

โมเดลที่ใช้ในการพยากรณ์

ในที่นี้จะทำการศึกษารเปรียบเทียบจากผลของการวิเคราะห์อนุกรมเวลา (Time Series Analysis) กับการวิเคราะห์ความถดถอย (Regression Analysis) ในการพยากรณ์ปริมาณการบริโภคและปริมาณผลผลิตน้ำคาลทรายภายใน ประเทศ เพื่อที่จะได้เลือกวิธีใดที่ให้ผลในการพยากรณ์ได้ดีกว่ากัน

๓.๑ การวิเคราะห์อนุกรมเวลา (Time Series Analysis)

การวิเคราะห์อนุกรมเวลา คือการศึกษาถึงความเคลื่อนไหวของข้อมูล ชุดหนึ่ง ๆ ตามงากระยะเวลา ข้อมูลต่าง ๆ อาทิ เช่น ข้อมูลของปริมาณการบริโภคน้ำคาลทรายภายในประเทศ และข้อมูลของปริมาณผลผลิตน้ำคาลทราย ภายในประเทศที่เก็บรวบรวมมาขึ้น เมื่อข้อมูลที่ได้อาจให้เรียงลำดับตามเวลา (ตามปี พ.ศ.) ที่เกิดขึ้นมาแล้ว

ส่วนประกอบของอนุกรมเวลา มีดังนี้ :-

๓.๑.๑ แนวโน้มตามลำดับเวลา (Secular Trend หรือ Trend หรือ Secular Movements) ซึ่งจะใช้ตัวย่อว่า "T" เป็นการเคลื่อนไหวของข้อมูล ในระยะเวลาที่ค่อนข้างจะยาว ตามปกติจะแสดงถึงทิศทางของอนุกรมเวลาชุดนั้น ๆ พุ่งไป สู่ ที่ซึ่งอาจจะมีลักษณะ เป็นเส้นตรง (Linear Trend) หรือ เส้นโค้ง (Non Linear Trend) ก็ได้ หน่วยของระยะเวลาตามปกติเป็นปี

๓.๑.๒ การเปลี่ยนแปลงตามฤดูกาล (Seasonal Variations หรือ Seasonal movements) จะใช้ตัวย่อว่า "S" เป็นการเปลี่ยนแปลงของข้อมูลที่มีลักษณะเปลี่ยนแปลงขึ้น ๆ ลง ๆ อย่างเดียวกันหรือ คล้าย ๆ กัน ในช่วงระยะเวลาที่กำหนด หน่วยของระยะเวลาอาจเป็นงวด ๓ เดือน รายเดือน รายสัปดาห์ หรือ รายวันก็ได้

๓.๑.๓ การเปลี่ยนแปลงตามวัฏจักร (Cyclical Variations หรือ

Cyclical Movements) จะใช้ตัวย่อ "C" เป็นการเปลี่ยนแปลงของข้อมูล ที่มีลักษณะคล้ายลูกคลื่น กล่าวคือจะมีลักษณะของข้อมูลพุ่งสูงขึ้น ในช่วงระยะเวลาหนึ่ง และติดตามด้วยการลดต่ำลงในอีกช่วงระยะเวลาหนึ่งสลับกันไป

๓.๑.๔ การเปลี่ยนแปลงเนื่องจากเหตุการณ์ผิดปกติหรือเป็นเชิงสุ่ม (Irregularity Variations หรือ Irregularity Movements หรือ Random Movements) จะใช้ตัวย่อว่า "I" เป็นการเปลี่ยนแปลงที่ไม่อาจคาดการณ์ได้ล่วงหน้า เหตุการณ์ผิดปกตินี้อาจเป็นเรื่องเล็ก ๆ น้อย ๆ เช่นการนัดหยุดงานของแรงงานไปจนถึงเรื่องใหญ่ ๆ เช่น การเกิดแผ่นดินไหว หรือการเกิดสงคราม ซึ่งเป็นปัจจัยที่มีใ้ อยู่ภายใต้กฎเกณฑ์ที่แน่นอนแต่อย่างใด

ในอนุกรมเวลาชุดหนึ่ง ๆ อาจมีส่วนประกอบเฉพาะบางส่วน หรือมีรวมกัน ทั้งหมด (T, S, C, I) ก็เป็นได้ และโมเดล (Model) ชั้นมูลฐานที่ใช้ในการแสดงส่วนประกอบของอนุกรมเวลามี ๒ โมเดล

ก. โมเดลเชิงบวก (Additive Model)

ได้แก่การนำเอาองค์ประกอบทั้ง ๔ มาบวกกันจะได้ $Y = T + S + C + I$

ข. โมเดลเชิงคูณ (Multiplicative Model)

ได้แก่การนำเอาองค์ประกอบทั้ง ๔ มาคูณกันจะได้ $Y = T \cdot S \cdot C \cdot I$

โมเดลที่จะใช้ในที่นี้เป็นโมเดลเชิงคูณ

๓.๒ วิธีการประมาณค่าแนวโน้มตามลำดับเวลา (หรือวิธีคำนวณหาเส้น แนว

โน้ม)

