

ศัพท์พัฒนาบริหารศาสตร์

DEVELOPMENT ADMINISTRATION GLOSSARY

ตัวแบบเชิงเส้นของการถดถอยและสหสัมพันธ์

Linear Regression and Correlation : A Note

ศัพท์นี้มีจุดมุ่งหมายที่จะอธิบายความหมายของตัวแบบเชิงเส้นของการถดถอยเชิงเดี่ยว (Linear Model of Simple Regression) และความหมายของสหสัมพันธ์ (Correlation) เพื่อแสดงความแตกต่างระหว่างสองอย่างนี้ และกล่าวถึงข้อจำกัด (limitations) ของการใช้ในการวิเคราะห์ การใช้เทคนิคทางสถิติทั้งสองอย่างนี้มีข้อผิดพลาดอยู่บ่อย ๆ เช่นนำไปใช้กับกรณีที่ไม่เหมาะสม และบางทีผู้ใช้ก็มีความอย่างผิด ๆ บทความนี้จึงมีเจตนาที่จะเน้นให้ผู้ใช้หลีกเลี่ยงข้อผิดพลาดดังกล่าว

ตัวแบบเชิงเส้นของสมการถดถอยเชิงเดี่ยว (Linear Model of Simple Regression)

ตัวแบบเชิงเส้นของการถดถอยเชิงเดี่ยวที่นักสถิติใช้อยู่เสมอ ๆ และที่นักสังคมศาสตร์ชอบใช้อยู่บ่อย ๆ อยู่ในรูปของ

$$Y_i = a + bX_i + u_j \quad (i = 1, 2, \dots, T) \quad (1)$$

โดยที่

Y_i = ตัวแปรตาม (dependent variable)

X_i = ตัวแปรอิสระ (independent variable)

u_j = ตัวคลาดเคลื่อน (error)

a, b = ตัวสัมประสิทธิ์ที่ต้องการหาค่า (unknown coefficients)

i = subscripts มีค่าเท่ากับ 1, 2, 3, ..., T ซึ่งเป็นหน่วยแจงนับในตัวอย่าง¹

ตัวแบบข้างบนนี้มีความหมายว่าตัวแปร X และตัวคลาดเคลื่อน (u) เป็น "มูลเหตุ" (causes) ที่ทำให้เกิดค่าของตัวแปร Y ² ตัวอย่าง เช่น ให้ Y แทนจำนวนผลผลิตของข้าวโพด

1 จากนั้นไปเราจะใช้ X_i แทนค่าของตัวแปร และใช้ X แทนตัวแปร สำหรับตัวอื่น ๆ ก็มีความหมายในทำนองเดียวกัน

