

ทฤษฎีที่ใช้ในการวิจัย

สมการถดถอยซึ่งประมาณโดยวิธี OLS ที่นำมาเปรียบเทียบมี 6 สมการด้วยกัน คือ

1. สมการถดถอย เมื่อไม่มีข้อมูลสูญหายเลยโดยจะเรียกว่าวิธีวิเคราะห์สมบูรณ์
2. สมการถดถอย เมื่อมีข้อมูลสูญหายโดยตัดชุดข้อมูลที่มีค่าสูญหายออกโดยจะเรียกว่าวิธีสูญหาย
3. สมการถดถอย เมื่อมีข้อมูลสูญหายโดยประมาณข้อมูลสูญหายจากสมการถดถอยเชิงเส้นโดยจะเรียกว่าวิธีวิเคราะห์ความถดถอย
4. สมการถดถอย เมื่อใช้ตัวประมาณ maximum likelihood ประมาณข้อมูลสูญหายโดยจะเรียกว่า วิธี MAXIMUM LIKELIHOOD
5. สมการถดถอย เมื่อใช้ค่าเฉลี่ยประมาณข้อมูลสูญหาย โดยจะเรียกว่าวิธีค่าเฉลี่ย
6. สมการถดถอย เมื่อใช้ค่ามัธยฐานประมาณข้อมูลสูญหาย โดยจะเรียกว่าวิธีค่ามัธยฐาน

การประมาณค่าสัมประสิทธิ์ความถดถอยด้วยวิธีกำลังสองน้อยที่สุด (ordinary least square : OLS)

วิธีการหาค่าประมาณพารามิเตอร์โดยวิธี OLS เป็นวิธีที่มีรากฐานมาจากทฤษฎีการประมาณค่าเชิงเส้นซึ่งเป็นวิธีที่คิดค้นโดย คาร์ล เฟรคิช เกาส์ (Karl Friedrich Gauss 1777 - 1855) และอังเดร แอนดรีวิช มาร์คอฟ (Andrei Andreevich Markov 1856 - 1922) โดยมีหลักเกณฑ์ดังนี้คือให้หาค่าตัวประมาณค่าพารามิเตอร์ที่ทำให้ผลบวกกำลังสองของผลต่างของค่าสังเกตกับค่าคาดหวังของตัวแปรที่มีค่าต่ำที่สุด

รูปแบบจำลองทั่วไปของสมการถดถอยเชิงเส้นคือ

$$Y = X\beta + \epsilon \quad \dots (2.3.1)$$

- เมื่อ  $\underline{Y}$  คือ เมตริกซ์ของตัวแปรตามขนาด  $N \times 1$   
 $X$  คือ เมตริกซ์ของตัวแปรอิสระขนาด  $N \times M$   
 $\underline{\beta}$  คือ เมตริกซ์สัมประสิทธิ์ความถดถอยเชิงเส้นขนาด  $M \times 1$   
 $\underline{e}$  คือ เมตริกซ์ของความคลาดเคลื่อนขนาด  $N \times 1$  โดยที่  
 $\underline{e} \sim N(0, \sigma^2 I_N)$  และ  
 $N$  คือ ขนาดตัวอย่าง  
 $M$  คือ จำนวนตัวแปรอิสระ + 1

ตัวประมาณค่าสัมประสิทธิ์ความถดถอยเชิงเส้น โดยวิธีกำลังสองน้อยที่สุดของ  $\underline{\beta}$  คือ  $\hat{\underline{\beta}}$  ที่ทำให้ผลบวกกำลังสองของความคลาดเคลื่อน (sum of squares error) มีค่าน้อยที่สุด