ปกติในการวิเคราะห์ แนวโน้มตามลำดับเวลา มักจะใช้ข้อมูลรายปีแทนที่จะเป็นงวด ๓ เดือน รายเดือน หรือสัปดาห์ ซึ่งจะทำให้เกิดความยุ่งยากในการคำนวณ โดยไม่จำเป็น โดยนำข้อมูลที่เตรียมไว้แล้วมาเขียนร่างในกระดาษกราฟก่อน เพื่อจะได้พิจารณาลักษณะของ แนวโน้มที่ควรจะเป็นว่าจะมีลักษณะเป็นเส้นตรงหรือ เส้นโค้งแล้วจะได้นำเอาเส้นแนวโน้มนั้น ๆ เพื่อใช้ในการพยากรณ์ต่อไป วิธีการประมาณค่าแนวโน้มตามลำดับเวลา มีอยู่หลายวิธีด้วยกัน แต่วิธีที่เห็นว่าเหมาะสม นิยมใช้มากที่สุด และ

ง่ายในการคำนวณ โดยใช้หลักที่ว่า ผลรวมของผลต่างระหว่างค่าของข้อมูลที่ได้ออกจากเส้นแนวโน้มกับค่าจริงของข้อมูล ยกกำลังสองแล้วได้ค่าน้อยที่สุด ซึ่งเรียกว่าวิธีกำลังสองน้อยที่สุด (Least Squares Method)

๓.๒.๑ กรณีที่ได้แนวโน้มมีลักษณะเป็นเส้นตรง (Linear Trend)

โมเดลที่ใช้แทนเส้นตรงคือ

$$Y = A + BX + E \dots\dots\dots (1)$$

A ค่าคงที่หรือเป็นค่า Y เมื่อ X = 0 (ระยะตัดแกน Y)

B ค่าความชัน (Slope) ของเส้นตรง

E ค่าผิดพลาด (Error หรือ Disturbance Term)

X เป็นตัวแปรอิสระ (ในที่นี้คือระยะเวลา มีหน่วยเป็น ปี)

Y เป็นค่าจริงของข้อมูล (เป็นตัวแปรตาม)

ค่าประมาณของ (1) คือ

$$\hat{Y} = a + b X \dots\dots\dots (2)$$

\hat{Y} คือค่าประมาณของ Y

a คือค่าประมาณของ A

b คือค่าประมาณของ B

$e = Y - \hat{Y} =$ ความแตกต่างระหว่างค่าจริงกับค่าประมาณหรือ ความคลาดเคลื่อนหรือ Residual

$$\sum e^2 = \sum (Y - \hat{Y})^2 \dots\dots\dots (3)$$

ในการคำนวณหาค่า a, b ซึ่งเป็นค่าคงที่ เพื่อให้ได้ความคลาดเคลื่อนน้อยที่สุดจะใช้วิธีกำลังสองน้อยที่สุด ก็จะได้ a, b เป็นค่า Unbiased Estimate และมี Standard Error น้อยที่สุดด้วยเมื่อเปรียบเทียบกับ Unbiased Estimate ตัวอื่น ๆ ที่คำนวณได้จากโมเดลเดียวกัน

การหาค่า a, b ทำได้โดยใช้ค่าเชิงลเกริเวทีฟ (Partial derivative)

ของ ③ เทียบกับ a, b แล้วให้ค่าเป็น 0

$$\frac{\partial \Sigma e^2}{\partial a} = 0, \quad \frac{\partial \Sigma e^2}{\partial b} = 0$$

จะได้ Normal Equation ดังนี้

$$\left. \begin{aligned} \Sigma Y &= Na + b \Sigma X \\ \Sigma XY &= a \Sigma X + b \Sigma X^2 \end{aligned} \right\} \dots\dots\dots ①'$$

เนื่องจากข้อมูลมี ๑๔ มีนึ่งเป็นเลขคู่ เพื่อง่ายแก่การคำนวณจึงเลือกให้มีดังนี้

- ให้มี ๒๕๐๘ มีค่า X = - 13
- ๒๕๐๖ มีค่า X = - 11
- ๒๕๐๓ มีค่า X = - 9
- •
- •
- •
- •
- •
- ๒๕๐๘ มีค่า X = 13

ดังนั้น $\Sigma X = 0, \quad \Sigma X^3 = 0$

และ $\Sigma X^2 = 910, \quad \Sigma X^4 = 105,742$

เพื่อความสะดวกในการคำนวณ จาก ①' เขียนได้ใหม่เป็น

$$\Sigma Y = 14 a, \quad N = 14$$

$$\Sigma XY = 910 b$$

หรือ
$$\left. \begin{aligned} a &= \Sigma Y / 14 \\ b &= \Sigma XY / 910 \end{aligned} \right\} \dots\dots\dots (2)'$$

๓.๒.๒ กรณีที่ได้แนวโน้มมีลักษณะเป็นเส้นโค้ง (Non Linear Trend)

กรณีที่ต่ำกว่าลักษณะข้อมูลที่ได้มีลักษณะเป็นเส้นโค้ง เมื่อทดลองนำไปเขียนโดยพิจารณา จาก Scatter Diagram การเลือกลักษณะโมเดลโดยทั่ว ๆ ไป จะอยู่ใน ๒ ลักษณะดังนี้คือ

ก. โมเดลที่อยู่ในรูปพหุนามกำลังสอง (Quadratic Polynomial)

$$\hat{Y} = a + bX + cX^2 \dots\dots\dots (4)$$

โดยที่ a เป็นค่าของ \hat{Y} เมื่อ X มีค่า = 0

b เป็นค่าของความชัน (Slope) ของแนวโน้ม

c เป็นอัตราการเปลี่ยนของความชัน

การหาค่า a, b, c ทำโดยวิธีใช้กำลังสองน้อยที่สุดเช่นกัน จะได้ Normal Equation ดังนี้

$$\left. \begin{aligned} \Sigma Y &= Na + b \Sigma X + c \Sigma X^2 \\ \Sigma XY &= a \Sigma X + b \Sigma X^2 + c \Sigma X^3 \\ \Sigma X^2 Y &= a \Sigma X^2 + b \Sigma X^3 + c \Sigma X^4 \end{aligned} \right\} \dots\dots\dots (4)'$$

