



การวิจัยในครั้งนี้สิ่งที่สนใจศึกษาคือ การตรวจสอบความเหมาะสมของตัวแบบถดถอยเชิงเส้นโดยการแบ่งข้อมูลด้วยวิธีดูเพล็กซ์ ซึ่งประกอบด้วย การแบ่งข้อมูลด้วยวิธีดูเพล็กซ์ การประมาณค่าสัมประสิทธิ์ถดถอยเชิงเส้นโดยวิธีกำลังสองน้อยที่สุด การประมาณค่ารากที่ p ของสัดส่วนระหว่างดีเทอร์มิแนนท์ของตัวแปรอิสระที่ได้จากข้อมูลที่ใช้ในการประมาณค่ากับดีเทอร์มิแนนท์ของตัวแปรอิสระที่ได้จากข้อมูลที่ใช้ในการพยากรณ์ $(|X'X|_{est}/|X'X|_{pre})^{1/p}$ และการทดสอบความคงที่ของสัมประสิทธิ์ถดถอยเชิงเส้นโดยใช้การทดสอบเข้า ซึ่งในบทนี้จะกล่าวถึงรายละเอียดของแต่ละเรื่องส่วนตอนท้ายของบทจะนำเสนอผลงานวิจัยที่เกี่ยวข้อง โดยมีรายละเอียดต่าง ๆ เป็นดังนี้

2.1 การแบ่งข้อมูลด้วยวิธีดูเพล็กซ์

การแบ่งข้อมูลด้วยวิธีดูเพล็กซ์เป็นวิธีการที่คิดขึ้นโดย อาร์. ดับบลิว. เคนนาร์ด (R.W. Kennard) ซึ่งวิธีนี้พัฒนามาจากวิธีคาเด็กซ์ (CADEX Algorithm) ของ อาร์. ดับบลิว. เคนนาร์ด และแอล. เอ. สโตน (L. A. Stone) โดยทั้ง 2 ท่านได้เสนอวิธีคาเด็กซ์ ในปี ค.ศ. 1969 แต่ก็ยังไม่เป็นที่แพร่หลายนัก จนกระทั่งในปี ค.ศ. 1977 โรนัลด์. ดี. สไน ได้นำวิธีดูเพล็กซ์มาใช้ในการแบ่งข้อมูลออกเป็น 2 ชุดคือ ชุดที่ใช้ในการประมาณค่าและชุดที่ใช้ในการพยากรณ์ เพื่อใช้สำหรับตรวจสอบความเหมาะสมของข้อมูลในตัวแบบถดถอย โดยนำเสนอวิธีการพร้อมทั้งยกตัวอย่างจากข้อมูลลักษณะต่าง ๆ สำหรับรายละเอียดของแต่ละขั้นตอนในการแบ่งข้อมูลด้วยวิธีดูเพล็กซ์จะนำเสนอตามลำดับ แต่สิ่งสำคัญที่ผู้วิจัยต้องพึงระลึกอยู่เสมอ ก่อนที่จะทำการแบ่งข้อมูลด้วยวิธีนี้ก็คือ จำนวนข้อมูลที่จะแบ่งนั้นต้องมีไม่น้อยกว่า $2P + 20$ ถึง $2P + 30$ (P เป็นจำนวนสัมประสิทธิ์ถดถอยเชิงเส้นในสมการ) เพื่อให้หองค่าอิสระของความคลาดเคลื่อนมีเพียงพอสำหรับการทดสอบสมมติฐาน นอกจากนี้ไม่จำเป็นต้องแบ่งให้ข้อมูลที่ใช้ในการประมาณค่ามีจำนวนเท่ากับข้อมูลที่ใช้ในการพยากรณ์ และการศึกษาในครั้งนี้จะทำการศึกษารณของความถดถอยเชิงเส้นอย่างง่าย และความถดถอยเชิงเส้นเชิงพหุ ซึ่งวิธีการแบ่งข้อมูลของ

ทั้ง 2 กรณีจะใช้หลักการเดียวกัน เพียงแต่ในกรณีแรกไม่ยุ่งยากเท่ากรณีหลังเพราะเหตุว่า กรณีความถดถอยเชิงเส้นเชิงพหุมีตัวแปรอิสระมากกว่า 1 ตัว จำเป็นต้องแปลงตัวแปรอิสระทุกตัวให้เป็นมาตรฐานและทำให้ตั้งฉากกัน (orthogonal independent variable) ก่อนที่จะแบ่งข้อมูลด้วยวิธีดิวาลีซ์ เพื่อให้ตัวแปรอิสระทุกตัวมีหน่วยเดียวกัน มีความแปรปรวนเดียวกัน (unit variance) และมีความเป็นอิสระจากกัน สำหรับรายละเอียดของแต่ละขั้นตอนเป็นดังนี้

2.1.1 การแปลงตัวแปรอิสระให้เป็นมาตรฐานเดียวกัน

ในกรณีที่ตัวแปรอิสระมีมากกว่า 1 ตัว บางครั้งตัวแปรอิสระแต่ละตัวอาจมีหน่วยที่แตกต่างกัน ดังนั้นเพื่อให้ตัวแปรอิสระทุกตัวมีหน่วยเดียวกันและความแปรปรวนเดียวกัน จึงจำเป็นต้องแปลงตัวแปรอิสระทุกตัวให้เป็นมาตรฐาน โดยใช้สูตรดังนี้