จากแบบจำลอง (2.3.1) ผลรวมของความคลาดเคลื่อนกำลังสองคือ

$$\begin{aligned} S(\underline{\beta}) &= \underline{e}'\underline{e} = (\underline{Y} - X\underline{\beta})'(\underline{Y} - X\underline{\beta}) \\ &= \underline{Y}'\underline{Y} - \underline{\beta}'X'\underline{Y} - \underline{Y}'X\underline{\beta} + \underline{\beta}'X'X\underline{\beta} \\ &= \underline{Y}'\underline{Y} - 2\underline{\beta}'X'\underline{Y} + \underline{\beta}'X'X\underline{\beta} \end{aligned} \quad \dots (2.3.2)$$

การประมาณค่าสัมประสิทธิ์ความถดถอยโดยวิธีกำลังสองน้อยที่สุด จากตัวแบบ

(2.3.1) อาศัยวิธีการทางแคลคูลัสโดยการดิฟเฟอเรนเชียล (Differentiate) สมการ

(2.3.2) เทียบกับ  $\underline{\beta}$  แล้วให้เท่ากับศูนย์

$$\begin{aligned} \frac{\partial S}{\partial \underline{\beta}} \bigg|_{\hat{\underline{\beta}}} &= -2X'\underline{Y} + 2X'X\underline{\beta} = 0 \\ X'X\underline{\beta} &= X'\underline{Y} \end{aligned} \quad \dots (2.3.3)$$

สมการ (2.3.3) เรียกว่าสมการปกติและตัวประมาณกำลังสองน้อยที่สุดคือ

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

ดังนั้นสมการถดถอยที่ใช้ประมาณ คือ

$$\hat{Y} = X\hat{\beta}$$

โดยที่  $E(\hat{\beta}) = \beta$  และ  $V(\hat{\beta}) = (X'X)^{-1} \hat{\sigma}^2$

รายละเอียดเกี่ยวกับวิธีประมาณค่าสัมประสิทธิ์การถดถอยแต่ละวิธีเป็นดังนี้

### 2.1 วิธีวิเคราะห์ถดถอย

จะทำการประมาณค่าสัมประสิทธิ์การถดถอยจากชุดข้อมูลครบถ้วนโดยวิธี OLS ดังนี้

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

### 2.2 วิธีสุ่ม

จะทำการประมาณค่าสัมประสิทธิ์การถดถอยจากชุดข้อมูลที่ตัดค่าสุ่มออกแล้วโดยวิธี OLS ดังนี้

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

เมื่อ  $X''$  และ  $Y''$  คือชุดข้อมูลที่ตัดค่าสุ่มทิ้ง

### 2.3 วิธีวิเคราะห์ความถดถอย

จะทำการประมาณค่าสุ่มของข้อมูลตัวแปรอิสระตามขั้นตอนต่อไปนี้

2.3.1 ตัดข้อมูลสุ่มทิ้งแล้วประมาณค่าสัมประสิทธิ์การถดถอยโดยวิธี OLS ดังนี้

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

2.3.2 นำสมการถดถอยที่ได้จาก 2.3.1 มาประมาณข้อมูลสูญหายโดยพิจารณาจากสมการถดถอยเชิงเส้น

$$\hat{Y}_{i,j} = \hat{\alpha} + \hat{\beta}_1^{(LS)} X_{i,j}$$

ดังนั้น 
$$X_{i,j}^{(EST)} = \frac{\hat{Y}_{i,j} - \hat{\alpha}}{\hat{\beta}_1^{(LS)}}$$

เมื่อ 
$$\hat{\alpha} = \bar{Y} - \bar{X}_1 \hat{\beta}_1^{(LS)}$$

2.3.3 นำค่า  $X_{i,j}^{(EST)}$  แทนในค่าสูญหายของแต่ละตัวแปรอิสระ แล้วทำการประมาณค่าสัมประสิทธิ์ความถดถอยใหม่ด้วยวิธี OLS

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

ดังนั้นจะได้สมการถดถอยที่จะใช้ประมาณ คือ

$$\hat{Y} = X\hat{\beta}$$

## 2.4 วิธี MAXIMUM LIKELIHOOD

Anderson ได้ทำการประมาณตัวพารามิเตอร์ของตัวแปรที่มีการแจกแจงแบบปกติและมีข้อมูลบางตัวของตัวแปรอิสระสูญหายโดยวิธี maximum likelihood ดังนี้