2 E. Malinvaud, *Statistical Methods of Econometrics* (Amsterdam: North - Holland Publishing Co., 1966), p. 74

ต่อไร่ X แทนจำนวนปุ๋ยที่ใช้ต่อไร่ ในระยะเวลาเพาะปลูก เราก็อาจตีความหมายตัวแบบนี้ได้ว่า จำนวนปุ๋ยที่ใช้เป็น “สาเหตุ” ที่ทำให้เกิดผลผลิตในปริมาณต่าง ๆ กัน ในที่นี้เรากล่าวถึงปุ๋ย แต่เพียงปัจจัยเดียว (ทั้งนี้เพื่อความง่ายในการคำนวณ หรือเราอาจจะมีความสนใจแต่ความสัมพันธ์ระหว่างปุ๋ยและผลผลิตเท่านั้น) ความจริงแล้วปริมาณผลผลิตของข้าวโพดยังขึ้นอยู่กับ น้ำ อากาศ และปัจจัยอื่นๆ อีก ทั้งที่สามารถวัดค่าได้ (quantifiable) และที่ไม่สามารถวัดค่าได้ (unquantifiable) เพราะฉะนั้นเราจึงต้องพิจารณาสิ่งเหล่านี้ด้วย ตัวคลาดเคลื่อน (u) นี้แหละที่แทนปัจจัยที่เราตัดทิ้งไปจากสมการ นอกจากนี้เราพบว่า X และ Y มีความสัมพันธ์กันอย่างเชิงเส้น แต่ความสัมพันธ์ที่แท้จริงอาจเป็นอย่างอื่น เช่น ในรูปของ quadratic function หรือ polynomial function หรือ exponential function ก็ได้ ในตัวอย่างข้างบนจะเห็นได้ว่า ปริมาณผลผลิตของข้าวโพดจะมีทางเพิ่มขึ้นได้น้อยที่จะมีความสัมพันธ์อย่างเชิงเส้นกับปริมาณปุ๋ยที่ใช้ ทั้งนี้เพราะว่าเมื่อเราเพิ่มปริมาณปุ๋ยต่อไร่ ผลผลิตจะเพิ่มขึ้น แต่จะถึงระยะหนึ่งเมื่อเพิ่มปุ๋ยมากขึ้นจะไม่ทำให้ผลผลิตเพิ่มขึ้นเลยหรืออาจทำให้ผลผลิตลดลงเพราะใช้ปุ๋ยปริมาณมากเกินไป ปริมาณของปุ๋ยและผลผลิตจะไม่เพิ่มด้วยกันเมื่อพ้นระยะหนึ่งไป เพราะฉะนั้นตัวคลาดเคลื่อน (u) จึงอาจใช้แทนความแตกต่างระหว่างตัวแบบที่ใช้กับตัวแบบที่แท้จริงได้อีกอย่างหนึ่ง

วิธีที่ใช้ประมาณค่าของ a และ b ในตัวแบบ $Y_i = a + bX_i + u_i$ ($i = 1, 2, \dots, T$) มีหลายวิธีด้วยกัน แต่วิธีที่รู้จักกันดีที่สุดคือ วิธีกำลังสองน้อยที่สุด (The Method of Least Squares) ซึ่งนักคณิตศาสตร์ชาวฝรั่งเศสชื่อ Adrain Legendre เป็นผู้เสนอ ตั้งแต่ปี ค.ศ. 1806 วิธีนี้เป็นที่นิยมใช้กันมาก เพราะเป็นวิธีที่ง่ายและถ้าสมมติฐานต่าง ๆ ของตัวแบบนี้เป็นจริง ตัวประมาณโดยวิธีกำลังสองน้อยที่สุด (least squares estimators) ก็จะมีคุณสมบัติเป็น Best Linear Unbiased Estimators³

หลักการของวิธีกำลังสองน้อยที่สุดมีอยู่ว่า จาก $Y_i = a + bX_i$ เราประมาณค่าด้วย $Y_i = \hat{a} + \hat{b}X_i$ (\hat{a} และ \hat{b} เป็นตัวประมาณของ a และ b ตามลำดับ) ความแตกต่างระหว่างค่าจริงและค่าประมาณของ Y (เรียกว่า residual) จะเท่ากับ $e_i = Y_i - \hat{Y}_i$ วิธีกำลังสองน้อยที่สุดต้องการทำให้

$$\sum_{i=1}^T e_i^2 = \sum_{i=1}^T (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^T (Y_i - \hat{a} - \hat{b}X_i)^2 \text{ น้อยที่สุด}$$

³ Ibid., ch. 3

ซึ่งเราสามารถแสดงว่าค่าประมาณของ a และ b จะได้ตามสูตร

$$\hat{a} = \bar{Y} - b\bar{X} \quad (2)$$

$$\hat{b} = \frac{\sum_{i=1}^T x_i y_i}{\sum_{i=1}^T x_i^2} \quad (3)$$

โดยที่

$$x_i = X_i - \bar{X}; \bar{X} = \frac{1}{T} \sum_{i=1}^T X_i$$

$$y_i = Y_i - \bar{Y}; \bar{Y} = \frac{1}{T} \sum_{i=1}^T Y_i$$

และเราสามารถหาความแปรปรวนของตัวประมาณทั้งสองได้จากสูตร

$$\text{Var}(\hat{b}) = \frac{\sigma_u^2}{\sum x_i^2}, \quad \text{Var}(\hat{a}) = \frac{\sigma_u^2 \sum X_i^2}{T \sum x_i^2} \quad (4)$$

