf{M}_{12} = \mu'_{11} \mu_{12}$$

ซึ่งบริเวณที่สนใจนี้จะมีค่าอยู่ระหว่าง $[-1, 1]$ หรือ $g \leq 1$ และสามารถคำนวณค่า g ได้ดังนี้

$$g^2 = n / [(d_1 + d_2) * (n - n_0)]$$

หรือ

$$g^2 = n / 3 * (n - n_0)$$

เมื่อ n คือจำนวนตัวอย่างทั้งหมด

n_0 คือจำนวนตัวอย่าง (การทำซ้ำ) ที่จุดศูนย์กลาง

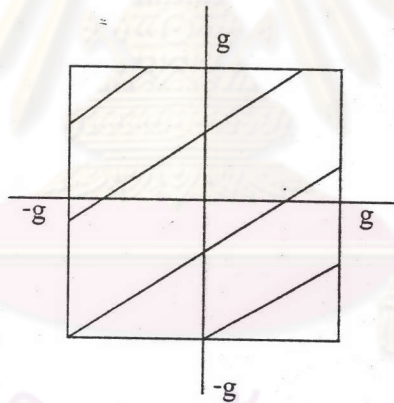
d_1 คือจำนวนของกำลังที่ใช้ในการประมาณ ในที่นี้เท่ากับ 1

และ d_2 เท่ากับ $d_1 + 1$ ซึ่งในที่นี้เท่ากับ 2

โดยสามารถเขียนเมทริกซ์ของแผนแบบการทดลอง (Design matrix) ได้ดังนี้

$$D = \begin{bmatrix} \pm g & \pm g & \dots & \pm g \\ 0 & 0 & \dots & 0 \end{bmatrix}$$

ในกรณีการทดลองแบบแฟคทอเรียลขนาด 2^2 จะมีบริเวณที่สนใจดังนี้



การทดสอบความเหมาะสมของรูปแบบ (lack of fit test)

การเลือกแผนแบบการทดลองที่เหมาะสมจะมีการพิจารณาจากหลาย ๆ อย่าง เช่น อันดับแรกจะพิจารณาจากความคลาดเคลื่อนกำลังสองเฉลี่ยมีค่าน้อยที่สุด และการพิจารณาแผนแบบที่มีการสืบค้นความไม่พอเพียงของรูปแบบ (good detection of mode inadequacy) โดยสามารถพิจารณาได้จากค่ากำลังของการทดสอบความเหมาะสมของรูปแบบ (Power of the lack of fit test) ซึ่งในที่นี้จะดูได้จากค่าพารามิเตอร์ไร้ศูนย์กลาง คือ

$$\lambda = n\beta_2'L\beta_2/\sigma^2$$

เมื่อ

$$L = M_{22} - M'_{12}M^{-1}_{11}M_{12}$$

β_2 คือสัมประสิทธิ์การถดถอยอันดับที่สอง (second order coefficient)

และ n คือขนาดตัวอย่าง

ในปี ค.ศ. 1972 แอทกินสัน (Atkinson) ได้เสนอแผนแบบการทดลองที่ทำให้ค่า $|L|$ มีค่ามากที่สุด ซึ่งถ้ารูปแบบความสัมพันธ์ของตัวแปรตามกับตัวแปรในแผนแบบการทดลองเป็นรูปแบบกำลังสอง นั่นคือ $\eta(x) = X_1\beta_1 + X_2\beta_2$ เป็นจริง ค่าของ λ จะมีค่าสูง เพราะ λ มีการแจกแจงไคสแควร์ ซึ่งเป็นค่าผลรวมกำลังสองของความถดถอย (Sum Squares of Regression (SSR)) ของสัมประสิทธิ์การถดถอยอันดับที่สอง (β_2) ซึ่งถ้ามีค่ามากแสดงว่า β_2 จะมีค่าไม่เท่ากับศูนย์ และแผนแบบการทดลองแบบใดที่ทำให้ค่า λ มีค่ามาก แสดงว่าแผนแบบการทดลองนั้นจะให้ค่ากำลังของการทดสอบความเหมาะสมของรูปแบบสูง