ถ้าข้อมูลมีลักษณะ

$$y_1, \dots, y_n, y_{n+1}, \dots, y_m$$

$$x_1, \dots, x_n$$

Bivariate density function ของ  $x$  และ  $y$  สามารถเขียนในรูป

$$n(y, x | \mu_y, \mu_x; \sigma_y^2, \sigma_x^2; \rho) = n(y | \mu_y, \sigma_y^2) \cdot$$

$$n(x | \mu_x + \beta_{xy}y, \sigma_{xy}^2)$$

$$\text{เมื่อ } \left. \begin{aligned} J &= \mu_x - \beta_{xy}\mu_y \\ \beta_{xy} &= r\sigma_x / \sigma_y \\ \sigma_{x.y}^2 &= \sigma_x^2(1-r^2) \end{aligned} \right\} \dots (A)$$

Likelihood function คือ

$$\begin{aligned} & \prod_{\alpha=1}^n n(y_{\alpha}, x_{\alpha} | \mu_y, \mu_x; \sigma_y^2, \sigma_x^2; r) \prod_{\alpha=1}^N n(y_{\alpha} | \mu_y, \sigma_y^2) \\ &= \prod_{\alpha=1}^N n(y_{\alpha} | \mu_y, \sigma_y^2) \prod_{\alpha=1}^n n(x_{\alpha} | J + \beta_{xy}y_{\alpha}, \sigma_{x.y}^2) \dots (1) \end{aligned}$$

ค่าประมาณ maximum likelihood ของ  $\mu_y$ ,  $\sigma_y^2$ ,  $J$ ,  $\beta_{xy}$  และ  $\sigma_{x.y}^2$  สามารถหาได้โดยการ differentiate สมการที่ (1) เทียบกับ  $\mu_y$ ,  $\sigma_y^2$ ,  $J$ ,  $\beta_{xy}$ ,  $\sigma_{x.y}^2$  แล้วเทียบให้เท่ากับ 0 ซึ่งจะได้ผลลัพธ์ดังนี้

$$\hat{\mu}_y = \bar{y} = \frac{\sum_{\alpha=1}^N y_{\alpha}}{N}$$

$$\hat{\sigma}_y^2 = \frac{\sum_{\alpha=1}^N (y_{\alpha} - \bar{y})^2}{N}$$

$$\hat{\beta}_{xy} = \frac{\sum_{\alpha=1}^n (x_{\alpha} - \bar{x})(y_{\alpha} - \bar{y})}{\sum_{\alpha=1}^n (y_{\alpha} - \bar{y})^2}$$

$$\hat{J} = \bar{x} - \hat{\beta}_{xy} \bar{y}$$

$$\hat{\sigma}_{x.y}^2 = \frac{\sum_{\alpha=1}^n (x_{\alpha} - \bar{x})^2 - \hat{\beta}_{xy}^2 \sum_{\alpha=1}^n (y_{\alpha} - \bar{y})^2}{n - 1}$$

$$\hat{\sigma}_{x.y}^2 = \frac{\sum_{\alpha=1}^n (x_{\alpha} - \bar{x})^2 - \hat{\beta}_{xy}^2 \sum_{\alpha=1}^n (y_{\alpha} - \bar{y})^2}{n - 1}$$

$$\hat{\sigma}_{x.y}^2 = \frac{\sum_{\alpha=1}^n (x_{\alpha} - \bar{x})^2 - \hat{\beta}_{xy}^2 \sum_{\alpha=1}^n (y_{\alpha} - \bar{y})^2}{n - 1}$$

$$\bar{y}^m = \frac{\sum_{\alpha=1}^n y_{\alpha}}{n}, \quad \bar{x}^m = \frac{\sum_{\alpha=1}^n y_{\alpha}}{n}$$

