

ระเบียบวิธีวิจัย

ในการคัดเลือกพื้นที่เพื่อประกาศเขตปฏิรูปที่ดินนั้น ผู้วิจัยได้นำข้อมูลที่เก็บรวบรวมได้มาทำการวิเคราะห์ตามระเบียบวิธีทางสถิติโดยใช้โปรแกรมสำเร็จรูป SPSS (Statistical Package for the Social Science) จะสร้างแบบจำลองที่ใช้แสดงสภาพการพัฒนาในแต่ละอำเภอ โดยพิจารณาจากตัวแปรต่าง ๆ ทางด้านเศรษฐกิจและสังคม ซึ่งได้กล่าวถึงตัวแปร และความหมายของตัวแปรไว้ในบทที่ 2 แล้ว ส่วนวิธีการวิเคราะห์ที่ใช้มีดังนี้

3.1 การคัดเลือกตัวแปร เพื่อจะนำไปใช้ในการวิเคราะห์

ในขั้นแรกจะคัดเลือกตัวแปรที่จะนำไปใช้ในการวิเคราะห์ด้วยวิธีทางสถิติต่าง ๆ โดยพิจารณาตัวแปรที่มีความสำคัญที่จะอธิบายถึงการจำแนกพื้นที่เพื่อประกาศเขตปฏิรูปที่ดินได้ ซึ่งจะพิจารณาจากความสัมพันธ์ของตัวแปรแต่ละตัวในเมตริกความสัมพันธ์ (Correlation Matrix) เพื่อขจัดอิทธิพลของพหุสัมพันธ์ในตัวแปรอิสระ (Multicollinearity) ของตัวแปรอิสระออกไปบ้าง ได้กำหนดสัญลักษณ์ของตัวแปรต่าง ๆ ดังนี้

- Y_i = รายได้เฉลี่ยต่อครัวเรือนต่อปีของประชากรในอำเภอที่ i (บาท/ครัวเรือน/ปี)
- X_1 = เนื้อที่ถือครองเพื่อการเกษตรเฉลี่ยต่อครัวเรือน (ไร่/ครัวเรือน)
- X_2 = ร้อยละของเนื้อที่ถือครองเพื่อการเกษตร
- X_3 = ร้อยละของเนื้อที่เช่า
- X_4 = ร้อยละของผู้ถือครองที่มีการเช่า
- X_5 = ร้อยละของผู้ถือครองที่ปราศจากที่ดิน
- X_6 = ร้อยละของผู้ถือครองที่มีขนาดการถือครองตั้งแต่ 100 ไร่
- X_7 = ร้อยละของเนื้อที่ชลประทาน

- X_8 = ร้อยละของเนื้อที่ที่ปลูกพืชได้มากกว่า 1 ครั้ง
 X_9 = เนื้อที่ปลูกข้าวเฉลี่ยต่อครัวเรือน
 X_{10} = ร้อยละของครัวเรือนที่ปลูกข้าว
 X_{11} = ร้อยละของครัวเรือนที่มีอาชีพหลักทางการเกษตร
 X_{12} = ร้อยละของครัวเรือนที่มีรายได้เฉลี่ยต่อปีตั้งแต่ 5,168 บาทขึ้นไป
 X_{13} = ร้อยละของครัวเรือนที่นำที่ดินไปจ้างหรือขายฝากพ่อค้าเอกชน
 X_{14} = ราคาที่ดินเฉลี่ยต่อไร่
 X_{15} = จำนวนธนาคารข้าว
 X_{16} = จำนวนธนาคารโคกระบือ
 X_{17} = จำนวนกลุ่มเกษตรกร
 X_{18} = จำนวนกลุ่มยุวเกษตรกร
 X_{19} = จำนวนกลุ่มสหกรณ์ต่าง ๆ
 X_{20} = จำนวนกลุ่มชลประทานราษฎร์
 X_{21} = จำนวนธนาคารต่าง ๆ ได้แก่ ธนาคารออมสิน ธนาคารเพื่อการเกษตรและสหกรณ์การเกษตร ธนาคารพาณิชย์
 X_{22} = จำนวนโรงสีข้าวขนาดใหญ่
 X_{23} = จำนวนโรงสีข้าวขนาดกลาง
 X_{24} = จำนวนโรงสีข้าวขนาดเล็ก
 X_{25} = อัตราการเพิ่มของประชากร
 X_{26} = ความหนาแน่นของประชากร
 X_{27} = จำนวนนักเรียนต่อครู 1 คน
 X_{28} = จำนวนประชากรในอำเภอเฉลี่ยต่อสถานีอนามัย 1 โรง