โดย σ_u^2 หาได้จากตัวประมาณ

$$\begin{aligned} \frac{\sum e_i^2}{T-2} &= \frac{\sum y_i^2 - b \sum x_i y_i}{T-2} \end{aligned} \quad (5)$$

ตัวแบบเชิงเส้นของการถดถอยอย่างง่ายเป็นที่นิยมใช้กันมาก เพราะง่ายต่อการคำนวณในการศึกษาที่ต้องการทราบแต่เพียงแนวโน้มของตัวแปรเพียง 2 ตัว การใช้ตัวแบบนี้ก็เป็นประโยชน์มาก แต่การใช้ตัวแบบนี้ก็มีขอบเขตจำกัด จะนำไปใช้กับทุกกรณีไม่ได้

ประการแรกเราควรใช้ตัวแบบนี้สำหรับกรณีที่ตัวแปร Y และ X มีความสัมพันธ์กันอย่างจริง ๆ หรือที่คิดว่าน่าจะมีความสัมพันธ์กัน ทั้งนี้เพราะเราต้องการใช้ตัวแบบนี้ในการพยากรณ์ ถ้าเราทราบค่า X ก็สามารถทำนายค่า Y ได้ ถ้า X และ Y ไม่มีความสัมพันธ์กันในลักษณะ 'เหตุ' และ 'ผล' การใช้ตัวแบบนี้ก็ให้ประโยชน์น้อย หรือไม่มีประโยชน์เลย จริงอยู่ที่เรามีวิธีทางสถิติที่จะทดสอบ สมมติว่า X และ Y มีความสัมพันธ์กันหรือไม่ (นั่นคือ $b = 0$ หรือ $b \neq 0$) แต่นี่ก็เป็นเพียงผลทางสถิติเท่านั้น เพราะว่าการทดสอบอาจจะบอกได้ว่า X และ Y มีความสัมพันธ์กัน ค่าของ r^2 (coefficient of determination) ซึ่งจะกล่าวถึงในหัวข้อต่อไป ก็

อาจมีค่าสูงแต่ในทางปฏิบัติใช้ไม่ได้ เพราะเราอาจเลือกตัวแปรสองตัวโดยเตาสุ่มแล้วนำมาเข้าตัวแบบ ก็มี probability สูงมากที่โดยทางสถิติเราต้องยอมรับว่าตัวแปรทั้งสองมีความสัมพันธ์กัน แต่เราจะได้ผลที่น่าทึ่งมาก เช่น ให้ Y เป็นจำนวนอาชญากรรมที่เกิดขึ้นในประเทศ X เป็นปริมาณถั่วเขียวที่ผลิตได้ในประเทศ ผู้เขียนเชื่อว่าตัวแปรทั้งสองจะเข้ากับตัวแบบได้ดี แต่คงไม่มีใครยอมรับได้ว่าเราจะสามารถทำนายจำนวนอาชญากรรมได้จากการผลิตถั่วเขียว หรือถั่วเขียวเป็นเหตุให้เกิดอาชญากรรม