แทนค่า N, ΣX , ΣX^2 , ΣX^3 , ΣX^4 จะได้

$$\left. \begin{aligned} \Sigma Y &= 14a + 910c \\ \Sigma XY &= 910b \\ \Sigma X^2 Y &= 910a + 105,742c \end{aligned} \right\} \text{--- (5)}$$

ดังนั้น

$$\left. \begin{aligned} a &= \frac{105,742 (\Sigma Y) - 910 (\Sigma X^2 Y)}{14 (105,742) - (910)^2} \\ b &= \Sigma XY / 910 \\ c &= \frac{14 (\Sigma X^2 Y) - 910 (\Sigma Y)}{14 (105,742) - (910)^2} \end{aligned} \right\} \text{--- (6)}$$

พิจารณา $\frac{14 (105,742) - (910)^2}{14} = 105,742 - 59,150$

$$= 46,592$$

$$\frac{105,742}{14} = 7,553$$

$$910/14 = 65$$

$$7,553 / 46,592 = 0.162109$$

$$65 / 46,592 = 0.001395$$

$$1 / 46,592 = 0.0000214$$

ดังนั้นจาก

(6) จะได้

$$\left. \begin{aligned} a &= 0.162109 \Sigma Y - 0.001395 \Sigma X^2 Y \\ b &= 0.001099 \Sigma XY \\ c &= 0.0000214 \Sigma X^2 Y - 0.001395 \Sigma Y \end{aligned} \right\} \text{--- (7)}$$

ข. โมเดลที่อยู่ในรูปเอกโพเนนเชียล (Exponential Trend)

มีรูป $Y = ab^X \dots \dots \dots (5)$

- a คือค่า \hat{Y} เมื่อ $X = 0$
 b อัตราการเพิ่ม (rate of increase)
 b ถ้ามีค่าอยู่ระหว่าง 0 กับ 1 แสดงว่า Y จะมีค่าลดลงเมื่อ X มีค่าเพิ่มขึ้น
 b ถ้ามากกว่า 1 แสดงว่า Y จะมีค่าเพิ่ม เมื่อ X มีค่าเพิ่มขึ้น

ในการหาค่า a, b จะยุ่งยาก วิธีที่สะดวกทำโดยการ Take \log สมการ (5) จะได้

$$\log \hat{Y} = \log a + X \log b \quad \text{-----} \quad (6)$$

การหาค่า $\log a$, $\log b$ จะเหมือนกรณีเส้นตรง โดยใช้วิธีกำลังสองน้อยที่สุดเช่น ก็จะได้ Normal Equation ดังนี้

$$\left. \begin{aligned} \sum \log Y &= N \log a + (\log b) (\sum X) \\ \sum X \log Y &= (\log a) (\sum X) + (\log b) (\sum X^2) \end{aligned} \right\} \quad (7)$$

แทนค่า $N, \sum X, \sum X^2$ จะได้

$$\sum \log Y = 14 \log a$$

$$\sum X \log Y = 910 \log b$$

$$\left. \begin{aligned} \text{ดังนั้น} \quad \log a &= (\sum \log Y) / 14 \\ \log b &= (\sum X \log Y) / 910 \end{aligned} \right\} \quad \text{-----} \quad (8)$$

เมื่อเปิดค่า Anti log จะได้ค่า a, b ตามต้องการ

การพิจารณาจาก Scatter Diagram จะไม่สามารถบอกรายละเอียดได้มากนัก จึงเห็นสมควรทำการทดสอบค่าลักษณะของเส้นแนวโน้มที่ควรจะเป็นใน ๓ โมเดล คือ :-

๑. โมเดลของเส้นตรง
๒. โมเดลของพหุนาม

๓. โมเดลของเฮกซ์โพเนนเชียล

จากนั้นจึงพิจารณาถึง ผลบวกกำลังสองของผลต่างระหว่างข้อมูลเดิมกับที่คำนวณได้จากเส้นแนวโน้ม (Sum of Square) ในแต่ละโมเดล เพื่อหาค่าที่คำนวณได้จากโมเดลใด มีค่าน้อยที่สุดก็เลือกเอาโมเดลนั้นเป็นตัวแทนของ Time Series Analysis หรือการเลือกโมเดล โดยการคำนวณหาค่า Standard Error of Estimate of Y on X or Standard Deviation of Regression

$$S_{Y.X} = \sqrt{\frac{\sum(Y - \hat{Y})^2}{(n - k)}} \dots\dots\dots (9)$$

$$n = \text{จำนวน ข้อมูล} = 14$$

$$k = \text{จำนวนพารามิเตอร์} = 2 \quad (\text{สำหรับเส้นตรง})$$

และคำนวณหาค่าเปอร์เซ็นต์ของการกระจายของ Y (Dependent Variable) ซึ่งคำนวณได้จาก

$$\frac{\sum(\hat{Y} - \bar{Y})^2}{\sum(Y - \bar{Y})^2} \times 100 \dots\dots\dots (10)$$