- X_{29} = ร้อยละของครัวเรือนที่มีไฟฟ้าใช้
 X_{30} = ร้อยละของครัวเรือนที่มีเครื่องรับวิทยุ
 X_{31} = ร้อยละของครัวเรือนที่มีเครื่องรับโทรทัศน์
 X_{32} = ผลผลิตข้าวเฉลี่ยต่อไร่ (กก./ไร่)
 X_{33} = เนื้อที่เข้าเฉลี่ยต่อครัวเรือน (ไร่/ครัวเรือน)
 X_{34} = ร้อยละของผู้ถือครองที่เป็นเจ้าของที่ดิน

โดยจะได้แปลงข้อมูลทั้งหมดให้อยู่ในรูปของค่ามาตรฐาน Z (Standardized data) เพื่อให้ข้อมูลอยู่ในหน่วยเดียวกัน และสะดวกในการเปรียบเทียบ

$$ZY = \frac{Y - \mu_Y}{\sigma_Y}$$

$$ZX_i = \frac{X_i - \mu_{X_i}}{\sigma_{X_i}}$$

ในทางปฏิบัติจะไม่ทราบค่า μ_Y , μ_{X_i} , σ_Y และ σ_{X_i} ดังนั้นจากตัวอย่างที่ได้

จะได้ค่าประมาณของ μ_Y , μ_{X_i} , σ_Y และ σ_{X_i} เป็น \bar{Y} , \bar{X}_i , S_Y และ S_{X_i} ตามลำดับ

โดยที่

$$\bar{Y} = \frac{1}{n} \sum_{j=1}^n Y_j$$

$$\bar{X}_i = \frac{1}{n} \sum_{j=1}^n X_{ij}$$

$$S_Y = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (Y_j - \bar{Y})^2}$$

$$S_{X_i} = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (X_{ij} - \bar{X}_i)^2}$$

ในการคัดเลือกตัวแปร เพื่อจะใช้ในการวิเคราะห์ จะพิจารณาเป็น 2 ตอนคือ

- 1) พิจารณาค่าความสัมพันธ์ของตัวแปรตาม (Dependent Variable ; ZY) กับตัวแปรอิสระใด ๆ (Independent Variable ; ZX_i) ถ้าค่าสัมบูรณ์ของค่าความสัมพันธ์ระหว่าง ZY กับ ZX_i ใด ๆ $|r_{ZY, ZX_i}| < .10$ จะตัด ZX_i นั้นออกจากการวิเคราะห์ เพราะถือว่า ZX_i ตัวนั้นมีความสัมพันธ์กับ ZY น้อยมาก
- 2) พิจารณาค่าความสัมพันธ์ของตัวแปรอิสระด้วยกัน (Independent Variables ; ZX_i, ZX_j) จะพิจารณาจากเมตริกความสัมพันธ์¹ (Examination of Correlation Matrix) ซึ่งเป็นวิธีหาค่าความสัมพันธ์ในตัวแปรอิสระวิธีหนึ่ง โดยดูที่ค่าของความสัมพันธ์ของ ZX_i กับ ZX_j ในเมตริกความสัมพันธ์ ถ้าตัวแปรอิสระ ZX_i กับ ZX_j ใดใดมีค่าสัมบูรณ์ของค่าความสัมพันธ์ $|r_{ZX_i, ZX_j}|$ เข้าใกล้ 1 ถือว่าตัวแปรอิสระคู่นั้นมีความสัมพันธ์กันค่อนข้างสูง จะใช้ตัวแปรอิสระตัวใดตัวหนึ่งแทนโดยตัดอีกตัวหนึ่งออกไปจากการวิเคราะห์ได้ ในการวิจัยครั้งนี้ ถ้า $|r_{ZX_i, ZX_j}| > .70$ จะตัดตัวแปรอิสระตัวใดตัวหนึ่งออกไป และยังสามารถถือว่าค่าสัมประสิทธิ์ของการตัดสินใจ (Coefficient of Determination; r_{ZX_i, ZX_j}^2) ซึ่งเป็นค่าแสดงอิทธิพลของ ZX_i, ZX_j ที่มีต่อกันถึงประมาณ 50% ขึ้นไป และเป็นค่าที่มากพอ