ประการที่สองตัวแบบที่กล่าวถึงนี้เป็น ตัวแบบเชิงเส้น เพราะเส้นที่ใช้ทำนายค่าของตัวแปรจะเป็นเส้นตรง ตัวแบบนี้เหมาะสำหรับกรณีที่ตัวแปรมีค่าเพิ่มก็เพิ่มด้วยกันหรือ มีค่าลดลงก็ลดด้วยกัน เพราะฉะนั้นตัวแบบนี้จึงไม่เหมาะสำหรับกรณีที่ตัวแปรมีความเปลี่ยนแปลงขึ้นลง (fluctuation) มาก เช่นในกรณีของผลิตผลทางเกษตรกรรมอย่าง ทั้งนี้เพราะผลิตผลเหล่านั้นขึ้นอยู่กับปัจจัยหลายอย่างซึ่งเปลี่ยนแปลงอยู่เรื่อย

ประการที่สาม วิธีที่ใช้ประมาณค่าของ a และ b จากตัวแบบเช่นวิธีกำลังสองน้อยที่สุดจะให้ตัวประมาณที่ดีก็ต่อเมื่อสมมติฐานของตัวแบบเป็นจริง ยกตัวอย่างเช่นเราตั้งข้อสมมติว่าตัวแปร X และ Y จะวัดค่าได้โดยไม่ผิดพลาดหรือความผิดพลาดมีน้อยมากจนตัดทิ้งไปได้ แต่ถ้าในทางปฏิบัติค่าของตัวแปรที่วัดได้ต่างกับค่าจริงมาก การใช้วิธีกำลังสองน้อยที่สุดตามธรรมดา ก็จะไม่ให้ตัวประมาณที่ดี หรือถ้าข้อสมมติที่ว่าความแปรปรวนของตัวกระจายคงที่ไม่เป็นจริง เราก็จำเป็นต้องคิดแปลงวิธีการประมาณใหม่

ประการที่สี่ ในตัวแบบ $Y_i = a + bX_i + u_i$ เราถือว่า X เป็นตัวแปรอิสระ และสมการนี้เป็นสมการเดี่ยวแต่ตัวมันเอง ถ้าสมมติฐานนี้และสมมติฐานอื่น ๆ ของตัวแบบเป็นจริง วิธีกำลังสองน้อยที่สุดก็จะให้ตัวประมาณที่ดี แต่ถ้าสมการนี้เป็นสมการหนึ่งใน Simultaneous Equation System การใช้วิธีกำลังสองน้อยที่สุด ตามธรรมดา ก็จะไม่ดี เพราะจะเกิด bias และ inconsistency ขึ้น⁴ ตัวอย่างของกรณีนี้ อาจเห็นได้จาก Consumption Function $C_i = a + bY_i + u_i$ โดย $C_i =$ ค่าใช้จ่ายในการบริโภค $Y_i =$ รายได้ ค่าของ C และ Y จะถูกตัดสินร่วมกัน การประมาณที่ดีจะต้องใช้ระบบ Simultaneous Equation การ

⁴ J. Johnston, *Econometric Methods* (New York: Mc-Graw-Hill, 1963), ch. 9.

ใช้วิธีกำลังสองน้อยที่สุดตามธรรมดาในการประมาณค่าของ $C_i = a + bY_i + u_i$ อย่างตรงไปตรงมา เราก็จะได้สิ่งที่ไม่มีความหมายมากนัก

สหสัมพันธ์ (Correlation)

ในการทำการวิเคราะห์ที่ใช้ตัวแบบการถดถอยมักจะให้ค่าของสัมประสิทธิ์สหสัมพันธ์ (Correlation Coefficient) r ด้วย ซึ่งอาจหาได้จากสูตร

$$r = \frac{\sum xy}{[\sum x^2 (\sum y^2)]^{1/2}} \quad (6)$$

ตัวสัมประสิทธิ์สหสัมพันธ์แสดงถึงทิศทางและความสัมพันธ์ระหว่างตัวแปรทั้งสอง ถ้า r มีค่าสูง หมายความว่าตัวแปรทั้งสองมีความสัมพันธ์กันมาก และเป็นตรงข้ามถ้าหาก r มีค่าต่ำ ส่วนเครื่องหมายของ r จะบอกว่า Y จะเพิ่มขึ้นหรือลดลงไปกับ X แต่สัมประสิทธิ์สหสัมพันธ์เป็นเครื่องวัดความสัมพันธ์ระหว่างตัวแปรที่ใช้ได้ในกรณีที่ตัวแปร X และ Y มีความสัมพันธ์อย่างเชิงเส้นเท่านั้น ถ้า $r = 0$ อาจเป็นเพราะ X และ Y ไม่มีความสัมพันธ์กันเลยหรือมีความสัมพันธ์ไม่ใช่เชิงเส้น (non-linear relationship)