\bar{Y} คือค่ามัธยฐานเลขคณิตของตัวแปร ไม่อิสระ

โมเดลใดที่ให้ค่า R ต่ำสุดและเปอร์เซ็นต์ของการกระจายมากที่สุดก็เลือกใช้โมเดลนั้นเป็นตัวแทน แต่บางครั้งก็อาจมีปัญหาเกิดขึ้น กล่าวคือเมื่อได้ทำการเลือกตัดสิน ตามกฎเกณฑ์ดังกล่าวแล้ว เมื่อนำไปใช้พยากรณ์ในอนาคตจะมีแนวโน้มที่ Diverge ค่อนข้างเร็ว ทำให้ค่าในอนาคตเปลี่ยนแปลงเร็วเกินความเป็นจริงบางครั้งจำเป็นต้องเลือกโมเดลอื่นที่เหมาะสมและสามารถให้เหตุผลพอประมาณในการพยากรณ์ได้

๓.๒.๓ วิธีการทดสอบโมเดล

การทดสอบความมีนัยสำคัญของโมเดล โดยใช้ **F - test**

เมื่อเลือกโมเดลได้ตามลักษณะของารพิจารณาข้างกล่าว จำเป็นที่จะต้องนำโมเดลที่ได้มาทดสอบนัยสำคัญ ภายใต้สมมติฐาน

$$H_0 : B = 0 \quad \text{หรือ} \quad B = C = 0$$

ค่าของ F พิจารณาได้จากตารางการวิเคราะห์ความแปรปรวน (Analysis of Variance) หรือตาราง ANOVA

ANOVA

Source of Variation	Degree of Freedom	Sum of Square	Mean Square	F_c
Regression	k	$\sum (\hat{Y} - \bar{Y})^2 = SSR$	$MSR = SSR/k$	$\frac{MSR}{MSE}$
Residual	n-k-1	$\sum (Y - \hat{Y})^2 = SSE$	$MSE = SSE / (n-k-1)$	
Total	n - 1	$\sum (Y - \bar{Y})^2 = SST$		

ต่อไปพิจารณาค่า F_c ที่คำนวณได้จาก ANOVA เปรียบเทียบกับค่า F_t จากตาราง F - distribution ที่ d.f = (k, n-k-1) ณ ระดับความเชื่อมั่นต่าง ๆ อาทิเช่น $\alpha = 0.01$ หรือ 0.05 ถ้าปรากฏว่า ค่า $F_c > F_t$ จะหมายความว่าเราจะปฏิเสธสมมติฐานที่ตั้งขึ้น นั่นคือเราจะยอมรับว่าโมเดลที่นำมาทำการทดสอบนั้นใช้ได้

หมายเหตุ	ค่าของ $\sum (Y - \bar{Y})^2$	เรียกว่า Total Variation
	ค่าของ $\sum (Y - \hat{Y})^2$	เรียกว่า Unexplained Variation
	ค่าของ $\sum (\hat{Y} - \bar{Y})^2$	เรียกว่า Explained Variation

Coefficient of determination: $(R^2) = \frac{\text{Explained Variation}}{\text{Total Variation}}$

$= 1 - \frac{\text{Unexplained Variation}}{\text{Total Variation}}$

๓.๓ การวิเคราะห์ความถดถอย (Regression Analysis)

เนื่องจากการวิเคราะห์ปริมาณการบริโภคน้ำคานและปริมาณผลผลิตน้ำคาน-ทรายภายในประเทศมีลักษณะเกี่ยวข้องกับหลายตัวแปร จึงจำเป็นต้องใช้การวิเคราะห์ความถดถอยเชิงซ้อน (Multiple Linear Regression) ซึ่งเป็นการวิเคราะห์ความสัมพันธ์ของข้อมูลตั้งแต่ ๒ ชุดขึ้นไป คือเมื่อกำหนดให้ข้อมูลชุดหนึ่งเป็นตัวแปรตาม และข้อมูลชุดอื่น ๆ เป็นตัวแปรอิสระ ในรูปสมการกำลังหนึ่ง

$$Y = B_0 + B_1 X_1 + B_2 X_2 + B_3 X_3 + \dots + B_k X_k + E$$

ค่าประมาณของสมการนี้คือ $\hat{Y} = a + b_1 X_1 + b_2 X_2 + b_3 X_3 + \dots + b_k X_k \dots (4)$

- X_i เป็นตัวแปรอิสระที่ i โดยที่ $i = 1, 2, 3, \dots, k$
- Y เป็นตัวแปรตาม
- B_i เป็นค่าสัมประสิทธิ์ของตัวแปรอิสระตัวที่ $i, i = 1, 2, 3, \dots, k$
- \hat{Y} เป็นค่าประมาณของ Y
- b_i เป็นค่าประมาณของ B_i
- a เป็นค่าประมาณของ B_0

๓.๓.๑ การคำนวณหาค่าสัมประสิทธิ์ที่แปรอิสระในสมการความถดถอยเชิงซ้อน

จากสมการ (A) ให้ Y เป็นค่าจริงของข้อมูล
 \hat{Y} เป็นค่าประมาณของข้อมูลใน (A)

$$e = \text{ความคลาดเคลื่อน} = Y - \hat{Y}$$

$$= \text{ความแตกต่างระหว่างค่าจริงกับค่าประมาณของข้อมูล}$$

$$e = Y - a - b_1X_1 - b_2X_2 - \dots - b_kX_k$$

ถ้ามี n ข้อมูลจะได้

$$\sum_{i=1}^n e^2 = \sum_{i=1}^n (Y_i - a - b_1X_{i1} - \dots - b_kX_{ik})^2 \dots \dots \dots (B)$$

โดยวิธีกำลังสองน้อยที่สุดจะได้ค่า a, b_1, b_2, \dots, b_k เป็นค่า Unbiased Estimate และมีค่า Standard Error น้อยที่สุดเมื่อเปรียบเทียบกับ Unbiased Estimate ทั่วอื่น ๆ ที่คำนวณได้จากโมเดลเดียวกัน

