

บทที่ 2

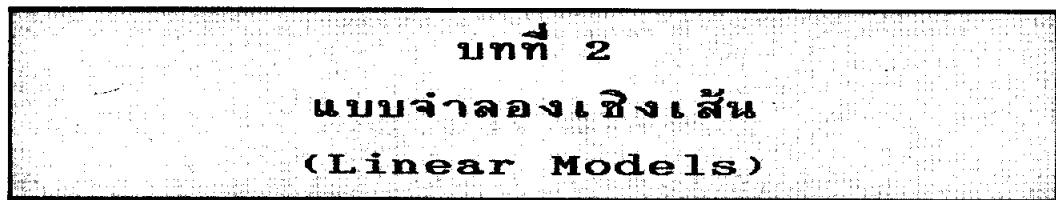
แบบจำลองเชิงเส้น (Linear Models)

หน้า

2.0 ค่าสำหรับ

2.1 ค่าจำากัดความ

2.2 การทำแทรกแบบจำลอง (ความวิธีของนักสถิติ)



2.0 ค่านำ

ความสัมพันธ์ในรูปฟังก์ชัน (Functional Relationship)

ถ้าให้ y เป็นปริมาณหนึ่งซึ่งเราต้องการท่านาย

x_1, x_2, \dots, x_n เป็นตัวแปรจำนวน n ตัว
และ ฟังก์ชัน g ซึ่งทำให้ $y = g(x_1, x_2, \dots, x_n)$

ถ้าเราสามารถสังเกตค่า x_1, x_2, \dots, x_n ได้ เราจะสามารถท่านายค่า y
ได้อีกต่อไป แต่เราจะกล่าวว่า

" y, x_1, x_2, \dots, x_n มีความสัมพันธ์ในรูปฟังก์ชัน"

ตัวอย่างเช่น เราต้องการท่านายความสูง (h) ของคนหนึ่ง เราทราบว่าเราไม่
อาจท่านายความสูงของคนคนนี้ได้อีกเมื่อจากกระบวนการน้ำหนัก (w) ของเขามาก
อย่างเดียว อย่างไรก็ได้เราอาจใช้ปริมาณอื่น ๆ มาช่วยได้ เช่น น้ำหนักของบิดา (x_1)
ความสูงของบิดา (x_2) น้ำหนักของมารดา (x_3) เป็นต้น

เรารسمุติว่าสมการอยู่ในรูป $h = g(w, x_1, x_2, \dots, x_n)$
หรืออาจสมมุติว่าเราสามารถเขียนสมการได้ในรูป

$$h = \alpha + \beta w + f(x_1, x_2, \dots, x_n)$$

และสมมุติอีกว่า สำหรับค่า w ที่กำหนดให้ ปริมาณ x_1, x_2, \dots, x_n จะเปลี่ยนไป และ
 $f(x_1, x_2, \dots, x_n)$ ทำหน้าที่เป็นความคลาดเคลื่อนสุ่ม (random error)

ถึงแม้ว่าเราทราบฟังก์ชัน g ข้างต้น และทราบว่า x_i ตัวใดบ้างที่ใช้ท่านายค่า y
เราที่ยังไม่สามารถท่านาย y ได้อีกต่อไปนั่น เพราะเราไม่สามารถวัดค่า x_i
บางตัวหรือทุกตัวได้อีกเมื่อ ทั้งนี้เพราะเราไม่สามารถแยกความคลาดเคลื่อนจากการวัด

ความคลาดเคลื่อนที่สัมภัย 2 ชนิดคือ

1) ความคลาดเคลื่อนจากการใช้สมการ (Equation Error)

เกิดจากการไม่ทราบพังก์ชัน g หรือการใช้พังก์ชัน g ที่ผิด เนื่องจากไม่ทราบว่า x_i ใด ข้างที่จะใช้กำหนด y จึงทำให้ใช้ x_i ไม่ถูกต้องหรือไม่ครบถ้วน

2) ความคลาดเคลื่อนจากการวัด (Measurement Error)

เกิดจากการที่ไม่สามารถวัด x_i ได้อย่างแม่นยำ

2.1 คำจำกัดความ

คำจำกัดความ 2.1 เมื่อกล่าวถึง แบบจำลอง (Model) เราหมายถึง สมการทางคณิตศาสตร์ (Mathematical equation) ซึ่งประกอบด้วย ตัวแปรเชิงสุ่ม (Random variables) ตัวแปรทางคณิตศาสตร์ (Mathematical variables) และ พารามิเตอร์ (Parameters)

คำจำกัดความ 2.2 เมื่อกล่าวถึง แบบจำลองเชิงเส้น (Linear Model) เราหมายถึง สมการซึ่งประกอบด้วยตัวแปรเชิงสุ่ม ตัวแปรทางคณิตศาสตร์ และ พารามิเตอร์ และเป็นเชิงเส้นในหมายความว่า

