

การวิเคราะห์การถดถอย

4.1 การวิเคราะห์การถดถอยกับงานวิจัย

เนื่องจากงานวิจัยนี้มีจุดประสงค์อย่างหนึ่ง คือ ใช้การถดถอยพหุคูณ (Multiple regression) วิเคราะห์ข้อมูลที่เก็บจากสถานที่ทำงานทั้ง 4 แห่ง เพื่อหาสมการโหนดเฉลี่ยของคนไทย (ในที่นี้เป็นตัวแปรตาม) ว่ามีตัวแปรอิสระอะไรบ้างที่มีนัยสำคัญ (Significance)

ตัวแปรตามและตัวแปรอิสระที่ใช้ในการวิเคราะห์ในงานวิจัยมีดังนี้
ตัวแปรตาม ได้แก่ PMV หรือ ดัชนีทำนายการโหนดเฉลี่ย
ตัวแปรอิสระมี 6 ตัว ได้แก่

- 1 ลักษณะงาน
- 2 ความต้านทานความร้อนของเสื้อผ้า หรือชุดแต่งกาย
- 3 อุณหภูมิอากาศในห้องทดสอบ
- 4 อุณหภูมิการแผ่ความร้อนเฉลี่ย
- 5 ความเร็วอากาศสัมพัทธ์
- 6 ความดันไอน้ำในอากาศ

ดังนั้นจะบอกกล่าวถึงทฤษฎีการถดถอยอย่างง่ายและพหุคูณ ดังหัวข้อต่อไปนี้

4.2 การถดถอยและสหสัมพันธ์ (Regression and correlation)

หมายถึง วิธีการอย่างหนึ่งที่ใช้ในการศึกษาตัวแปรหนึ่งเมื่อมีความสัมพันธ์ หรือได้รับอิทธิพลจากตัวแปรอื่นๆ ตัวแปรที่เกี่ยวข้องกับตัวแปรที่ต้องการศึกษาอาจมีตัวเดียว การวิเคราะห์ลักษณะนี้เรียกว่า การวิเคราะห์การถดถอยและสหสัมพันธ์อย่างง่าย (Simple regression and correlation analysis) แต่ในบางกรณีตัวแปรที่เกี่ยวข้องกับตัวแปรที่

ต้องการศึกษาจากมีตั้งแต่ 2 ตัวขึ้นไป การวิเคราะห์การถดถอย ลักษณะนี้ เรียกว่า การถดถอยพหุคูณ

ไม่ว่าจะเป็นการวิเคราะห์การถดถอยอย่างง่าย หรือ การถดถอยพหุคูณ วิธีการที่ใช้ คล้ายคลึงกัน เพียงแต่การวิเคราะห์การถดถอยพหุคูณ วิธีการจะยุ่งยากขึ้นเนื่องจากตัวแปรที่เกี่ยวข้องมีมากกว่า 1 ตัว ยิ่งตัวแปรที่เกี่ยวข้องมีมากเท่าใด ความยุ่งยากก็จะเพิ่มมากขึ้นเท่านั้น

4.3 ข้อตกลงเบื้องต้นของการวิเคราะห์การถดถอยเชิงเส้นอย่างง่าย

การวิเคราะห์การถดถอยอย่างง่ายมีข้อตกลงเบื้องต้นดังนี้

4.3.1 ตัวแปรอิสระ X และตัวแปรตาม Y จะต้องมีความสัมพันธ์ในแบบเส้นตรง

$$Y = a + bX + e$$

เมื่อ e คือค่าคลาดเคลื่อนสุ่มที่เกิดขึ้นในค่าที่สังเกตได้ของตัวแปร Y อันเนื่องมาจากตัวอย่างที่สุ่มมา

4.3.2 ตัวแปรตามต้องเป็นตัวแปรสุ่มชนิดต่อเนื่อง ในขณะที่ตัวแปรอิสระเป็นเซตของค่าต่าง ๆ ที่เลือกกำหนดได้