การคำนวณหาค่า a

โดยการใช้อนุพันธ์ย่อย (Partial derivative) ของ (B)

เทียบกับ a แล้วให้เท่ากับ 0 นั่นคือ

$$\left(\frac{\partial \sum e^2}{\partial a} \right) = 0$$

$$- 2 \sum (Y - a - b_1X_1 - b_2X_2 - \dots - b_kX_k) = 0$$

$$a = \bar{Y} - b_1\bar{X}_1 - b_2\bar{X}_2 - \dots - b_k\bar{X}_k$$

แทนค่า a ใน (B) จะได้

$$\Sigma e^2 = \Sigma \left\{ Y - \bar{Y} - b_1 (X_1 - \bar{X}_1) - \dots - b_k (X_k - \bar{X}_k) \right\}^2$$

ให้ $y_i = Y_i - \bar{Y}_1$

$x_i = X_i - \bar{X}_1$

∴ $\Sigma e^2 = \Sigma (y - b_1 x_1 - b_2 x_2 - \dots - b_k x_k)^2 \dots \dots \dots (C)$

การคำนวณหาค่า b_1

โดยการใช้พหุคูณย่อยเชิงอนุพันธ์ (Partial derivative) ของ (B) เทียบกับ

b_1 แล้วให้ค่า = 0 จะได้ k สมการซึ่งเป็น Normal Equation สำหรับ X_1

k ตัว

$$\left. \begin{aligned} \frac{\partial \Sigma e^2}{\partial b_1} = 0 ; & b_1 \Sigma x_1^2 + b_2 \Sigma x_1 x_2 + \dots + b_k \Sigma x_1 x_k = \Sigma x_1 y \\ & b_1 \Sigma x_2 x_1 + b_2 \Sigma x_2^2 + \dots + b_k \Sigma x_2 x_k = \Sigma x_2 y \\ & \dots \\ & \frac{\partial \Sigma e^2}{\partial b_k} = 0 ; b_1 \Sigma x_k x_1 + b_2 \Sigma x_k x_2 + \dots + b_k \Sigma x_k^2 = \Sigma x_k y \end{aligned} \right\} (D)$$

จาก (D) แก้สมการหาค่า b_1, b_2, \dots, b_k ได้จาก k สมการ และกรณีที่มีข้อมูลมี X_i มากกว่า ๓ ตัว ขึ้นไปควรจะหาค่า b_i โดยวิธีเมทริกซ์ (Matrix) จะสะดวกกว่า โดยวิธี Matrix

สมมติมี n ข้อมูลและมีตัวแปรอิสระ ๓ ตัวคือ $k = 3$ เขียนเรียงลำดับได้ในรูป

Matrix ดังนี้ :-

ลำดับ	Y	X_1	X_2	X_3
1	Y_1	X_{11}	X_{21}	X_{31}
2	Y_2	X_{12}	X_{22}	X_{32}
3	Y_3	X_{13}	X_{23}	X_{33}
.
.
n	Y_n	X_{1n}	X_{2n}	X_{3n}



หรือเขียนใหม่ได้เป็น

$$Y = X\beta + \epsilon$$

โดยให้ $Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}$, $X = \begin{bmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ \vdots & \vdots & \vdots \\ X_{n1} & X_{n2} & X_{n3} \end{bmatrix}$

$$\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix}, \quad \epsilon = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}$$

คำนวณหาค่า

$$X'X = \begin{bmatrix} \sum X_{i1}^2 & \sum X_{i1}X_{i2} & \sum X_{i1}X_{i3} \\ \sum X_{i2}X_{i1} & \sum X_{i2}^2 & \sum X_{i2}X_{i3} \\ \sum X_{i3}X_{i1} & \sum X_{i3}X_{i2} & \sum X_{i3}^2 \end{bmatrix}$$

และได้ $X'Y = \begin{bmatrix} \sum X_{i1}Y_i \\ \sum X_{i2}Y_i \\ \sum X_{i3}Y_i \end{bmatrix}$

ได้ Normal Equation

$$\begin{aligned} (X'X)\hat{\beta} &= X'Y \\ (X'X)^{-1} (X'X)\hat{\beta} &= (X'X)^{-1} (X'Y) \\ \hat{\beta} &= (X'X)^{-1} (X'Y) \\ \hat{\beta} &= C (X'Y) \end{aligned}$$

$$\text{ให้ } c = (X'X)^{-1} = \begin{bmatrix} c_{11} & c_{12} & c_{13} \\ c_{21} & c_{22} & c_{23} \\ c_{31} & c_{32} & c_{33} \end{bmatrix}$$

จะได้

$$\beta_1 = c_{11} \sum X_{i1} Y_i + c_{12} \sum X_{i2} Y_i + c_{13} \sum X_{i3} Y_i$$

$$\beta_2 = c_{21} \sum X_{i1} Y_i + c_{22} \sum X_{i2} Y_i + c_{23} \sum X_{i3} Y_i$$

$$\beta_3 = c_{31} \sum X_{i1} Y_i + c_{32} \sum X_{i2} Y_i + c_{33} \sum X_{i3} Y_i$$