$$z_{ij} = \frac{x_{ij} - \bar{x}_j}{s_j^{1/2}}, \quad i = 1, 2, \dots, n$$

$$j = 1, 2, \dots, k$$

เมื่อ $s_j = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}$

n = จำนวนข้อมูล

k = จำนวนตัวแปรอิสระ

2.1.2 การทำให้ตัวแปรอิสระตั้งฉากกัน

ภายหลังจากการแปลงตัวแปรอิสระให้เป็นมาตรฐาน เพื่อให้ตัวแปรอิสระทุกตัวมีความเป็นอิสระกัน ดังนั้นจึงต้องทำให้ตัวแปรอิสระตั้งฉากกันโดยมีขั้นตอนดังนี้

$$\text{ให้ } A = Z'Z = T'T$$

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} & \dots & A_{1k} \\ A_{21} & A_{22} & A_{23} & \dots & A_{2k} \\ A_{31} & A_{32} & A_{33} & \dots & A_{3k} \\ \vdots & \vdots & \vdots & & \vdots \\ A_{k1} & A_{k2} & A_{k3} & \dots & A_{kk} \end{bmatrix} = \begin{bmatrix} t_{11} & 0 & 0 & \dots & 0 \\ t_{12} & t_{22} & 0 & \dots & 0 \\ t_{13} & t_{23} & t_{33} & \dots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ t_{1k} & t_{2k} & t_{3k} & \dots & t_{kk} \end{bmatrix} \begin{bmatrix} t_{11} & t_{12} & t_{13} & \dots & t_{1k} \\ 0 & t_{22} & t_{23} & \dots & t_{2k} \\ 0 & 0 & t_{33} & \dots & t_{3k} \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \dots & t_{kk} \end{bmatrix}$$

สมาชิกของ เมตริกซ์ T คำนวณโดยใช้วิธีของโคลลีย์ (Graybill 1976 : 229)

ซึ่งมีขั้นตอนดังนี้

2.1.2.1 หาสมาชิกแถวที่ 1 คอลัมน์ที่ 1 ของเมตริกซ์ T โดย

$$t_{11} = \sqrt{A_{11}}$$

2.1.2.2 หาสมาชิกแถวที่ 1 คอลัมน์ที่ 2 ถึง คอลัมน์ที่ k ของเมตริกซ์

T โดย

$$t_{1j} = \frac{A_{1j}}{t_{11}} = \frac{A_{1j}}{\sqrt{A_{11}}}, \quad j = 2, 3, \dots, k$$

2.1.2.3 หาสมาชิกบนเส้นทแยงมุม (diagonal) ของเมตริกซ์ T

โดย

$$t_{ii} = \sqrt{A_{ii} - \sum_{p=1}^{i-1} t_{pi}^2}, \quad i = 2, 3, \dots, k$$

2.1.2.4 หาสมาชิกที่เหลือของเมตริกซ์ T โดย

$$t_{ij} = \frac{1}{t_{ii}} \left[A_{ij} - \sum_{p=1}^{i-1} t_{pi} t_{pj} \right], j > i \text{ and } i = 2, 3, \dots, k-1$$

$$t_{ij} = 0, j < i \text{ and } i = 2, 3, \dots, k$$

จากนั้นตัวแปรอิสระที่ถูกทำให้ตั้งฉากจะมีความแปรปรวนเดียวกัน แล้วนำมาแปลงให้อยู่ในรูปใหม่โดย

$$W = ZT^{-1}$$

เมื่อแปลงตัวแปรอิสระให้อยู่ในรูปของ W เรียบร้อยแล้ว จึงนำตัวแปรที่ได้เหล่านี้มาคำนวณหาระยะทาง เพื่อที่จะแบ่งข้อมูลออกเป็น ชุดที่ใช้ในการประมาณค่า และชุดที่ใช้ในการพยากรณ์โดยใช้สูตรในการคำนวณหาระยะทางดังนี้

$$D_{ij} = \sum_{h=1}^k (W_{hi} - W_{hj})^2, i = 1, 2, \dots, n; j = 1, 2, \dots, n, i \neq j$$

เมื่อ D_{ij} เป็นระยะทางจากค่าสังเกตที่ i ไปยังค่าสังเกตที่ j รวมระยะทางในทุก ๆ ตัวแปรอิสระ

W_{hi} เป็นตัวแปรอิสระที่ h ในค่าสังเกตที่ i

W_{hj} เป็นตัวแปรอิสระที่ h ในค่าสังเกตที่ j

k เป็นจำนวนตัวแปรอิสระ

2.1.3 ขั้นตอนในการแบ่งข้อมูลด้วยวิธีดูเพล็กซ์

2.1.3.1 หาค่าสังเกต 2 ค่าใด ๆ เพื่อเป็นข้อมูลที่ใช้ในการประมาณค่า 2 ค่าแรก โดยพิจารณาหาระยะทางที่ใกล้ที่สุดจากค่าสังเกตที่ i ใด ๆ ไปยังค่าสังเกตที่ j ใด ๆ

2.1.3.2 หาค่าสังเกต 2 ค่าใด ๆ เพื่อเป็นข้อมูลที่ใช้ในการพยากรณ์ 2 ค่าแรก จากข้อมูลที่เหลือจาก 2.1.3.1 โดยพิจารณาหาระยะทางที่ใกล้ที่สุดจากค่าสังเกตที่ i ใด ๆ ไปยังค่าสังเกตที่ j ใด ๆ