4.3.3 ความแปรปรวนของ Y สำหรับแต่ละค่าของ X ที่กำหนดจะต้องมีค่าเท่ากัน หรือเรียกว่า Homoscedasticity สามารถพิจารณาอย่างคร่าว ๆ ว่าข้อมูลที่ได้มาสอดคล้องกับคุณสมบัติข้อนี้หรือไม่อาจทำได้โดยการพิจารณาจากแผนภาพการกระจายของข้อมูลจากตัวอย่าง ถ้าการกระจายของจุดไปจากเส้นแนวโน้มมีพอ ๆ กัน สำหรับค่าต่าง ๆ ของ X ที่กำหนด แสดงว่าข้อมูลที่ได้มาสอดคล้องกับคุณสมบัติข้อนี้

4.3.4 ค่าที่สังเกตได้แต่ละค่าต้องไม่มีความสัมพันธ์กัน ถ้าเป็นตัวอย่างที่สุ่มมากก็จะสอดคล้องกับคุณสมบัติข้อนี้

4.4 การถดถอยเชิงเส้นอย่างง่าย (Simple linear regression)

มีสมการพื้นฐานคือ

$$Y' = a + bX \quad (4.1)$$

- เมื่อ Y' คือ ค่าทำนายตัวแปรตาม Y
 X คือ ค่าตัวแปรอิสระ
 a คือ ค่าคงที่ Intercept
 b คือ สัมประสิทธิ์การถดถอย

โดยที่ ค่า a หาได้จากสูตร

$$a = \bar{Y} - b\bar{X} \quad (4.2)$$

และ b หาได้จากสูตร

$$b = \frac{\sum xy}{\sum x^2} \quad (4.3)$$

เมื่อ x, y คือ ค่าเบี่ยงเบนจากค่าเฉลี่ยของ X, Y หรือ $x = X - \bar{X}$ $y = Y - \bar{Y}$ และ
 $\sum x^2 = \sum X^2 - (\sum X)^2 / N$, $\sum y^2 = \sum Y^2 - (\sum Y)^2 / N$, $\sum xy = \sum XY - (\sum X)(\sum Y) / N$
 และ N คือ จำนวนคู่ข้อมูล

ส่วนวิธีการลากเส้นถดถอย มีกฎอยู่ว่า ต้องลากผ่านค่าคงที่ Intercept และจุด X, Y

4.5 ความหมายและการหาค่าสัมประสิทธิ์สหสัมพันธ์ (r)

ในการวิเคราะห์การถดถอย จะมีสัมประสิทธิ์ตัวหนึ่งที่เป็นดัชนี แสดงขนาดและทิศ
 ทางของความสัมพันธ์ระหว่างตัวแปร ดัชนีตัวนั้นได้แก่ สัมประสิทธิ์สหสัมพันธ์ (r) ซึ่งค่า r มีค่า
 อยู่ระหว่าง $-1 < r < +1$ ถ้า $r = -1.0$ เมื่อนำมาพล็อตกราฟ ค่า X, Y ทุกค่าจะอยู่บนเส้น
 ถดถอย แต่ถ้า r เข้าใกล้ 0 จุด X, Y จะยิ่งกระจายห่างจากเส้นถดถอยมากขึ้น เครื่องหมาย
 + หรือ - ติดอยู่ แสดงว่า slope ของเส้นถดถอยเป็น + หรือ - ส่วนการหาค่าสัมประสิทธิ์ตัวนี้
 มีหลายวิธี แต่ละวิธีจะให้ค่าเท่ากันเช่น

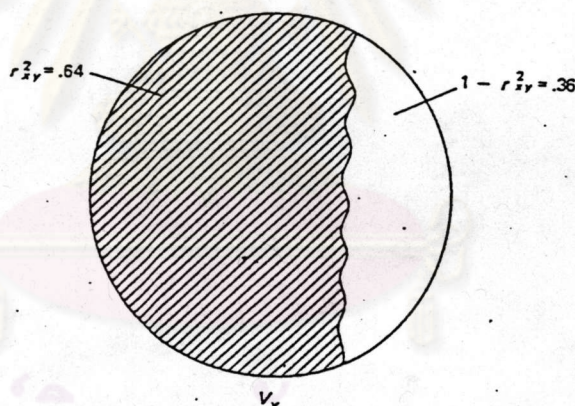
$$r_{xy} = \frac{\sum xy}{(\sum x^2 \sum y^2)^{0.5}} \quad (4.4)$$

$$r_{xy} = \frac{\sum XY}{N s_x s_y} \quad (4.5)$$

เมื่อ s_x คือ ค่าเบี่ยงเบนมาตรฐานของ X

4.6 ความหมายและการหาค่าสัมประสิทธิ์การตัดสินใจ (r^2)