ถ้ากำหนดให้ $\hat{y} = \hat{Y} - \bar{Y}$

$$e = Y - \hat{Y}$$

$$y = Y - \bar{Y}$$

$$= (\hat{Y} - \bar{Y}) + (Y - \hat{Y})$$

$$= \hat{y} + e$$

$$\text{ดังนั้น } \sum y^2 = \sum \hat{y}^2 + \sum e^2$$

จากสมการข้างบนจะเห็นว่า ผลบวกกำลังสองของ deviation ของ Y จาก mean ประกอบด้วย ๒ ส่วนดังนี้

ก. $\sum \hat{y}^2$ เรียกว่า Sum of Squares เนื่องจากความถดถอยหนึ่งเป็นผลบวกกำลังสองของ deviation \hat{Y} จาก mean

ข. $\sum e^2$ เป็นผลบวกกำลังสองของ deviation ของค่าจริงจากค่าประ-

มาณ

การคำนวณค่า F value เพื่อทดสอบโมเดลสมการถดถอยเชิงซ้อนได้

จากตาราง ANOVA สำหรับ Multiple Linear Regression

ANOVA

Source of Variation	d.f	Sum of Squares	Mean Square	F _c
Regression	k	$\sum (\hat{Y} - \bar{Y})^2 = \sum \hat{y}^2 = \text{RSS}$	A = RSS/k	
Residual (error)	n-k-1	$\sum (Y - \hat{Y})^2 = \sum e^2 = \text{SSE} = \text{TSS} - \text{RSS}$	B = SSE/n-k-1	A/B
Total	n-1	$\sum (Y - \bar{Y})^2 = \sum y^2 = \text{TSS} = \sum Y^2 - \frac{(\sum Y)^2}{n}$		

n = จำนวนข้อมูล

k = จำนวนตัวแปรอิสระ

เปรียบเทียบค่า F_c ที่ได้จากการคำนวณเทียบกับ F_t ที่ได้จากการวาง F-distribution ที่ d.f = (k, n-k-1) ณ ระดับความเชื่อมั่นต่าง ๆ เช่น α = 0.05, 0.01 ทำให้สามารถที่จะสรุปผลของการวิเคราะห์ได้ว่า มีอิทธิพลต่อ Y อย่างมีนัยสำคัญหรือไม่ การคำนวณค่า t - value เพื่อทดสอบค่า β_i

โดยการตั้งสมมติฐาน

$$H_0 : \beta_i = 0, i = 1, 2, \dots, k$$

$$t_c = \frac{b_i - \beta_i}{s_{b_i}}, \text{ d.f} = n-k-1$$

$$= \frac{b_i}{s_{b_i}}$$

b_i เป็นค่า Regression Coefficient ที่ต้องการทดสอบ

t_c เป็นค่าที่คำนวณได้จากสูตร

S_{b_i} เป็นค่า Standard Error ของ b_i

$$S_{b_i} = \sqrt{C_{ii} S_{Y.1 \dots k}^2}$$

C_{ii} เป็นค่าของ diagonal จาก $(X'X)^{-1}$

$S_{Y.1 \dots k}^2$ เป็นค่า Standard Error ของการประมาณค่า \hat{Y}

ซึ่งเป็นค่าของ Mean Square of deviation ในตาราง ANOVA ของ Multiple Linear Regression

$$S_{Y.1 \dots k}^2 = \frac{\sum (Y - \hat{Y})^2}{n - k - 1}$$

จากนั้นค่า t_c ที่คำนวณได้เปรียบเทียบกับค่า t จากตาราง t-distribution ที่ $d.f = n - k - 1$ ณ ระดับความเชื่อมั่นต่าง ๆ แล้วจึงจะสามารถสรุปผล

ได้ว่าจะยอมรับหรือปฏิเสธสมมติฐาน $H_0 : \beta_i = 0$

๓.๓.๒ สัมประสิทธิ์แห่งสหสัมพันธ์ (Coefficient of Correlation)

สัมประสิทธิ์แห่งสหสัมพันธ์อย่างง่าย (Simple Correlation Coefficient)

เป็นค่าสัมพัทธ์ (Relative Value) แสดงความสัมพันธ์ระหว่างตัวแปร (Variables)

ทั้งสองว่าที่อยู่มากน้อยเพียงใด ตัวแปรทั้งสองปกติจะใช้แทนด้วย X กับ Y และค่าของสัม

ประสิทธิ์แห่งสหสัมพันธ์จะแทนด้วย r_{XY} หรือ R_{XY}

$$r_{XY} = \frac{S_{XY}}{S_X \cdot S_Y}$$

$$S_X = \sqrt{\frac{\sum (X_i - \bar{X})^2}{(n - 1)}} \quad , \quad S_Y = \sqrt{\frac{\sum (Y_i - \bar{Y})^2}{(n - 1)}}$$

$$S_{XY} = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{(n - 1)}$$

S_X , S_Y เป็นค่าเบี่ยงเบนมาตรฐาน (Standard deviation) ของ X และ Y

S_{XY} เป็น Covariance ของ X กับ Y

$$r_{XY} = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum (X_i - \bar{X})^2 \sum (Y_i - \bar{Y})^2}}$$

การพิจารณาค่าของ r มี ๓ กรณี

๑. r มีค่าเป็นบวกหมายความว่าค่าของตัวแปรทั้งคู่จะเปลี่ยนไปในลักษณะเดียวกันคือ เมื่อ X มีค่าเพิ่มขึ้น Y ก็จะมีค่าเพิ่มขึ้นด้วย

๒. r มีค่าเป็นลบ หมายความว่าค่าของตัวแปรทั้งคู่จะเปลี่ยนไปในลักษณะตรงกันข้าม กล่าวคือ X มีค่าเพิ่มขึ้น ค่าของ Y ก็จะมีค่าลดลง