เพราะฉะนั้น ตัวประมาณ maximum likelihood ของ  $\mu_x$ ,  $\sigma_x^2$  และ  $\rho$  หาได้จาก การแทนค่า  $f = \hat{f}$ ,  $\beta_{xy} = \hat{\beta}_{xy}$ ,  $\sigma_{xy}^2 = \hat{\sigma}_{xy}^2$  ในสมการ (A) และแก้สมการได้ดังนี้

$$\hat{\mu}_x = \bar{x}^m + \beta_{xy} (y - \bar{y}^m)$$

$$\hat{\sigma}_x^2 = \frac{\sum_{\alpha=1}^n (X_{\alpha} - \bar{x}^m)^2 + \hat{\beta}_{xy}^2 \left[ \frac{\sum_{\alpha=1}^n (\hat{\sigma}_{xy}^2 - \sum_{\alpha=1}^n (y_{\alpha} - \bar{y}^m)^2)}{n} \right]}{n}$$

$$\hat{\rho} = \hat{\beta}_{xy} \hat{\sigma}_y / \hat{\sigma}_x$$

ถ้าในกรณีที่มีตัวแปรอิสระมากกว่า 2 ตัวแปร และมีข้อมูลอยู่ในลักษณะดังนี้

$Y_1$	$X_{11}$	$X_{31}$
$Y_2$	$X_{12}$	$X_{32}$
.	.	.
.	.	.
$Y_n$	$X_{1n}$	$X_{3n}$
$Y_{n+1}$	$X_{1n+1}$	$X_{4n+1}$
.	.	.
.	.	.
$Y_{n+m}$	$X_{1n+m}$	$X_{4n+m}$
$Y_{n+m+1}$		$X_{2n+m+1}$
.		.
.		.
$Y_N$		$X_{cN}$

การประมาณตัวพารามิเตอร์  $\mu_{x_1}, \mu_{x_2}, \mu_{x_3}, \mu_{x_4}$  โดยวิธี maximum likelihood จะทำในทำนองเดียวกับวิธีดังกล่าวข้างต้นโดยพิจารณาตัวแปรทีละคู่

ในการวิจัยครั้งนี้จะนำค่าประมาณพารามิเตอร์ไปแทนในข้อมูลสุ่มหายแล้วจึงทำการประมาณสมการถดถอยด้วยวิธี OLS

2.5 วิธี OLS เมื่อพิจารณาจาก ค่าเฉลี่ย ซึ่งมีขั้นตอนการดำเนินการดังนี้

ขั้นที่ 1 ประมาณค่าสุ่มหายของแต่ละตัวแปรโดยใช้ค่าเฉลี่ยของข้อมูลที่ไม่สุ่มหายของแต่ละตัวแปรนั้น ๆ

$$\bar{x}_i = \frac{\sum_{j=1}^{k_i} x_{i,j}}{k_i} \quad ; \quad i = 1, 2, \dots, M$$

โดยที่  $\bar{x}_i$  คือ ค่าเฉลี่ยของข้อมูลที่มีอยู่ของตัวแปรที่  $i$

$k_i$  คือ จำนวนข้อมูลที่ไม่สุ่มหายของตัวแปรที่  $i$

$M$  คือ จำนวนตัวแปรอิสระ

$x_{i,j}$  คือ ค่าของข้อมูลที่ไม่สุ่มหายของตัวแปรที่  $i$  ค่าสังเกตที่  $j$

ขั้นที่ 2 นำค่าเฉลี่ยแทนในข้อมูลสุ่มหายแล้วทำการประมาณสัมประสิทธิ์ความถดถอยด้วยวิธีกำลังสองน้อยที่สุด

2.6 วิธี OLS เมื่อพิจารณาจาก ค่ามัธยฐาน ซึ่งมีขั้นตอนการดำเนินการ ดังนี้

ขั้นที่ 1 ประมาณค่าสุ่มหายของแต่ละตัวแปรอิสระโดยใช้ค่ามัธยฐานของข้อมูลที่มีอยู่ของแต่ละตัวแปรนั้น ๆ โดยแยกพิจารณา 2 กรณี