การตีความหมายสัมประสิทธิ์สหสัมพันธ์ว่าเป็นเครื่องวัดความสัมพันธ์เชิงเส้น ระหว่างตัวแปรสองตัวนั้นเป็นคุณสมบัติทางคณิตศาสตร์เท่านั้น มิได้มีนัยต่อความสัมพันธ์แบบ "สาเหตุ" และ "ผล" แต่อย่างไร การที่ตัวแปรมีแนวโน้มที่เพิ่มขึ้นหรือลดต่ำด้วยกันมิได้หมายความว่าตัวแปรตัวหนึ่งจะมีอิทธิพลต่อตัวแปรอีกตัวหนึ่งโดยตรงหรือโดยอ้อมเสมอไป ในบางกรณีตัวแปรตัวหนึ่งจะมีอิทธิพลต่อตัวแปรอีกตัวหนึ่งจริง แต่ค่าของสัมประสิทธิ์สหสัมพันธ์มีอาจบอกเราได้ ตัวอย่างที่ยกมาแล้วข้างต้น คือจำนวนอากาศยานกรรม และปริมาณผลผลิตของถั่วเขียว ถ้าเรากำหนดหาค่าสัมประสิทธิ์สหสัมพันธ์ ผู้เขียนเชื่อว่าจะต้องมีค่าสูง ซึ่งแสดงว่าสิ่งทั้งสองอย่างต่างมีแนวโน้มสูงชันด้วยกัน ตัวสัมประสิทธิ์นี้มีอาจบอกอะไรเราได้มากกว่านี้

สัมประสิทธิ์อีกตัวหนึ่งที่น่าสนใจคือ r^2 ซึ่งเรียกว่า Coefficient of Determination หากได้จากสูตร

$$r^2 = 1 - \frac{\sum e_i^2}{\sum y_i^2} \quad (7)$$

ตัวสัมพันธ์กันนี้มีความสัมพันธ์กับความแปรปรวน (Variances) ของ X และ Y โดย

$$\begin{aligned} \frac{\sum \hat{y}_i^2}{\sum y_i^2} &= \frac{\sum (\hat{bx}_i)^2}{\sum y_i^2} \\ &= \hat{b} \frac{\sum x_i^2}{\sum y_i^2} = \frac{(\sum x_i y_i)^2}{(\sum x_i^2)^2} \cdot \frac{\sum x_i^2}{\sum y_i^2} \\ &= r^2 \end{aligned} \quad (8)$$

นั่นคือ $r^2 =$ สัดส่วนของความแปรปรวนของ Y ที่อธิบายโดยอิทธิพลของ X เช่น ถ้า r^2 มีค่า 0.8 หมายความว่า การลดลงโดยวิธีกำลังสองน้อยที่สุดของ Y ต่อ X จะอธิบาย 80% ของความแปรปรวนของ Y