2.1.3.3 หาข้อมูลที่ใช้ในการประมาณค่า ค่าต่อไป จากข้อมูลที่เหลือจาก

2.1.3.2 โดย

- ก. หาระยะทางจากทุก ๆ ค่าสังเกตที่เหลือจาก 2.1.3.2 ไปยังค่าสังเกตแรกของข้อมูลที่ใช้ในการประมาณค่า
- ข. หาระยะทางจากทุก ๆ ค่าสังเกตที่เหลือจาก 2.1.3.2 ไปยังค่าสังเกตที่ 2 ของข้อมูลที่ใช้ในการประมาณค่า
- ค. ในแต่ละค่าสังเกตเลือกค่าระยะทางที่น้อย (minimum) ระหว่างค่าระยะทางที่คำนวณได้จาก ก. และ ข. เพียงหนึ่งค่า
- ง. ในบรรดาค่าสังเกตที่เหลือจาก 2.1.3.2 ค่าสังเกต i ใด ๆ ที่ให้ค่าระยะทางมากที่สุด จากบรรดาค่าที่น้อย (maximin) จะถูกเลือกมาเป็นข้อมูลที่ใช้ในการประมาณค่า ค่าต่อไป

2.1.3.4 หาข้อมูลที่ใช้ในการพยากรณ์ ค่าต่อไป จากข้อมูลที่เหลือจาก

2.1.3.3 โดย

- ก. หาระยะทางจากทุก ๆ ค่าสังเกตที่เหลือจาก 2.1.3.3 ไปยังค่าสังเกตแรกของข้อมูลที่ใช้ในการพยากรณ์
- ข. หาระยะทางจากทุก ๆ ค่าสังเกตที่เหลือจาก 2.1.3.3 ไปยังค่าสังเกตที่ 2 ของข้อมูลที่ใช้ในการพยากรณ์
- ค. ในแต่ละค่าสังเกตเลือกค่าระยะทางที่น้อย (minimum) ระหว่างค่าระยะทางที่คำนวณได้จาก ก. และ ข. เพียงหนึ่งค่า
- ง. ในบรรดาค่าสังเกตที่เหลือจาก 2.1.3.3 ค่าสังเกต i ใด ๆ ที่ให้ค่าระยะทางมากที่สุดจากบรรดาค่าที่น้อย (maximin) จะถูกเลือกมาเป็นข้อมูลที่ใช้ในการพยากรณ์ ค่าต่อไป

- 2.1.3.5 หาข้อมูลที่ใช้ในการประมาณค่า ค่าต่อไป จากข้อมูลที่เหลือจาก
2.1.3.4 โดยวิธีเดิม (2.1.3.3)
- 2.1.3.6 หาข้อมูลที่ใช้ในการพยากรณ์ ค่าต่อไป จากข้อมูลที่เหลือจาก
2.1.3.5 โดยวิธีเดิม (2.1.3.4)
- 2.1.3.7 ทำเช่นนี้จนหมดข้อมูล จะได้ข้อมูลเป็น 2 ชุดคือ ชุดที่ใช้ในการ
ประมาณค่าและชุดที่ใช้ในการพยากรณ์

2.2 การประมาณค่าสัมประสิทธิ์ถดถอยเชิงเส้นโดยวิธีกำลังสองน้อยที่สุด

วิธีการประมาณค่าสัมประสิทธิ์วิธีนี้มีรากฐานมาจากทฤษฎีการประมาณเชิงเส้น (Theory of Linear Estimation) ที่คิดขึ้นโดย คาร์ล เฟรดริก เกาส์ (Karl Friedrich Gauss) ในปี ค.ศ. 1777-1855 และ อังเดร แอนดรีวิช มาร์คอฟ (Andrei Andreevich Markov) ในปี ค.ศ. 1856-1922 โดยมีหลักในการประมาณค่าสัมประสิทธิ์คือ ทำให้ผลบวกกำลังสองของความคลาดเคลื่อนมีค่าน้อยที่สุด ซึ่งแสดงรายละเอียดดังนี้

นิยาม 2.1 จากสมการ $Y = X\beta + \epsilon$ เมื่อ $\epsilon \sim N(0, \sigma^2 I)$ ตัวประมาณกำลังสองน้อยที่สุดของ β คือ $\hat{\beta}$ ที่ทำให้ผลบวกกำลังสองของความคลาดเคลื่อน (Sum Square Errors) หรือ SSE มีค่าน้อยที่สุด

จากนิยาม 2.1 จะทำการหาตัวประมาณกำลังสองน้อยที่สุดได้ดังนี้

$$\begin{aligned} \text{เนื่องจาก} \quad \text{SSE} &= \epsilon' \epsilon \\ &= (Y - X\hat{\beta})' (Y - X\hat{\beta}) \\ &= (Y'Y - 2\hat{\beta}'X'Y + \hat{\beta}'X'X\hat{\beta}) \end{aligned}$$

การหาค่าน้อยที่สุดของผลบวกกำลังสองของความคลาดเคลื่อนทำได้โดยการดิฟเฟอเรนเชียล (differentiate) เทียบกับ $\hat{\beta}$ แล้วกำหนดให้เท่ากับ 0 ดังนี้

$$\begin{aligned} \frac{\partial}{\partial \hat{\beta}} (Y'Y - 2\hat{\beta}'X'Y + \hat{\beta}'X'X\hat{\beta}) &= 0 \\ - 2X'Y + 2X'X\hat{\beta} &= 0 \\ \hat{\beta} &= (X'X)^{-1} X'Y \end{aligned}$$

ในการศึกษาครั้งนี้จะทำการประมาณค่าสัมประสิทธิ์ถดถอยเชิงเส้นจากข้อมูล 3 ชุด คือ ข้อมูลทั้งหมด ข้อมูลที่ใช้ในการประมาณค่า และข้อมูลที่ใช้ในการพยากรณ์ เพื่อที่จะนำไปใช้ในการทดสอบเข้า ซึ่งมีรายละเอียดดังนี้