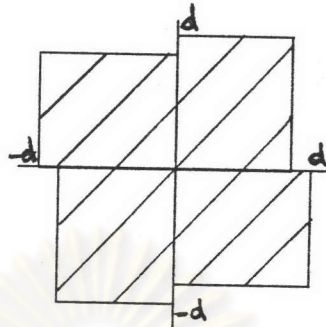
ถ้าจะพิจารณาจากเงื่อนไขดังกล่าวทั้ง 2 ข้อไปพร้อม ๆ กัน คือแผนแบบการทดลองที่มีค่าความเอนเอียงต่ำที่สุด และแผนแบบการทดลองที่ให้ค่าพารามิเตอร์ไร้ศูนย์กลางสูงที่สุด จะไม่สามารถหาค่า g ที่ทำให้เกิดเงื่อนไขทั้งสองข้อนี้พร้อมกันได้ ดังนั้นในปี ค.ศ. 1992 แพททริก ดีฟีโอ และเรย์มอนด์ เอช. ไมเยอร์ (Patrick DeFeo and Raymond H. Myers) จึงได้เสนอบริเวณที่สนใจในแผนแบบการทดลอง โดยบริเวณที่สนใจที่ใช้พิจารณานี้จะเป็นบริเวณที่ให้ค่าความเอนเอียงคงที่ และสามารถหาค่าเฉลี่ยของพารามิเตอร์ไร้ศูนย์กลางได้ในเทอมของค่าความเอนเอียงคงที่ โดยจะมีความสอดคล้องกับบริเวณที่สนใจของแผนแบบการทดลองแบบแฟกทอเรียล โดยการแปลง (transform) ด้วยเมทริกซ์ R ซึ่งในกรณีแผนแบบการทดลองแบบแฟกทอเรียลขนาด 2^2 สามารถสร้างเมทริกซ์ R ได้ดังนี้

$$R = \begin{bmatrix} \cos\theta & \sin\theta \\ -\sin\theta & \cos\theta \end{bmatrix}$$

เมื่อ θ คือมุมคงที่ และสามารถหาเมทริกซ์ของแผนแบบการทดลองได้ดังนี้

$$D' = \begin{bmatrix} -d & d & -1 & 1 & 0 \\ -1 & 1 & d & -d & 0 \end{bmatrix}$$

โดยบริเวณที่สนใจในแผนแบบการทดลองนี้เป็นดังนี้



โดย g และ d มีความสัมพันธ์กันดังนี้

$$d^2 = \sqrt{(2g^2) - 1}$$

และ $d \leq 1$

การหาค่าเฉลี่ยของพารามิเตอร์ไร้ศูนย์กลางสามารถทำได้ดังนี้

$$\begin{aligned} \lambda_{ave}^* &= \int_{\phi^*} \lambda d\beta_2 / \int_{\phi^*} d\beta_2 \\ &= \delta \text{tr}(\mathbf{T}^{-1}\mathbf{L}) / p_2 \end{aligned}$$

เมื่อ $\delta = \beta_2 \mathbf{T}' \beta_2 / \sigma^2$

β_2 คือสัมประสิทธิ์การถดถอยอันดับที่สอง (second order coefficient)

และ p_2 คือจำนวนมิติ (dimension) ของ β_2

การประมาณค่าสัมประสิทธิ์การถดถอยพหุคูณโดยวิธีกำลังสองน้อยที่สุดสามัญ

วิธีการประมาณค่าสัมประสิทธิ์การถดถอยพหุคูณนี้มีรากฐานมาจากทฤษฎีการประมาณเชิงเส้นที่คิดขึ้นโดยคาร์ล เฟร德里ก เกาส์ (Carl Friedrich Gauss) ในปี ค.ศ. 1777-1855 และอังเดร แอนดรีวิช มาร์คอฟ (Andrei Andreevich Markov) ในปี ค.ศ. 1855-1922 โดยหลักใน

การประมาณค่าสัมประสิทธิ์คือ ทำให้ผลรวมกำลังของความคลาดเคลื่อน (Sum of Squares Error (SSE)) มีค่าน้อยที่สุด ซึ่งแสดงรายละเอียดดังนี้

นิยามที่ 1

จากสมการ $y = X\beta + \varepsilon$ จะได้ว่าตัวประมาณกำลังสองน้อยที่สุดของ β คือ $\hat{\beta}$ ที่ทำให้ผลรวมกำลังสองของความคลาดเคลื่อนมีค่าน้อยที่สุด