๓. r มีค่าเท่ากับ ๐ หมายความว่า X กับ Y ไม่มีความสัมพันธ์กัน ค่าของ r จะมีค่าอยู่ระหว่าง -1 กับ +1 ; $-1 \leq r \leq 1$

และค่า $|r|$ เป็นค่าสัมบูรณ์ของ r ถ้ามีค่าเข้าใกล้ 1 แสดงว่า X กับ Y มีความสัมพันธ์ต่อกันมาก ถ้ามีค่าเข้าใกล้ ๐ แสดงว่า X กับ Y เกือบไม่มีความสัมพันธ์กัน

ค่าที่คำนวณหรือดัชนี 'ระสิทธิภาพการตัดสินใจ (Coefficient of determination)

$$r^2 = \frac{\left\{ \sum (X_i - \bar{X})(Y_i - \bar{Y}) \right\}^2}{\sum (X_i - \bar{X})^2 \sum (Y_i - \bar{Y})^2}$$

r^2 เรียกว่าดัชนีกำหนด ส่วนที่เหลือ $1 - r^2$ เรียกว่าดัชนีไม่กำหนด

(Coefficient of non determination)

เมื่อให้ $x = (X - \bar{X})$, $y = (Y - \bar{Y})$

ดังนั้น $r^2 = \frac{(\sum xy)^2}{\sum x^2 \sum y^2}$
 $r^2 = \frac{(\sum xy)^2 / \sum x^2}{\sum y^2}$

หรือ $r^2 = \frac{(\sum xy)^2}{\sum x^2 \sum y^2}$

ค่าสัมประสิทธิ์สหสัมพันธ์เชิงส่วน (Partial Correlation Coefficient)

ในกรณีที่มีตัวแปรอิสระหลายตัวอยู่ในสมการถดถอย นอกจากจะหาความสัมพันธ์กันเป็นคู่ ๆ คือหาค่าสหสัมพันธ์อย่างง่ายระหว่างตัวแปร X_i กับ X_j มาแล้ว ถ้าตัวแปรต่าง ๆ ทุกตัวผันแปรพร้อมกันแบบการแจกแจงปกติ (Normal Distribution) เราก็อาจจะหาความสัมพันธ์ในลักษณะสหสัมพันธ์เชิงส่วนของตัวแปรได้ เช่น

$$Y = b_0 + b_1 X_1 + b_2 X_2 + b_3 X_3$$

นอกจากจะได้ค่า r_{YX_1} , r_{YX_2} , r_{YX_3} แล้วเรายังอาจหาค่าสหสัมพันธ์เชิงส่วนได้ เช่น

- $r_{YX_1 \cdot X_2}$ หรือ $r_{Y1.2}$ คือค่าสหสัมพันธ์เชิงส่วนระหว่างตัวแปร Y กับ X_1 โดยถือว่า X_2 คงที่
- $r_{YX_2 \cdot X_1}$ หรือ $r_{Y2.1}$ คือค่าสหสัมพันธ์เชิงส่วนระหว่างตัวแปร Y กับ X_2 โดยถือว่า X_1 คงที่
- $r_{YX_3 \cdot X_1}$ หรือ $r_{Y3.1}$ คือค่าสหสัมพันธ์เชิงส่วนระหว่างตัวแปร Y กับ X_3 โดยถือว่า X_1 คงที่
- $r_{YX_2 \cdot X_3}$ หรือ $r_{Y2.3}$ คือค่าสหสัมพันธ์เชิงส่วนระหว่างตัวแปร Y กับ X_2 โดยถือว่า X_3 คงที่

$$r_{12.3} = \frac{r_{12} - r_{13} r_{23}}{[(1 - r_{13}^2)(1 - r_{23}^2)]^{1/2}}$$

เรียก $r_{12.3}$ ว่าเป็น Partial Correlation ค่าของ r และเครื่องหมายมีความหมายเช่นเดียวกับกรณีค่าสหสัมพันธ์อย่างง่าย

$$r_{Y1.2} = \frac{r_{Y1} - r_{Y2} r_{12}}{[(1 - r_{Y2}^2)(1 - r_{12}^2)]^{1/2}}$$

$$r_{Y3.12} = \frac{r_{Y3.1} - r_{Y2.1} r_{32.1}}{[(1 - r_{Y2.1}^2)(1 - r_{32.1}^2)]^{1/2}}$$

ค่าสัมประสิทธิ์สหสัมพันธ์เชิงซ้อน (Multiple Correlation Coefficient) โดยทั่วไปเป็นค่าสหสัมพันธ์ระหว่าง Y และ X_1, X_2, \dots, X_k หรือเป็นค่าสหสัมพันธ์อย่างง่ายระหว่าง Y กับ \hat{Y} ให้ R

$$R = r_{Y\hat{Y}} = \frac{\text{Cov}(Y, \hat{Y})}{[V(Y) V(\hat{Y})]^{1/2}} ; Y = \hat{Y} + e$$

$$R^2 = r_{Y\hat{Y}}^2$$

$$\begin{aligned} \text{ค่าของ Cov}(Y, \hat{Y}) &= \sum y \hat{y} \\ &= \sum \hat{y}^2 \end{aligned}$$

$$\text{ดังนั้น } r_{Y\hat{Y}} = \frac{\sum \hat{y}^2}{[\sum y^2 \sum \hat{y}^2]^{1/2}}$$

$$R^2 = \frac{\sum \hat{y}^2}{\sum y^2}$$

$$1 - R^2 = 1 - \frac{\sum \hat{y}^2}{\sum y^2}$$

การทดสอบนัยสำคัญของค่าสัมประสิทธิ์เชิงเส้นโดยการตั้งสมมติฐาน

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0$$

แล้วใช้ F - test ซึ่งมีค่า

$$F = \frac{(n - k - 1) R^2}{k (1 - R^2)} \quad \text{ใช้} \quad d.f = (k, n - k - 1)$$