2.2.1 ประมาณค่าสัมประสิทธิ์ถดถอยเชิงเส้นจากข้อมูลทั้งหมด

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

เมื่อ Y เป็นเวกเตอร์ของตัวแปรตามขนาด $(n \times 1)$

X เป็นเมตริกซ์ของตัวแปรอิสระขนาด $(n \times p)$

n เป็นจำนวนข้อมูลทั้งหมด

p เป็นจำนวนสัมประสิทธิ์ถดถอยเชิงเส้น

2.2.2 ประมาณค่าสัมประสิทธิ์ถดถอยเชิงเส้นจากข้อมูลที่ใช้ในการประมาณค่า

$$\hat{\beta}_e = (X'_e X_e)^{-1} X'_e Y_e$$

เมื่อ Y_e เป็นเวกเตอร์ของตัวแปรตามขนาด $(m \times 1)$

X_e เป็นเมตริกซ์ของตัวแปรอิสระขนาด $(m \times p)$

m เป็นจำนวนข้อมูลที่ใช้ในการประมาณค่า

2.2.3 การประมาณค่าสัมประสิทธิ์ถดถอยเชิงเส้นจากข้อมูลที่ใช้ในการพยากรณ์

$$\hat{\beta}_{pr} = (X'_{pr} X_{pr})^{-1} X'_{pr} Y_{pr}$$

เมื่อ Y_{pr} เป็นเวกเตอร์ของตัวแปรตามขนาด $(l \times 1)$

X_{pr} เป็นเมตริกซ์ของตัวแปรอิสระขนาด $(l \times p)$

l เป็นจำนวนข้อมูลที่ใช้ในการพยากรณ์

$$\text{และ } m + l = n$$



2.3 การประมาณค่าราคา P ของสัดส่วนระหว่างดีเทอร์มิแนนท์ของตัวแปรอิสระที่ได้จากข้อมูลที่ใช้ในการประมาณค่ากับดีเทอร์มิแนนท์ของตัวแปรอิสระที่ได้จากข้อมูลที่ใช้ในการพยากรณ์

ในการพิจารณาว่าตัวแปรอิสระของข้อมูลที่ใช้ในการประมาณค่าและของข้อมูลที่ใช้ในการพยากรณ์มีลักษณะคล้ายคลึงกันมากน้อยเพียงใด โรนัลด์.ดี. สนิ ได้แนะนำถึงการนำใช้คุณสมบัติทางสถิติของข้อมูลทั้ง 2 ชุดมาเปรียบเทียบกัน โดยพิจารณาจาก

$$(|X'X|_{est} / |X'X|_{pre})^{1/p}$$

เมื่อ $|X'X|_{est}$ เป็นดีเทอร์มิแนนท์ของตัวแปรอิสระที่ได้จากข้อมูลที่ใช้ในการประมาณค่า

$|X'X|_{pre}$ เป็นดีเทอร์มิแนนท์ของตัวแปรอิสระที่ได้จากข้อมูลที่ใช้ในการพยากรณ์

p เป็นจำนวนสัมประสิทธิ์ถดถอยในสมการ

โดยค่าที่ประมาณได้จะเข้าใกล้ 1 ถ้าตัวแปรอิสระของข้อมูลทั้ง 2 ชุดมีลักษณะคล้ายคลึงกัน แต่ถ้าหากค่าที่ประมาณได้ต่างไปจาก 1 มาก แสดงว่าตัวแปรอิสระของข้อมูล 2 ชุดมีลักษณะที่แตกต่างกัน ไม่เหมาะสมที่จะนำข้อมูลทั้งหมดมาพิจารณาในการสร้างตัวแบบถดถอยเพื่อการพยากรณ์ เพราะจะทำให้เกิดความเสี่ยงสูง

2.4 การทดสอบความคงที่ของสัมประสิทธิ์ถดถอยเชิงเส้นโดยใช้การทดสอบเขา

จุดประสงค์ของการตรวจสอบความเหมาะสมของข้อมูลในตัวแบบถดถอยเชิงเส้นด้วยวิธีดูเพล็กซ์ก็คือ เพื่อพิจารณาว่าตัวประมาณที่ได้จากข้อมูลทั้งหมดสมควรจะนำไปใช้พยากรณ์หรือไม่ ซึ่งจะตรวจสอบโดยทำการแบ่งข้อมูลด้วยวิธีดูเพล็กซ์ หลังจากนั้นจะได้ข้อมูลชุดที่ใช้ในการประมาณค่าและชุดที่ใช้ในการพยากรณ์ จากข้อมูลทั้ง 2 ชุดนำมาหาตัวประมาณคือ

$$\hat{y}_e = X_e \hat{\beta}_e$$

$$\hat{y}_{pr} = X_{pr} \hat{\beta}_{pr}$$

เมื่อ \hat{y}_e	เป็นเวกเตอร์ของตัวประมาณที่ได้จากข้อมูลที่ใช้ในการประมาณค่า ขนาด $(m \times 1)$
X_e	เป็นเมตริกซ์ของตัวแปรอิสระของข้อมูลที่ใช้ในการประมาณค่า ขนาด $(m \times p)$
$\hat{\beta}_e$	เป็นสัมประสิทธิ์ถดถอยเชิงเส้นที่ได้จากข้อมูลที่ใช้ในการประมาณ ค่าขนาด $(p \times 1)$
\hat{y}_{pr}	เป็นเวกเตอร์ของตัวประมาณที่ได้จากข้อมูลที่ใช้ในการพยากรณ์ ขนาด $(l \times 1)$
X_{pr}	เป็นเมตริกซ์ของตัวแปรอิสระของข้อมูลที่ใช้ในการพยากรณ์ขนาด $(l \times p)$
$\hat{\beta}_{pr}$	เป็นสัมประสิทธิ์ถดถอยเชิงเส้นที่ได้จากข้อมูลที่ใช้ในการพยากรณ์ ขนาด $(p \times 1)$
m	เป็นจำนวนข้อมูลที่ใช้ในการประมาณค่า
l	เป็นจำนวนข้อมูลที่ใช้ในการพยากรณ์
และ $m + l = n$	หรือจำนวนข้อมูลทั้งหมด