จากนิยามที่ 1 เราจะทำการหาตัวประมาณกำลังสองน้อยที่สุดได้ดังนี้

กำหนด $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)'$ จะได้ว่า $\hat{y} = X\hat{\beta}$ ดังนั้น ε เป็นเวกเตอร์ของความคลาดเคลื่อน ซึ่งผลบวกกำลังสองของความคลาดเคลื่อนคือ

$$\begin{aligned} SSE &= (\hat{\varepsilon}'\hat{\varepsilon}) \\ &= (y - X\hat{\beta})'(y - X\hat{\beta}) \\ &= (y' - \hat{\beta}'X')(y - X\hat{\beta}) \\ &= y'y - y'X\hat{\beta} - \hat{\beta}'X'y + \hat{\beta}'X'X\hat{\beta} \\ &= y'y - 2\hat{\beta}'X'y + \hat{\beta}'X'X\hat{\beta} \end{aligned}$$

การหาค่าน้อยที่สุดของผลรวมกำลังสองของความคลาดเคลื่อนทำได้โดยการหาอนุพันธ์ (differentiate) เทียบกับ β_i ; $i = 0, 1, \dots, p$ (p เป็นจำนวนตัวแปร) แล้วกำหนดให้เท่ากับ 0

$$\frac{\partial(\varepsilon'\varepsilon)}{\partial\beta} = \frac{\partial(y'y - 2\hat{\beta}'X'y + \hat{\beta}'X'X\hat{\beta})}{\partial\beta} = 0$$

$$\text{จะได้ว่า } -2X'y + 2X'X\hat{\beta} = 0$$

$$(X'X)\hat{\beta} = X'y$$

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

และเมทริกซ์ความแปรปรวนร่วมของตัวประมาณ β คือ

$$\begin{aligned} \text{cov}(\beta) &= \text{cov}[(X'X)^{-1}X'y] \\ &= (X'X)^{-1}X' \text{cov}(y) X(X'X)^{-1} \\ &= (X'X)^{-1}X' \sigma^2 X(X'X)^{-1} \\ &= \sigma^2 (X'X)^{-1} \end{aligned}$$

การประมาณค่าสัมประสิทธิ์การถดถอยพหุคูณโดยวิธีกำลังสองน้อยที่สุดสามัญกรณีที่เมทริกซ์ X มีลำดับชั้นไม่เต็ม

ในกรณีที่ X มีลำดับชั้นไม่เต็ม (not full rank) จะมีผลทำให้ไม่สามารถหาเมทริกซ์ผกผันของ $X'X$ ได้ ทำให้สมการปกติมีคำตอบได้หลายคำตอบ (no unique) วิธีหนึ่งที่จะหา คำตอบหนึ่งในนั้นได้ คือการหาเมทริกซ์ผกผันทั่วไป (Generalized inverse matrix) ซึ่งในที่นี้แทนด้วย G แล้วนำเมทริกซ์ผกผันทั่วไปที่ได้นี้ไปแทนค่าในสมการปกติ ซึ่งวิธีหาเมทริกซ์ผกผันทั่วไปสามารถทำเป็นขั้นตอนได้ดังนี้

1. ในเมทริกซ์ A ซึ่งมีลำดับชั้นเท่า r และหาเมทริกซ์ไม่เป็นเอกฐาน (non-singular) ไมเนอร์ r เรียกว่าเมทริกซ์ M
2. หาเมทริกซ์ผกผันของเมทริกซ์ M (M^{-1}) และหาเมทริกซ์สลับเปลี่ยน (transpose matrix) ของเมทริกซ์ผกผันของเมทริกซ์ M ($(M^{-1})'$)
3. แทนสมาชิกแต่ละตัวของเมทริกซ์ A ด้วยสมาชิกแต่ละตัวของ $(M^{-1})'$ ในตำแหน่งของเมทริกซ์ M
4. แทนสมาชิก (element) ตำแหน่งอื่นๆ ของเมทริกซ์ A ด้วย 0
5. หาเมทริกซ์สลับเปลี่ยนของเมทริกซ์ที่ได้ในข้อ 4.
6. เมทริกซ์ที่ได้ในข้อ 5. คือเมทริกซ์ผกผันทั่วไป (G) ของเมทริกซ์ A