1. เมื่อตัวแปรอิสระมีจำนวนข้อมูลที่ไม่สูญหายเป็นเลขคู่

$$x_{M_i} = \frac{1}{2} \left[ \frac{x_{(i, k_i)}}{2} + \frac{x_{(i, k_i + 1)}}{2} \right]$$

โดยที่  $x_{(i, 1)} < x_{(i, 2)} < \dots < x_{(i, k_i)}$

เมื่อ  $x_{M_i}$  คือ ค่ามัธยฐานของข้อมูลที่มียู่ของตัวแปรที่  $i$

$k_i$  คือ จำนวนข้อมูลที่ไม่สูญหายของตัวแปร  $X_i$

$M$  คือ จำนวนตัวแปรอิสระ

$\frac{x_{(i, k_i)}}{2}$  คือ ค่าของข้อมูลที่ไม่สูญหายของตัวแปรที่  $i$  ค่าสังเกตที่  $\frac{k_i}{2}$

$\frac{x_{(i, k_i + 1)}}{2}$  คือ ค่าของข้อมูลที่ไม่สูญหายของตัวแปรที่  $i$  ค่าสังเกตที่  $\frac{k_i + 1}{2}$

2. เมื่อตัวแปรอิสระมีจำนวนข้อมูลที่ไม่สูญหายเป็นเลขคี่

$$x_{M_i} = \frac{x_{(i, k_i + 1)}}{2}$$

โดยที่  $x_{(i, 1)} < x_{(i, 2)} < \dots < x_{(i, k_i)}$

เมื่อ  $x_{M_i}$  คือ ค่ามัธยฐานของข้อมูลที่มียู่ของตัวแปรที่  $i$

$k_i$  คือ จำนวนข้อมูลที่ไม่สูญหายของตัวแปร  $X_i$

$M$  คือ จำนวนตัวแปรอิสระ

$\frac{x_{(i, k_i + 1)}}{2}$  คือ ค่าของข้อมูลที่ไม่สูญหายของตัวแปรที่  $i$  ค่าสังเกตที่  $\frac{k_i + 1}{2}$

ขั้นที่ 2 นำค่ามัธยฐานแทนในข้อมูลสูญหายแล้วทำการประมาณค่าสัมประสิทธิ์ที่มีความถดถอย  
ด้วยวิธีกำลังสองน้อยที่สุด



เกณฑ์ที่ใช้ในการเปรียบเทียบสมการถดถอยเชิงซ้อน เมื่อใช้การประมาณข้อมูลสูญหายของตัวแปรอิสระโดยวิธีต่าง ๆ คือ ค่าความคลาดเคลื่อนกำลังสองเฉลี่ย (MSE) ของสมการถดถอยเมื่อไม่มีข้อมูลสูญหาย

ค่าความคลาดเคลื่อนกำลังสองเฉลี่ย (mean square error : MSE) คำนวณจาก

$$MSE = \frac{\sum_{j=1}^M (Y_{i,j} - \hat{Y}_{i,j})^2}{N-M-1}$$

โดยที่  $Y_{i,j}$  คือ ค่าจริงของข้อมูลของตัวแปรตามตัวที่  $i$

$\hat{Y}_{i,j}$  คือ ค่าประมาณของข้อมูลของตัวแปรตามตัวที่  $i$

$N$  คือ ขนาดตัวอย่าง

$M$  คือ จำนวนตัวแปรอิสระที่ใช้ในการวิเคราะห์

นอกจากการใช้เกณฑ์ดังกล่าวข้างต้น ในการวิจัยครั้งนี้ยังใช้อัตราส่วนของความคลาดเคลื่อนรวมของค่าประมาณที่ได้จากสมการถดถอยเมื่อแทนค่าข้อมูลสูญหายของตัวแปรอิสระด้วยวิธีต่าง ๆ เมื่อเทียบกับวิธีที่ 1 ซึ่งอัตราส่วนดังกล่าวสามารถบอกได้ถึงความเสี่ยงของค่าประมาณของตัวแปรตามของชุดข้อมูลครบถ้วน และชุดข้อมูลที่ได้ตัดค่าสูญหายทิ้งและประมาณข้อมูลสูญหายในการวิเคราะห์การถดถอยด้วยวิธีการต่าง ๆ ซึ่งคำนวณจาก