ในการพิจารณาว่าตัวประมาณที่ได้จากข้อมูลทั้งหมดสมควรนำไปใช้ในการพยากรณ์หรือไม่จะพิจารณาจากความคงที่ของสัมประสิทธิ์ถดถอยเชิงเส้นในสัมภาระ โดยการทดสอบสมมติฐานที่ว่าสัมประสิทธิ์ถดถอยเชิงเส้นที่ประมาณจากข้อมูลที่ใช้ในการประมาณค่า ($\hat{\beta}_e$) และสัมประสิทธิ์ถดถอยเชิงเส้นที่ประมาณจากข้อมูลที่ใช้ในการพยากรณ์ ($\hat{\beta}_{pr}$) มีความแตกต่างกันอย่างมีนัยสำคัญหรือไม่ ได้มีผู้ที่ศึกษาเรื่องการทดสอบสมมติฐานเกี่ยวกับความเท่ากันของสัมประสิทธิ์ถดถอยเชิงเส้นจากตัวอย่าง 2 ชุดคือ เช่า (Chow 1960:591-605) โดยทำการศึกษาเกี่ยวกับความถดถอยเชิงเส้นที่แสดงถึงความสัมพันธ์ทางด้านเศรษฐกิจค่าสตรีใน เรื่องของฟังก์ชันการบริโภค (consumption function) ที่ขึ้นอยู่กับตัวแปรต่าง ๆ เช่น รายได้ ระดับราคา ความต้องการของผู้บริโภค และปริมาณการผลิต เป็นต้น เมื่อมีการใช้ความถดถอยเชิงเส้นในการแสดงความสัมพันธ์ทางด้านเศรษฐกิจค่าสตรีมากขึ้น จึงทำให้เกิดคำถามที่ว่าความสัมพันธ์ที่เกิดขึ้นจะ

ยังคงใช้ได้หรือไม่ในช่วงระยะเวลาที่ต่างกัน เขาได้ให้หลักสถิติในการตอบคำถามเหล่านี้ โดยการทดสอบตัวแบบถดถอยเชิงเส้นที่ได้จากกลุ่มของข้อมูล 2 ชุด ชุดแรกคือ ข้อมูลที่มีอยู่เดิม ส่วนชุดที่ 2 คือข้อมูลที่เก็บในระยะเวลามา จะแตกต่างกันอย่างมีนัยสำคัญหรือไม่ การทดสอบสมมติฐานจะทำให้ทราบว่าในช่วงระยะเวลาที่ต่างกันลักษณะหรือรูปแบบของความสัมพันธ์จะเปลี่ยนแปลงไปหรือไม่ เบิร์ค (Berk 1984 : 331-338) ได้แนะนำถึงการทดสอบเชื่อว่าเหมาะสมที่จะนำมาทดสอบความคงที่ของสัมประสิทธิ์ถดถอยเชิงเส้น ดังนั้นจึงทำให้เกิดแนวความคิดในการทดสอบค่าสัมประสิทธิ์ถดถอยเชิงเส้นที่ได้จากข้อมูลที่ใช้ในการประมาณค่าและค่าสัมประสิทธิ์ถดถอยเชิงเส้นที่ได้จากข้อมูลที่ใช้ในการพยากรณ์ว่าจะแตกต่างกันอย่างมีนัยสำคัญหรือไม่

ข้อตกลงเบื้องต้นของการทดสอบเขา

1. ตัวแปรอิสระแต่ละตัวเป็นค่าคงที่
2. ความคลาดเคลื่อนเป็นอิสระซึ่งกันและกัน และมีการแจกแจงแบบปกติที่มีค่าเฉลี่ยเป็น 0 และความแปรปรวน σ^2 ซึ่งการวิจัยในครั้งนี้จะฝ่าฝืนข้อตกลงนี้คือ ความคลาดเคลื่อนจะมีการแจกแจงแบบโลจิสติก ดับเบิลเอ็กซ์โพเนนเชียล และปกติปลอมปน
3. จำนวนข้อมูลที่ใช้ในการประมาณค่ามีมากกว่าจำนวนสัมประสิทธิ์ถดถอยเชิงเส้นในสมการ

การทดสอบในที่นี้จะพิจารณาว่าค่าสัมประสิทธิ์ถดถอยเชิงเส้นที่ได้จากข้อมูลที่ใช้ในการประมาณค่า และค่าสัมประสิทธิ์ถดถอยเชิงเส้นที่ได้จากข้อมูลที่ใช้ในการพยากรณ์ แตกต่างกันอย่างมีนัยสำคัญหรือไม่ โดยพิจารณาว่าข้อมูลที่ใช้ในการประมาณค่าคือข้อมูลเดิมที่มีอยู่ ส่วนข้อมูลที่ใช้ในการพยากรณ์คือข้อมูลที่เก็บในระยะเวลามา

สมมติฐานที่ใช้ในการทดสอบ

H_0 : สัมประสิทธิ์ถดถอยเชิงเส้นที่ได้จากข้อมูลทั้ง 2 ชุด มีค่าเท่ากัน

$$(\beta_e = \beta_{pr})$$

H_a : สัมประสิทธิ์ถดถอยเชิงเส้นที่ได้จากข้อมูลทั้ง 2 ชุด มีค่าแตกต่างกัน

$$(\beta_e \neq \beta_{pr}) \quad \text{สถิติที่ใช้ในการทดสอบคือ}$$

$$F = \frac{\|x_{e\hat{\beta}_e} - x_e\hat{\beta}_e\|^2 + \|x_{pr\hat{\beta}_{pr}} - x_{pr}\hat{\beta}_{pr}\|^2}{\|y_e - x_e\hat{\beta}_e\|^2 + \|y_{pr} - x_{pr}\hat{\beta}_{pr}\|^2} \cdot \frac{(m+l-2p)}{p}$$