เป็นสัมประสิทธิ์ที่จะบอกให้ทราบถึงความแปรปรวนของ Y ที่สามารถอธิบายได้ด้วยตัวแปร X ได้มากน้อยเท่าไร เช่น



รูปที่ 4.1 สัดส่วนความแปรปรวนของ Y

จากรูปที่ 4.1 พื้นที่แรเงาคือสัดส่วนของความแปรปรวน Y ที่สามารถอธิบายได้ด้วยตัวแปร $X = 0.64$ ไม่สามารถอธิบายด้วยตัวแปร $X = 0.36$ เมื่อพื้นที่วงกลม (V_y) คือ ความแปรปรวนทั้งหมดของ $Y = 1$

4.7 การทดสอบนัยสำคัญของการถดถอยอย่างง่าย

เนื่องจากเส้นถดถอยสำหรับชุดข้อมูลแต่ละชุดเป็นเพียงเส้นหนึ่งที่ถูกเลือกขึ้นมา ดังนั้นความคลาดเคลื่อนที่เกิดขึ้นจึงเป็นแบบสุ่ม จึงจำเป็นต้องมีการทดสอบนัยสำคัญของสัมประสิทธิ์ถดถอย

โดยการตั้งสมมติฐาน เมื่อนำมาใช้กับการวิเคราะห์การถดถอยซึ่งมีรูปสมการดังนี้

$$SS_t = SS_{reg} + SS_{res} \quad (4.6)$$

เมื่อ SS_t คือ ผลรวมทั้งหมดกำลังสองของ Y

$$= \sum y_t^2 = \sum Y^2 - (\sum Y)^2/N$$

SS_{reg} คือ ผลรวมกำลังสองของ Y เนื่องจากการถดถอย

$$= \sum y'^2 = \sum Y'^2 - (\sum Y')^2/N$$

SS_{res} คือ ผลรวมกำลังสองของ Residuals

$$= \sum d^2 = \sum (Y - Y')^2$$

และการทดสอบหาค่าสำคัญจะใช้ F-ratio ดังสูตร

$$F = (SS_{reg}/df_1) / (SS_{res}/df_2) \quad (4.7)$$

เมื่อ df_1 คือ องศาอิสระ (degree of freedom) ของการถดถอย (=K)

df_2 คือ องศาอิสระของ residuals (=N-K-1)

K คือ จำนวนตัวแปรอิสระ

เมื่อหาค่า F ได้แล้วจะเปิดตาราง F-distribution ในภาคผนวกเพื่อตรวจสอบดูว่ามีนัยสำคัญหรือไม่

4.8 การวิเคราะห์การถดถอยพหุคูณ (Multiple regression analysis)

วิธีวิเคราะห์ที่จะกล่าวต่อไปนี้ เป็น General method ของการวิเคราะห์ ที่มีตัวแปรอิสระที่ตัวก็ได้ สามารถใช้ได้กับการถดถอยทุกรูปแบบไม่ว่าตัวแปรจะเป็นแบบ Continuous, Category หรือ Codes