$$\text{อัตราส่วนของค่าความคลาดเคลื่อนรวมของวิธีที่ } i = \frac{\sum_{j=1}^N (\hat{Y}_{i,j} - \hat{Y}_{1,j}) / \hat{Y}_{1,j}}{\text{เมื่อเทียบกับวิธีที่ 1}}$$

เมื่อ  $i = 2, 3, 4, 5, 6$

วิธีดำเนินการวิจัย

การวิจัยครั้งนี้ต้องการศึกษาเปรียบเทียบวิธีการประมาณค่าสหุหายในการวิเคราะห์การถดถอยเมื่อมีสมการถดถอย 6 สมการ คือ

- ก. สมการถดถอย เมื่อไม่มีข้อมูลสหุหายเลย
- ข. สมการถดถอย เมื่อมีข้อมูลสหุหายโดยตัดชุดข้อมูลที่มีค่าสหุหายออก
- ค. สมการถดถอย เมื่อประมาณข้อมูลสหุหายจากสมการถดถอยเชิงเส้น
- ง. สมการถดถอย เมื่อประมาณข้อมูลสหุหายด้วยวิธี MAXIMUM LIKELIHOOD
- จ. สมการถดถอย เมื่อใช้ค่าเฉลี่ยประมาณข้อมูลสหุหาย
- ฉ. สมการถดถอย เมื่อใช้ค่ามัชฌิมฐานประมาณข้อมูลสหุหาย

โดยจะพิจารณาเปรียบเทียบค่าคลาดเคลื่อนกำลังสองเฉลี่ยของสมการถดถอยเมื่อได้ประมาณข้อมูลสหุหายด้วยวิธีต่าง ๆ กับของสมการถดถอยเมื่อไม่มีข้อมูลสหุหายเมื่อประชากรมีการแจกแจงปกติ (normal distribution) ที่มีการกระจาย 3 ระดับคือ C.V. = 0.05, 0.20 และ 1.00, ขนาดตัวอย่างที่สนใจในการศึกษามี 3 ขนาด คือ 30, 70 และ 100, ค่าเบี่ยงเบนมาตรฐานมี 4 ระดับคือ 5, 10, 20 และ 25, จำนวนตัวแปร มี 4 ขนาดคือ 2, 3, 5 และ 7 และการสหุหายของข้อมูลมี 3 ระดับ คือ 5%, 10% และ 15% ทั้งนี้เทคนิคที่ใช้ในการสร้างแบบจำลองข้อมูล คือ วิธีมอนติคาร์โล

เนื่องจากวิธีมอนติคาร์โลเป็นเทคนิคที่ใช้ในการวิจัยครั้งนี้ ดังนั้นในตอนแรกของบทนี้ จะกล่าวถึงวิธีมอนติคาร์โลก่อน แล้วจึงเป็นขั้นตอนการวิจัย และโปรแกรมที่ใช้ในการวิเคราะห์สมการถดถอยเชิงซ้อนตามลำดับ

3.1 วิธีมอนติคาร์โล (Monte Carlo Method)

เทคนิคที่ใช้สำหรับแก้ปัญหาในการคำนวณทางคณิตศาสตร์นั้นมีอยู่หลายวิธี วิธีมอนติคาร์โลเป็นเทคนิคอย่างหนึ่งที่ใช้แก้ปัญหานี้ Hammevsley และ Handscomb (1964 : 2) กล่าวว่า