หรือ

$$F = \frac{(A-B-C)}{p} / \frac{(B+C)}{(m+l-2p)}$$

เมื่อ

$$A = (y_e - x_e\hat{\beta}_e)' (y_e - x_e\hat{\beta}_e) \text{ หรือผลบวกกำลังสองของความคลาดเคลื่อนของข้อมูลทั้งหมดโดยใช้ค่าพารามิเตอร์ที่ประมาณจากข้อมูลทั้งหมด}$$

$$B = (y_e - x_e\hat{\beta}_e)' (y_e - x_e\hat{\beta}_e) \text{ หรือผลบวกกำลังสองของความคลาดเคลื่อนของข้อมูลที่ใช้ในการประมาณค่าโดยใช้ค่าพารามิเตอร์ที่ประมาณจากข้อมูลที่ใช้ในการประมาณค่า}$$

$$C = (y_{pr} - x_{pr}\hat{\beta}_{pr})' (y_{pr} - x_{pr}\hat{\beta}_{pr}) \text{ หรือผลบวกกำลังสองของความคลาดเคลื่อนของข้อมูลที่ใช้ในการพยากรณ์โดยใช้ค่าพารามิเตอร์ที่ประมาณจากข้อมูลที่ใช้ในการพยากรณ์}$$

$$p = \text{จำนวนสัมประสิทธิ์ถดถอยเชิงเส้นในสมการ}$$

$$m = \text{จำนวนข้อมูลที่ใช้ในการประมาณค่า}$$

$$l = \text{จำนวนข้อมูลที่ใช้ในการพยากรณ์}$$

$$m+l=n \text{ หรือจำนวนข้อมูลทั้งหมด}$$

เกณฑ์การตัดสินใจ

จะปฏิเสธสมมติฐาน H_0 เมื่อ $F > F(p, m+l-2p)$

เมื่อ $F(p, m+l-2p)$ เป็นค่าวิกฤตที่ได้จากตารางเอฟ (F table) ที่องศา

อิสระ p กับ $m+l-2p$ ณ ระดับนัยสำคัญ α

การพิสูจน์สูตรแสดงไว้ดังนี้

การทดสอบความเท่ากันของสัมประสิทธิ์ถดถอยเชิงเส้นของตัวอย่าง 2 ชุด โดยตัวแบบถดถอยเชิงเส้นของตัวอย่างทั้ง 2 ชุด มีรูปแบบดังนี้

$$\begin{aligned} \tilde{y}_1 &= x_1\beta_1 + \varepsilon_1 = x_1\beta_1 + 0\beta_2 + \varepsilon_1 \\ \tilde{y}_2 &= x_2\beta_2 + \varepsilon_2 = 0\beta_1 + x_2\beta_2 + \varepsilon_2 \end{aligned} \quad (1)$$

หรือ

$$\begin{bmatrix} \tilde{y}_1 \\ \tilde{y}_2 \end{bmatrix} = \begin{bmatrix} x_1 & 0 \\ 0 & x_2 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \end{bmatrix}$$

ภายใต้สมมติฐาน $H_0 : \beta_1 = \beta_2 = \beta$ จะได้ว่า

$$\begin{bmatrix} \tilde{y}_1 \\ \tilde{y}_2 \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \beta + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \end{bmatrix} \quad (2)$$

ถ้าสมมติฐาน $H_0 : \beta_1 = \beta_2 = \beta$ เป็นจริง ตัวประมาณที่ได้จากวิธีกำลังสอง

น้อยที่สุดของ β คือ

$$\begin{aligned} \hat{\beta}_0 &= \left[\begin{matrix} x_1' & x_2' \end{matrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \right]^{-1} \begin{bmatrix} x_1' & x_2' \end{bmatrix} \begin{bmatrix} \tilde{y}_1 \\ \tilde{y}_2 \end{bmatrix} \quad (3) \\ &= \left[\begin{matrix} x_1' & x_1 & x_2' & x_2 \end{matrix} \right]^{-1} \begin{bmatrix} x_1' & x_2' \end{bmatrix} \begin{bmatrix} \tilde{y}_1 \\ \tilde{y}_2 \end{bmatrix} = \beta + \left[\begin{matrix} x_1' & x_1 & x_2' & x_2 \end{matrix} \right]^{-1} \begin{bmatrix} x_1' & x_2' \end{bmatrix} \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \end{bmatrix} \end{aligned}$$

ดังนั้นความคลาดเคลื่อนที่เกิดขึ้นคือ

$$\begin{aligned} \begin{bmatrix} \tilde{y}_1 \\ \tilde{y}_2 \end{bmatrix} - \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \hat{\beta}_0 &= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \beta + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \end{bmatrix} - \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \beta - \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \left[\begin{matrix} x_1' & x_1 & x_2' & x_2 \end{matrix} \right]^{-1} \begin{bmatrix} x_1' & x_2' \end{bmatrix} \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \end{bmatrix} \quad (4) \\ &= \left[\begin{matrix} I & - \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \begin{pmatrix} x_1' & x_1 & x_2' & x_2 \end{pmatrix}^{-1} \begin{pmatrix} x_1' & x_2' \end{pmatrix} \end{matrix} \right] \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \end{bmatrix} \end{aligned}$$