ตัวอย่าง ถ้า χ_0, χ_1, χ_2 เป็นพารามิเตอร์ที่ไม่ทราบค่า

$$\chi_0 + \chi_1 X + \chi_2 Y = 0 \quad \text{เป็นแบบจำลองเชิงเส้น}$$

$$\text{และ } \chi_0 + \chi_1 e^x + \chi_2 Y \cos x = 0 \quad \text{เป็นแบบจำลองเชิงเส้น}$$

$$\text{แต่ } \chi_0^2 + X \sin \chi_1 + \chi_2 = 0 \quad \text{ไม่เป็นแบบจำลองเชิงเส้น}$$

ในบทนี้สัญลักษณ์ต่าง ๆ ที่ใช้คือ

ให้ U^*, V^*, X^*, \dots แทนตัวแปรทางคณิตศาสตร์ที่สังเกตค่าไม่ได้

U, V, X, \dots แทนตัวแปรทางคณิตศาสตร์ที่สังเกตค่าได้

$y_i^*, u_i^*, v_i^*, \dots$ แทนตัวแปรเชิงสัมที่สังเกตค่าไม่ได้

y_i, u_i, v_i, \dots แทนตัวแปรเชิงสัมที่สังเกตค่าได้

a, b, d, e, \dots แทนตัวแปรเชิงสัมที่ไม่สามารถสังเกตค่าได้

และ อักษรกรีก แทนพารามิเตอร์ที่ไม่ทราบค่า

2.2 การจำแนกแบบจำลอง (ตามวิธีของนักสถิติ)

แบบจำลอง 1 General Linear Hypothesis of Full Rank

$$y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + e$$

หรือเขียนแบบจำลองในรูปเมตริกซ์ดังนี้คือ

$$Y = X\beta + e$$

โดยที่ X_1, X_2, \dots, X_k เป็นตัวแปรทางคณิตศาสตร์ที่ทราบค่า

y เป็นตัวแปรเชิงสัมที่สังเกตค่าได้

$\beta_0, \beta_1, \dots, \beta_k$ เป็นพารามิเตอร์ที่ไม่ทราบค่า

e เป็นตัวแปรเชิงสัมที่สังเกตค่าไม่ได้ และมี mean เป็น 0

แบบจำลอง 2 Functionally Related Models

(with Variables Subject to Measurement Error)

แบบจำลองนี้แสดงความสัมพันธ์ในรูปฟังก์ชันของตัวแปรทางคณิตศาสตร์
ถ้า Y^*, X_1^*, \dots, X_k^* สังเกตค่าไม่ได้

y, X_1, \dots, X_k สังเกตค่าได้ โดยที่ $y = Y^* + e$

$$X_i = X_i^* + e_i, \quad i = 1, \dots, k$$

e และ e_i 's คือความคลาดเคลื่อนจากการวัด (Measurement error)

$$\text{จากความสัมพันธ์ } Y^* = \beta_0 + \beta_1 X_1^* + \dots + \beta_k X_k^*$$

$$\text{หรือ } y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + e - \beta_1 e_1 - \beta_2 e_2 - \dots - \beta_k e_k$$

ตัวอธิบาย เช่น กฎของโอห์ม (Ohm's Law) กล่าวว่า จำนวนไฟล์ V^* ในวงจร
หนึ่งมีค่าเท่ากับผลคูณของความต้านทาน μ ของเส้น導 คับกระแสไฟฟ้า I

สมการคือ $V^* = \mu I$ สมมุติว่าผู้ทดลองต้องการหาความต้านทาน μ ของวงจร สมมุติว่า
เขายังไม่สามารถวัด V^* แต่เขาวัดได้ค่า $v = V^* + e$ แล้วสมการกล้ายเป็น

$$v = \mu I + e \quad \text{ในที่นี้ } v \text{ และ } e \text{ เป็นตัวแปรเชิงสุ่ม}$$

I เป็นตัวแปรทางคณิตศาสตร์ และ

μ คือพารามิเตอร์ที่ไม่ทราบค่าซึ่งเราจะทำการประมาณค่า

แบบจำลอง 3 Regression Models

y, x_1, \dots, x_k เป็นเซ็ตของตัวแปรเชิงสุ่มซึ่งมีฟังก์ชันการแจกแจงร่วม

ที่มี

$$E(Y | x_1 = X_i) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k \quad , \quad i = 1, \dots, k$$

$$\text{หรือ } y_x = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + e \quad \text{โดย } E(e) = 0$$

หมายเหตุ ข้อแตกต่างระหว่างแบบจำลอง 1 และแบบจำลอง 3 คือ

x_i ในแบบจำลอง 1 เป็นตัวแปรทางคณิตศาสตร์

x_i ในแบบจำลอง 3 เป็นค่าของตัวแปรเชิงสุ่ม

ตัวอธิบาย เช่น สมมุติว่าเราต้องการทำนายอุณหภูมิ y ในบริเวณหนึ่งโดยการใช้เพียง
ค่าความชื้น x เราสมมุติว่า x และ y มาจากฟังก์ชันการแจกแจงร่วมแบบ bivariate
เราสามารถใช้แบบจำลองการ回帰และทำนายอุณหภูมิเหลือสำหรับค่าความชื้นที่กำหนดไว้ที่
ค่าหนึ่ง ๆ