ในกรณีที่เมทริกซ์ X คือ เมทริกซ์ของแผนแบบการทดลองและมีลำดับชั้นไม่เต็ม ซึ่งทำให้ $X'X$ ซึ่งมี $(p-r)$ แถว ที่มีค่าเป็น 0 และลำดับชั้นของ $X'X$ เท่ากับ r สามารถสร้างเมทริกซ์ผกผันทั่วไป (G) ของ $X'X$ ได้ดังนี้

$$X'X = \begin{bmatrix} X_1'X_1 & X_1'X_2 \\ X_2'X_1 & X_2'X_2 \end{bmatrix}$$

$$G = \begin{bmatrix} (X_1'X_1)^{-1} & 0 \\ 0 & 0 \end{bmatrix}$$

เนื่องจาก $X_1'X_1$ เป็นเมทริกซ์แนวทแยงมุมทำให้เมทริกซ์สลับเปลี่ยนของ $X_1'X_1$ เท่ากับ $X_1'X_1$ โดยที่เมทริกซ์ $X_1'X_1$ มีลำดับชั้นเต็ม (full rank) ซึ่งเท่ากับลำดับชั้นของ X จะได้ว่า G เป็นเมทริกซ์ผกผันทั่วไปของ $X'X$ โดยที่ $X = [X_1 \ X_2]$ เมื่อ X_1 มีลำดับชั้นเต็ม

และ $X_2 = X_1 M$ สำหรับบางเมทริกซ์ M เพราะสมาชิกแต่ละตัวของ X คือ -1 0 และ 1 ดังนั้น M จะมีสมาชิกแต่ละตัวเป็น -1 0 และ 1 และเมื่อหาเมทริกซ์ผกผันทั่วไปได้แล้ว นำมาแทนค่าในสมการปกติดังนี้

$$\begin{aligned}\hat{\beta}^0 &= GX'y \\ E(\hat{\beta}^0) &= GX'E(y) \\ &= GX'X\beta \\ &= H\beta\end{aligned}$$

ค่าคาดหวังของ $\hat{\beta}^0$ คือ $H\beta$ เมื่อ $H = GX'X$ ซึ่งในที่นี้จะทำให้ $\hat{\beta}^0$ เป็นตัวประมาณที่ไม่เอนเอียงของ $H\beta$ ซึ่งไม่ใช่ของ β และมีค่าความแปรปรวนดังนี้

$$\begin{aligned}V(\hat{\beta}^0) &= V(GX'y) \\ &= GX'V(y)XG' \\ &= GX'XG'\sigma^2\end{aligned}$$

โดยที่ค่าคาดหวังของ y สามารถประมาณได้ดังนี้

$$\begin{aligned}E(y) &= y \\ &= X\hat{\beta}^0 \\ &= XGX'y\end{aligned}$$

การประมาณค่าความคลาดเคลื่อนกำลังสอง

$$\begin{aligned}SSE &= (y - X\hat{\beta}^0)'(y - X\hat{\beta}^0) \\ &= y'(I - XGX')(I - XGX')y \\ &= y'(I - XGX')y\end{aligned}$$

เพราะว่า $(I - XGX')$ เป็นไอเดมโพเทนต์ และเป็นเมทริกซ์สมมาตรเพราะ XGX' เท่ากับ G ดังนั้นผลบวกกำลังสองของความคลาดเคลื่อนจึงสามารถหาได้จากสมการปกติที่เกิดจาก $\hat{\beta}^0$ ซึ่งสามารถเขียนได้ดังนี้

$$\begin{aligned}\text{SSE} &= \mathbf{y}'(\mathbf{I} - \mathbf{XGX}')(\mathbf{I} - \mathbf{XGX}')\mathbf{y} \\ &= \mathbf{y}'\mathbf{y} - \mathbf{y}'\mathbf{XGX}'\mathbf{y} \\ &= \mathbf{y}'\mathbf{y} - \hat{\beta}^0\mathbf{X}'\mathbf{y}\end{aligned}$$



ศูนย์วิทยทรัพยากร
จุฬาลงกรณ์มหาวิทยาลัย