ผลบวกกำลังสองของความคลาดเคลื่อนคือ

$$\begin{aligned} \left\| \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} - \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} b_0 \right\|^2 &= \left[\begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} - \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} b_0 \right]' \left[\begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} - \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} b_0 \right] \quad (5) \\ &= \begin{bmatrix} \varepsilon_1' & \varepsilon_2' \end{bmatrix} \begin{bmatrix} I - \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} (X_1' X_1 + X_2' X_2)^{-1} (X_1' X_2') \end{bmatrix} \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \end{bmatrix} \end{aligned}$$

เพราะเหตุว่าความคลาดเคลื่อนได้จากค่าสังเกต $m+l$ ค่า ซึ่งมีตัวแปร p ตัว

ดังนั้นอันดับ (rank) ของความคลาดเคลื่อนคือ $m+l - p$ (Kempthorne 1952:54-59)

ถ้าสมมติฐานแย้ง $H_a : \beta_1 \neq \beta_2$ เป็นจริง โดยพิจารณาสมการ (1) ตัวประมาณที่ได้จากวิธีกำลังสองน้อยที่สุดของ β_1 และ β_2 คือ

$$\begin{bmatrix} b_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} X_1' X_1 & 0 \\ 0 & X_2' X_2 \end{bmatrix}^{-1} \begin{bmatrix} X_1' & 0 \\ 0 & X_2' \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} = \begin{bmatrix} (X_1' X_1)^{-1} X_1' Y_1 \\ (X_2' X_2)^{-1} X_2' Y_2 \end{bmatrix} \quad (6)$$

ความคลาดเคลื่อนที่เกิดขึ้นภายใต้สมมติฐานแย้ง H_a คือ

$$\begin{bmatrix} Y_1 - X_1 b_1 \\ Y_2 - X_2 b_2 \end{bmatrix} = \begin{bmatrix} [I - X_1 (X_1' X_1)^{-1} X_1'] \varepsilon_1 \\ [I - X_2 (X_2' X_2)^{-1} X_2'] \varepsilon_2 \end{bmatrix} \quad (7)$$

ผลบวกกำลังสองของความคลาดเคลื่อนคือ

$$\begin{aligned} \left\| \begin{bmatrix} Y_1 - X_1 b_1 \\ Y_2 - X_2 b_2 \end{bmatrix} \right\|^2 &= \left\| Y_1 - X_1 b_1 \right\|^2 + \left\| Y_2 - X_2 b_2 \right\|^2 \quad (8) \\ &= \varepsilon_1' [I - X_1 (X_1' X_1)^{-1} X_1'] \varepsilon_1 + \varepsilon_2' [I - X_2 (X_2' X_2)^{-1} X_2'] \varepsilon_2 \end{aligned}$$

อันดับของเทอมทางขวามือคือ $m-p$ และ $l-p$ ตามลำดับและเนื่องจาก ξ_1 และ ξ_2 เป็นอิสระต่อกัน ดังนั้นอันดับของเทอมทางซ้ายมือคือ $m+l-2p$

จากสมการ (5) ภายใต้สมมติฐาน H_0 จะเขียนให้อยู่ในรูปของสมการ (8) ภายใต้สมมติฐานแย้ง H_a บวกกับผลบวกกำลังสองของผลต่าง

$$\begin{bmatrix} x_1 b_1 - x_1 b_0 \\ x_2 b_2 - x_2 b_0 \end{bmatrix} \text{ และ } \begin{bmatrix} x_1 b_1 - x_1 b_0 \\ x_2 b_2 - x_2 b_0 \end{bmatrix}$$

ดังนั้นสมการที่ได้คือ

$$\begin{bmatrix} x_1 - x_1 b_0 \\ x_2 - x_2 b_0 \end{bmatrix} = \begin{bmatrix} y_1 - x_1 b_1 \\ y_2 - x_2 b_2 \end{bmatrix} + \begin{bmatrix} x_1 b_1 - x_1 b_0 \\ x_2 b_2 - x_2 b_0 \end{bmatrix} \quad (9)$$

ผลบวกกำลังสองของสมาชิกในแต่ละเทอมของสมการ (9) คือ

$$\left\| \begin{bmatrix} x_1 - x_1 b_0 \\ x_2 - x_2 b_0 \end{bmatrix} \right\|^2 = \left\| \begin{bmatrix} y_1 - x_1 b_1 \\ y_2 - x_2 b_2 \end{bmatrix} \right\|^2 + \left\| \begin{bmatrix} x_1 b_1 - x_1 b_0 \\ x_2 b_2 - x_2 b_0 \end{bmatrix} \right\|^2 \quad (10)$$

หรือ

$$Q_1 = Q_2 + Q_3 \quad (11)$$

จะแสดงว่าอันดับของ Q_3 มีค่ามากที่สุดคือ p , จากสมการ (3) และ (6) จะได้ว่า

$$\begin{bmatrix} x_1' x_1 + x_2' x_2 \end{bmatrix} b_0 = x_1' y_1 + x_2' y_2 = x_1' x_1 b_1 + x_2' x_2 b_2 \quad (12)$$

หรือ

$$b_2 - b_0 = - (x_2' x_2)^{-1} x_1' x_1 (b_1 - b_0) \quad (13)$$

นำสมการ (13) ไปใช้ใน Q_3 จะได้เป็น

$$Q_3 = \left\| \begin{bmatrix} x_1 (b_1 - b_0) \\ -x_2 (x_2' x_2)^{-1} x_1' x_1 (b_1 - b_0) \end{bmatrix} \right\|^2 \quad (14)$$

$$= \begin{bmatrix} b'_1 - b'_0 \\ b'_1 - b'_0 \end{bmatrix} \begin{bmatrix} x'_1 - x'_1 x_1 (x'_2 x_2)^{-1} x'_2 \\ -x_2 (x'_2 x_2)^{-1} x'_1 x_1 \end{bmatrix} \begin{bmatrix} x_1 \\ b_1 - b_0 \end{bmatrix}$$