4.9 การวิเคราะห์สมการที่มีตัวแปรอิสระ 3 ตัวขึ้นไป

ในกรณีที่มีตัวแปรอิสระ 3 ตัว จำเป็นต้องใช้ Matrix หา สปส.การถดถอย (แต่ถ้ามีตัวแปรอิสระ 2 ตัว ไม่จำเป็นต้องใช้ matrix ก็สามารถใช้ Solve ปัญหาได้ ในที่นี้จะไม่บอกกล่าวถึง) รูปแบบของสมการการถดถอยพหุคูณทั่วไปเป็นดังนี้

$$Y' = a + b_1X_1 + b_2X_2 + \dots + b_kX_k \quad (4.8)$$

ค่า a, b_1, \dots, b_k เป็นค่าคงที่ที่ต้องหาเพื่อให้สมการมี residuals น้อยที่สุด

4.10 ข้อสมมติฐานในการวิเคราะห์การถดถอยพหุคูณ

ในการวิเคราะห์การถดถอยพหุคูณเชิงเส้นจะมีสมมติฐานดังนี้

4.10.1 ตัวแปรอิสระและตัวแปรตามมีความสัมพันธ์กันในแบบเส้นตรง สมการที่มีตัวแปรอิสระ k ตัว จะอยู่ในรูป

$$Y = a + b_1X_1 + b_2X_2 + \dots + b_kX_k + e \quad (4.9)$$

เมื่อ e คือ ความคลาดเคลื่อนสุ่มซึ่งเกิดจากการสุ่มตัวอย่าง

4.10.2 ตัวแปรตามเป็นตัวแปรสุ่มต่อเนื่อง ในขณะที่ตัวแปรอิสระเป็นตัวแปรควบคุมและเป็นเชิงของค่าซึ่งไม่ใช่ค่าสุ่ม ดังนั้น ค่า e ในรูปแบบเชิงเส้นจึงเป็นความคลาดเคลื่อนของตัวแปรตามอันเกิดจากการสุ่มตัวอย่าง การวิเคราะห์การถดถอยพหุคูณสามารถนำมาใช้ได้แม้ว่าตัวแปรอิสระจะเป็นตัวแปรสุ่มเช่นเดียวกัน

4.10.3 ความแปรปรวนของตัวตามเมื่อกำหนดค่าต่าง ๆ ของตัวแปรอิสระจะมีค่าเท่ากันหรือเรียกว่ามีคุณสมบัติของ Homoscedasticity

4.10.4 ค่าของตัวแปรสุ่มตัวที่อยู่ติดกัน ไม่สัมพันธ์กัน ถ้าค่าของตัวแปรสุ่มไม่สอดคล้องกับคุณสมบัติข้อนี้ เรียกว่า ข้อมูลมี Autocorrelation กัน

4.10.5 ตัวแปรอิสระแต่ละตัวจะต้องไม่มีความสัมพันธ์กัน

ถ้าจะกล่าวกันโดยแท้จริงแล้ว ข้อสมมติฐานทั้งห้า เป็นสิ่งที่จะสมบูรณ์ได้ยาก หากข้อสมมติฐานดังกล่าวเป็นสิ่งที่ใกล้เคียงกับความจริง แบบแผนของความสัมพันธ์ก็จะเป็สิ่งที่มิใช่ประโยชน์มาก

4.11 Solution ทั่วไปของสมการการถดถอย

จากสมการ (4.8) เราต้องมีสมการ k สมการในการหาค่า b

จาก Calculus จะได้สมการปกติ (Normal equation) สำหรับตัวแปรอิสระ 3 ตัว ดังนี้

$$\begin{aligned} \beta_1 + r_{12} \beta_2 + r_{13} \beta_3 &= r_{y1} \\ r_{21} \beta_1 + \beta_2 + r_{23} \beta_3 &= r_{y2} \\ r_{31} \beta_1 + r_{32} \beta_2 + \beta_3 &= r_{y3} \end{aligned} \quad (4.10)$$

หมายเหตุ

1. $r_{12} = r_{21}$, $r_{13} = r_{31}$ ฯลฯ และ r_{11} , r_{22} , r_{33} จะไม่เขียนในสมการเพราะมีค่าเท่ากับ 1
2. β คือ สัมประสิทธิ์การถดถอย population

จาก (4.10) สามารถเขียนเป็น matrix ได้ดังนี้

$$\begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} = \begin{bmatrix} r_{y1} \\ r_{y2} \\ r_{y3} \end{bmatrix} \quad (4.11)$$