แบบจำลอง 4 Experimental Design Models

(General Linear Hypothesis of Less than Full Rank)

$$y = \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + e$$

โดยที่ y เป็นตัวแปรเชิงสุ่ม

X_1, X_2, \dots, X_k มีค่า 0 หรือ 1

$\beta_1, \beta_2, \dots, \beta_k$ เป็นพารามิเตอร์ที่ไม่ทราบค่า

e เป็นความคลาดเคลื่อนสุ่ม

แบบจำลอง 4 เป็นกรณีพิเศษของแบบจำลอง 1 แต่เนื่องจากความสำคัญของแบบจำลองนี้จึงแยกออกมาพิจารณาต่างหาก

โดยทั่ว ๆ ไปอาจเขียนแบบจำลองได้ดังนี้

$$y_{i,j,\dots,k} = \mu + \beta_i + \alpha_j + \dots + e_{i,j,\dots,k}$$

โดยที่ β_i, α_j, \dots เป็นพารามิเตอร์ที่ไม่ทราบค่า
 $e_{i,j,\dots,k}$ เป็นความคลาดเคลื่อนสุ่ม

แบบจำลอง 5 Component-of-Variance Models

a_1, a_2, \dots, a_p เป็นตัวแปรเชิงสุ่มซึ่งสังเกตค่าไม่ได้จากการแจกแจงซึ่งมีค่าเฉลี่ย (mean) เท่ากับ 0 และ ความแปรปรวน (variance) เท่ากับ σ_a^2

$b_{11}, b_{12}, \dots, b_{pq}$ เป็นตัวแปรเชิงสุ่มซึ่งสังเกตค่าไม่ได้จากการแจกแจงซึ่งมีค่าเฉลี่ยเท่ากับ 0 และ ความแปรปรวนเท่ากับ σ_b^2

$y_{i,j}$ เป็นตัวแปรเชิงสุ่มซึ่งสังเกตค่าได้

$$y_{i,j} = \mu + a_i + b_{ij}$$

โดยที่ μ เป็นพารามิเตอร์ที่ไม่ทราบค่า

$$\text{ดังนั้น } \text{Var}(y_{i,j}) = \sigma_a^2 + \sigma_b^2$$

เราอาจเขียนแบบจำลองนี้ในรูป

$$y = aX_1 + bX_2 + e$$

แบบจำลองนี้คล้ายกับแบบจำลอง 4 ตรงที่ X_i 's มีค่า 0 หรือ 1 เท่านั้น
ส่วน a และ b ไม่ใช้พารามิเตอร์ แต่เป็นตัวแปรเชิงสุ่มจากการแจกแจงซึ่งมี
ความแปรปรวนเท่ากับ σ_a^2 และ σ_b^2 ตามลำดับ

จุดประสงค์ของแบบจำลองชนิดนี้คือ การสังเกตค่าของ y แล้วประมาณค่า
 σ_a^2 , σ_b^2 และ $\text{Var}(e) = \sigma^2$

ตัวอย่างเช่น ในการวัดปริมาณในโครงการนึงในนี้ของตัวอย่างนี้ ที่มาหลัก ๆ
ของความแปรปรวนมีอยู่ 2 ที่มา คือ จากความแปรปรวนของใบไม้บนต้นไม้ และ
จากความแปรปรวนเนื่องจากการวัดหรือการทดลอง สมมุติว่าเราเอาใบไม้มา n ใบจาก
ต้นไม้ซึ่งมีปริมาณในโครงการที่แท้จริงของใบไม้ใบที่ i เท่ากับ a_i ในกรณี a_i เป็น
ตัวแปรเชิงสุ่มจากการแจกแจงซึ่งมีความแปรปรวน σ_a^2 สมมุติว่า ทำการวัด m ครั้งจาก
ใบไม้แต่ละใบเพื่อหาปริมาณในโครงการนึงใบไม้ใบนั้น

ถ้าให้ $y_{i,j}$ แทนปริมาณในโครงการที่วัดได้ในครั้งที่ j ของใบไม้ใบที่ i

$e_{i,j}$ เป็นตัวแปรเชิงสุ่มจากการแจกแจงซึ่งมีความแปรปรวน σ^2

เราเขียนแบบจำลองได้ในรูป

$$y_{i,j} = \mu + a_i + e_{i,j}$$

μ เป็นค่าคงที่ ซึ่งคือค่าเฉลี่ยของ $y_{i,j}$

จุดประสงค์หนึ่งของแบบจำลองชนิดนี้คือ การสังเกตค่าของ y แล้วประมาณค่า σ_a^2
และ $\text{Var}(e) = \sigma^2$

โปรดสังเกตว่าความแปรปรวนของ $y_{i,j}$ เท่ากับ $\sigma_a^2 + \sigma^2$ ซึ่งเป็นพังก์ชัน
เชิงเส้นของความแปรปรวน

ในบทที่ 3 - 5 จะกล่าวถึงแบบจำลอง 1 ในบทที่ 6 จะกล่าวถึงแบบจำลอง 3
ส่วนในบทที่ 7 จะกล่าวถึงแบบจำลอง 4 และกล่าวถึงแบบจำลอง 5 เพื่อแก้ไข