สมมติว่า Y เกี่ยวข้องกับตัวแปรอิสระที่สำคัญ ๓ ตัว เช่น X_2, X_4 และ X_5
ค่าของ Multiple Standard Error of Estimate

$$S_{Y.245} = \sqrt{\frac{\sum (Y - \hat{Y})^2}{(n-k-1)}} = \sqrt{\frac{\sum (Y - \hat{Y})^2}{10}}$$

n = จำนวนข้อมูล = ๑๔

k = จำนวนตัวแปรอิสระ = ๓

๓.๓.๓ การสร้างสมการถดถอยเชิงเส้นโดยคอมพิวเตอร์

ในการหาโมเดลที่เหมาะสมสำหรับปริมาณการผลิตน้ำตาลทรายและปริมาณการบริโภคน้ำตาลทรายภายในประเทศ จะอาศัยวิธี Stepwise Multiple Regression โดยเลือกตัวแปรอิสระ ที่มีอิทธิพลต่อตัวแปรตามมากที่สุด และรอง ๆ ลงมาเข้ามาในสมการเพื่อได้สมการที่แสดงค่าประมาณของตัวแปรตามได้ ใกล้เคียงความเป็นจริงมากที่สุด

วิธีการเลือกตัวแปรอิสระเข้าในสมการหรือโมเดลกระทำดังนี้ :-

- พิจารณาเลือกค่า X_i เข้าในสมการทีละตัวตามลำดับความสัมพันธ์ระหว่างตัวแปรตามและตัวแปรอิสระ X_i กล่าวคือจะต้องเป็น X_i (ตัวแปรอิสระ) ตัวที่ให้ค่าดัชนีกำหนด (r^2) สูงสุดเข้าเป็นตัวแรก

เช่น สมมติให้ X_1 เป็นตัวที่ให้ค่าสัมประสิทธิ์กำหนดสูงสุดคือ $V_1 = r_{YX_1}^2$ มีค่าสูงสุดถึงนั้น

จะได้สมการในขั้นแรกเป็น $\hat{Y} = a^I + b_1^I X_1$

๒. ตัวแปรอิสระตัวที่สองที่เลือกเข้าไปในโมเดล คือตัวแปรอิสระที่ให้ค่า ผลคูณระหว่างค่าสัมประสิทธิ์กำหนด ของตัวมันเองกับค่าสัมประสิทธิ์กำหนดที่เหลือจากตัวแปรอิสระตัวแรก สูงสุด สมมติว่าเลือกได้เป็น X_2 นั่นคือ X จะให้ค่า

$$V_2 = (1 - r_{YX_1}^2) r_{YX_2 \cdot X_1}^2 \text{ สูงสุด}$$

ดังนั้นจะได้สมการใหม่ในขั้นที่ ๒ เป็น $\hat{Y} = a^{II} + b_1^{II} X_1 + b_2^{II} X_2$

๓. ตัวแปรอิสระตัวต่อไปสมมติให้ X_3 นั่นคือ X_3 จะให้ค่า

$$V_3 = (1 - r_{YX_1}^2) (1 - r_{YX_2 \cdot X_1}^2) r_{YX_3 \cdot X_1 X_2}^2 \text{ สูงสุด}$$

ได้สมการใหม่เป็น $\hat{Y} = a^{III} + b_1^{III} X_1 + b_2^{III} X_2 + b_3^{III} X_3$

ทำเช่นนี้ต่อไปจนครบ k ตัวจะได้

$$V_k = (1 - r_{YX_1}^2) (1 - r_{YX_2 \cdot X_1}^2) \dots (1 - r_{YX_{k-1} \cdot X_1 X_2 \dots X_{k-2}}^2) r_{YX_k \cdot X_1 X_2 \dots X_{k-1}}^2$$

โดยมีสมการเป็น

$$\hat{Y} = a^{(k)} + b_1^{(k)} X_1 + b_2^{(k)} X_2 + \dots + b_k^{(k)} X_k$$

เนื่องจากตัวแปรอิสระทุกตัวไม่ได้มีอิทธิพลอย่างมีนัยสำคัญต่อโมเดลทุกตัว ดังนั้นโมเดลในขั้นสุดท้าย ซึ่งมีตัวแปรอิสระ k ตัวจึงอาจจะไม่ใช่โมเดลที่มีนัยสำคัญเสมอไป วิธีการที่จะดีกว่าสมการในขั้นตอนใดควรจะต้องได้ จะยังคงใช้หลักในการพิจารณาเช่นเดียวกับกรณีการเลือกหาความโน้มตามลำดับเวลาดกล่าวคือ :-

๑. พิจารณาว่าสมการขั้นใดที่ให้ค่า Standard error ที่ต่ำสุดและ ค่าสัมประสิทธิ์กำหนดสูงสุด (เมื่อเพิ่มตัวแปรอิสระเข้าไปอีกค่า S.E. จะเพิ่มขึ้น)

๒. พิจารณาจากการทดสอบนัยสำคัญของสมการ โดยพิจารณา F-test จาก ANOVA

ANOVA

๓. พิจารณาจากการทดสอบนัยสำคัญของ Partial Regression
Coefficient โดย t-test