สมการ (14) คือรูปแบบกำลังสอง (quadratic form) ใน $b_1 - b_0$ ซึ่งอันดับจะมีค่าไม่เกิน p และจากสมการ (3) $b_1 - b_0$ คือการแปลงเชิงเส้น (linear transformation) ของความคลาดเคลื่อน (ε) ดังนั้นจะได้ว่า

$$b_1 - b_0 = \beta_1 - \beta_2 \left\{ \begin{bmatrix} (x'_1 x_1)^{-1} x'_1 0 \\ -[x'_1 x_1 + x'_2 x_2]^{-1} [x'_1 x_2] \end{bmatrix} \right\} \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \end{bmatrix} \quad (15)$$

ภายใต้สมมติฐาน $H_0 : \beta_1 = \beta_2 = \beta$, Q_3 ก็คือรูปแบบกำลังสองของความคลาดเคลื่อน ซึ่งมีค่าอันดับไม่เกิน p จากสมการ (15) จะเห็นได้ว่า Q_3 จะมีค่ามากถ้า สมมติฐาน H_0 เป็นเท็จ

เนื่องจากอันดับของ Q_2 คือ $m+l-2p$ อันดับของ Q_3 คือ p และอันดับของ Q_1 มีค่าน้อยกว่าหรือเท่ากับ อันดับของ Q_2 บวกอันดับของ Q_3 ดังนั้นภายใต้สมมติฐาน H_0 , Q_2 และ Q_3 จะมีการแจกแจงอย่างอิสระเป็น $\chi^2 (m+l-2p) \sigma^2$ และ $\chi^2 (p) \sigma^2$ ซึ่งเราสามารถทดสอบสมมติฐาน H_0 โดยใช้ค่าสัดส่วนเอฟ (F-ratio)

$$F(p, m+l-2p) = \frac{Q_3 / p}{Q_2 / (m+l-2p)}$$

$$= \frac{\|x_1 b_1 - x_1 b_0\|^2 + \|x_2 b_2 - x_2 b_0\|^2}{\|x_1 - x_1 b_1\|^2 + \|x_2 - x_2 b_2\|^2} \cdot \frac{(m+l-2p)}{p}$$

2.5 ผลงานวิจัยที่เกี่ยวข้อง

ในการศึกษาเรื่องความเหมาะสมของตัวแบบถดถอยเชิงเส้นนั้น มีผู้ศึกษาไว้มากนัก ดังนั้นผลงานวิจัยที่เกี่ยวข้อง จึงมีอยู่น้อย แต่ก็ยังมีนักวิจัยบางท่านที่ได้ศึกษาเกี่ยวกับวิธีต่าง ๆ ที่ใช้ในการตรวจสอบความเหมาะสมของตัวแบบถดถอยเชิงเส้น ซึ่งวิธีที่น่าสนใจที่แนะนำเล่นพร้อมทั้งข้อสรุปต่าง ๆ มีดังนี้คือ

2.5.1 ตรวจสอบตัวแบบที่ใช้ในการพยากรณ์และสัมประสิทธิ์ของตัวแบบ (Check on Model Prediction and Coefficients) ในบางครั้งค่าสัมประสิทธิ์ถดถอยเชิงเส้นที่ประมาณได้จากข้อมูล ($\hat{\beta}$) โดยวิธีกำลังสองน้อยที่สุดอาจไม่ดีเท่าที่ควร จากการศึกษาของโรนัลด์ ดี. สนิ พบว่าถ้าตัวแปรอิสระมีความสัมพันธ์ต่อกัน (multicollinearity) ควรใช้วิธีรีดจ์รีเกรสชัน (Ridge Regression) ในการประมาณค่าสัมประสิทธิ์ถดถอยเชิงเส้นแล้วจึงนำตัวแบบที่ได้มาตรวจสอบความเหมาะสม

2.5.2 การตรวจสอบความเหมาะสมของตัวแบบถดถอยด้วยข้อมูลใหม่ (Validating Regression Model with New Data) โดยใช้ข้อมูลที่มีอยู่เป็นข้อมูลที่ใช้ในการประมาณค่าและใช้ข้อมูลใหม่ (new data) เป็นข้อมูลที่ใช้ในการพยากรณ์ จากการศึกษาของเบิร์คพบว่าในการตรวจสอบความเหมาะสมของตัวแบบถดถอยสิ่งที่สำคัญอีกอย่างหนึ่งก็คือ วิธีการประมาณค่าสัมประสิทธิ์ถดถอย ควรสอดคล้องกับลักษณะของข้อมูลคือ ถ้าตัวแปรอิสระแต่ละตัวมีความสัมพันธ์กับตัวแปรตามในลักษณะเดียวกัน วิธีการประมาณค่าพารามิเตอร์ของลินด์ลีย์-สมิท เบย์ส์ (Lindley-Smith Bayes) จะให้ผลดี ส่วนวิธีการคัดเลือกตัวแปร (Subset Selection) จะให้ผลดีถ้าตัวแปรอิสระมีไม่มาก และทราบว่าตัวแปรอิสระตัวใดมีความสัมพันธ์กับตัวแปรตามอย่างเห็นได้ชัด แต่อย่างไรก็ตามวิธีการตรวจสอบความเหมาะสมของตัวแบบโดยวิธีนี้ก็ยังมีข้อบกพร่องดังที่ได้กล่าวมาแล้วในบทที่ 1