สำหรับตัว r^2 นี้ ผู้ทำการวิเคราะห์ส่วนมากมักอธิบายว่าเป็นตัวที่แสดง goodness of fit แต่ความจริงแล้ว r^2 มิได้ให้ความรู้เกี่ยวกับเรื่องนี้เลย r^2 เพียงแต่บอกรูปลักษณะของการรวมกลุ่มของจุดต่างๆ เท่านั้น การที่ r^2 มีค่าสูงเราพอที่จะถือว่าตัวแบบเชิงเส้นที่เราใช้มีความเหมาะสม ในตัวอย่างเล็กๆ ซึ่งปกติใช้ในการวิเคราะห์ทางเศรษฐศาสตร์หรือสังคมศาสตร์แขนงอื่นนั้น r^2 จะมีค่าสูงได้ถึงแม้ความเคลื่อนคลาดมาตรฐาน (standard error) ของ \hat{a} และ \hat{b} จะมีค่าสูง นั่นคือเราอาจมีช่วงทำนายผลกว้างมากเกินไป จนไม่มีประโยชน์อะไร ทั้ง ๆ ที่ r^2 มีค่าสูงในทางตรงข้ามตัวอย่างใหญ่จะให้ค่าของ \hat{a} และ \hat{b} ที่ใกล้เคียงกับ a b มาก แต่ r^2 อาจมีค่าต่ำ⁵ Malinvaud ได้มีตัวอย่างดังนี้

$$\begin{aligned} \hat{Y} &= -1.78 + 1.20 X \\ &\quad (0.40) \quad (0.10) \\ r^2 &= 0.59 ; T = 112 \end{aligned} \quad (9)$$

$$\begin{aligned} \hat{Y} &= -6.20 + 0.115 X \\ &\quad (1.88) \quad (0.012) \\ r^2 &= 0.95 ; T = 11 \end{aligned} \quad (10)$$

⁵ Malinvaud, *op. cit.*, p. 83

จะเห็นได้ว่าใน (9) ค่า r^2 ต่ำกว่า (10) มากทั้งที่อัตราส่วนระหว่างค่าของ \hat{a} , \hat{b} กับความเปลี่ยนแปลงมาตรฐานของแต่ละตัว (ตัวเลขในวงเล็บข้างล่าง ตัวสัมประสิทธิ์ \hat{a} , \hat{b}) จะใกล้เคียงกัน ทั้งนี้เพราะใน (9) เรามีข้อมูลถึง 112 ส่วน ใน (10) มีเพียง 11 เท่านั้น เพราะฉะนั้นจึงพอสรุปได้ว่าการใช้ r^2 แต่เพียงอย่างเดียวในการเลือกตัวแบบที่ดีจากตัวแบบที่ให้ค่า a และ b ต่าง ๆ กันนั้นไม่เป็นการเพียงพอ

หากที่กล่าวมาทั้งหมดนี้ความแตกต่างระหว่างตัวแบบการถดถอยและสหสัมพันธ์ก็คงชัดเจนพอสมควร ความหมายของตัวแบบทั้งสองมีความสำคัญมาก แต่มักขาดการเน้นในการสอนสถิติเบื้องต้นแก่นักศึกษา ซึ่งเป็นการเสี่ยงต่อการใช้แบบผิด ๆ สิ่งที่น่าสนใจให้ผู้เขียน เขียนเรื่องนี้ก็คือ การที่ได้อ่านวิทยานิพนธ์ชั้นปริญญาโทของมหาวิทยาลัยบางแห่งและบทความในหนังสือปริทัศน์ ตลอดจนผลงานวิจัยซึ่งใช้วิธีการทางสถิติทั้งสองอย่างนี้ อย่างผิด ๆ และตีความผิดพลาด

ความสำเร็จของการใช้เทคนิคทั้งสองอย่างนี้หรือเทคนิคอื่น ๆ ทางสถิติขึ้นอยู่กับความเข้าใจในคุณสมบัติของตัวแบบ ความหมาย ข้อจำกัด และการใช้ถูกต้อง เหมาะสมกับภาวะเหตุการณ์ การสอนสถิติเบื้องต้นจึงน่าจะเน้นหนักไปในเรื่องเหล่านี้บ้าง มิใช่แต่ให้นักศึกษาท่องสูตรแล้วนำไปใช้อย่างผิด ๆ

จำลอง อติกุล