หรือ

$$[R_{ij}] [\beta_j] = [R_{yj}]$$

เป็นที่ทราบแล้วว่า matrices สามารถ บวก, ลบ, คูณ ได้ แต่ไม่สามารถนำมาหารได้แต่จะใช้ inverse matrix แทน

จาก (4.11) จะได้

$$[\beta_j] = [R_{ij}]^{-1} [R_{yj}] \quad (4.12)$$

เมื่อ $[R_{ij}]^{-1}$ คือ inverse ของ $[R_{ij}]$

จาก (4.12) จะหาค่า b จากสูตร

$$b_j = \beta_j (s_y / s_j) \quad (4.13)$$

เมื่อ b คือ regression weight ที่ $j=1,2,3 \dots$
 s_y คือ ค่าเบี่ยงเบนมาตรฐานของตัวแปรตาม Y
 s_j คือ ค่าเบี่ยงเบนมาตรฐานของตัวแปรอิสระ

4.12 ความหมายและการหาค่าสัมประสิทธิ์การตัดสินใจพหุคูณ (R^2)

คือสัมประสิทธิ์ที่บอกถึงสัดส่วนของความแปรปรวนในตัวแปรตาม Y ที่อธิบายได้ด้วยตัวแปรอิสระต่าง ๆ ในเส้นถดถอย หาได้จากสูตร

$$R^2 = SS_{reg} / SS_t \quad (4.14)$$

เมื่อ SS_{reg} หาได้จาก Σy^2
 SS_{res} หาได้จาก Σd^2

4.13 ความหมายและการหาค่าสัมประสิทธิ์สหสัมพันธ์พหุคูณ (R)

คือ สัมประสิทธิ์สหสัมพันธ์เชิงพหุคูณที่อยู่ในรูปสัมบูรณ์หรือไม่มีเครื่องหมายกำกับ เพราะอาจเป็นไปได้ที่ตัวแปรตัวหนึ่งมีความสัมพันธ์ในทางตรงกับตัวแปรตามในขณะที่ตัวแปรอื่น ๆ มีความสัมพันธ์แบบผกผันกับตัวแปรตาม ดังนั้นความสัมพันธ์ระหว่างตัวแปรทั้งหมดจึงไม่สามารถบอกได้ชัดเจนว่าเป็นไปในทางบวกหรือลบ จะหาได้จากสูตร

$$R = (R^2)^{0.5} \quad (4.15)$$

เช่น $R^2 = 0.66$
 จะได้ $R = 0.81$

4.14 การทดสอบนัยสำคัญของการถดถอยพหุคูณ

หลังจากตั้งสมมติฐานแล้วจะใช้การทดสอบ F-ratio ตรวจสอบการมีนัยสำคัญ โดยใช้สูตร

$$F = \frac{R^2/K}{(1-R^2)/(N-K-1)} \quad (4.16)$$

เมื่อ K = จำนวนตัวแปรอิสระ

เมื่อกำหนดค่า F ได้แล้วจะเปิดตาราง F-distribution ในภาคผนวกเพื่อตรวจสอบว่า มีนัยสำคัญหรือไม่

4.15 การทดสอบนัยสำคัญทางสถิติของตัวแปรที่เพิ่มเข้ามาในสมการถดถอย

จากหัวข้อ 4.14 เรายังไม่ทราบว่า X ตัวไหนมีอิทธิพลต่อ Y มากหรือน้อย (เนื่องจากมี X หลายตัว) ดังนั้น จึงต้องมีการทดสอบว่า X ใดมีนัยสำคัญบ้าง โดยลดตัวแปรอิสระลงทีละตัว แล้วตรวจสอบนัยสำคัญของตัวแปรนั้นว่ามีหรือไม่ ถ้าไม่มีแสดงว่าเราสามารถตัดตัวแปรนั้นออกจากสมการได้

หมายเหตุ เนื่องจากการวิเคราะห์การถดถอย เป็นวิธีการทางสถิติที่นำมาใช้ในงานวิจัยนี้ เพื่อให้มีความเข้าใจได้ง่ายขึ้น โปรดดูตัวอย่างการวิเคราะห์ ในภาคผนวก ข.1