การทดสอบการกระจายของข้อมูล (Test for Goodness of fit) ว่า
ข้อมูลมีการกระจายแบบปกติ และมีค่าเฉลี่ย = 0 มีความแปรปรวน = 1 โดยใช้การทดสอบโคโมโกรอฟ-สมีนอฟ (Kolmogorov-Smirnov Test)²

คุณสมบัติของการทดสอบโคโมโกรอฟ-สมีนอฟ

- (1) ข้อมูลที่จะใช้ทดสอบ จะต้องเป็นข้อมูลเชิงปริมาณ สามารถให้ลำดับได้ (Ordinal)
- (2) ทราบฟังก์ชันความน่าจะเป็นของข้อมูล (probability density function ; pdf.) และ pdf. เป็นแบบต่อเนื่อง และถ้าไม่ทราบค่าพารามิเตอร์ จะประมาณ

¹Douglas C. Montgomery, Introduction to linear regression analysis, (New York : John Wiley & Son, 1981), p. 296-299.

²Sidney Seigel, Nonparametric Statistics for Behavioral Sciences, (Japan : McGraw Hill. Kogakusha, LTD., 1956). p. 47-52.

ค่าพารามิเตอร์

(3) ใช้เมื่อจำนวนตัวอย่างมีน้อย

หลักการของวิธีการนี้คือ วิธีการทดสอบโคโมโกรอฟ-ส്മิโนฟ เป็นการเปรียบเทียบระหว่างความถี่สะสม ที่ได้ทางทฤษฎีกับการที่ได้จากการสังเกตในกรณีที่มีสมมติฐานดังนี้

H_0 : ข้อมูลที่ใช้มีการกระจายแบบปกติ

$$\text{ตัวสถิติที่ใช้ทดสอบคือ } D = \text{maximum} | F_0(X) - S_n(X) |$$

โดยที่ $F_0(X)$ แทนฟังก์ชันความน่าจะเป็นสะสมของ X ภายใต้ H_0

$S_n(X)$ แทนความถี่สะสมของ X ในการสุ่มตัวอย่าง n ครั้ง

$$= \frac{k}{n} \text{ เมื่อ } k = \text{จำนวนครั้งที่สังเกตได้}$$

ซึ่งจะมีค่าน้อยกว่าหรือเท่ากับ X

อาณาเขตวิกฤตที่ระดับนัยสำคัญ α คือ $D \geq D_{n,\alpha}$

เมื่อ $D_{n,\alpha}$ คือค่าที่เปิดจากตาราง

n คือจำนวนตัวอย่าง

α คือระดับนัยสำคัญ

ดังนั้น จะปฏิเสธ H_0 เมื่อ $D \geq D_{n,\alpha}$

3.2 วิธีวิเคราะห์การถดถอยเชิงซ้อน

จากแนวความคิดเรื่องการวัดความยากจนและการศึกษาของธนาคารโลกเกี่ยวกับความยากจน ซึ่งใช้รายได้ของประชากรมาเป็นตัววัดสภาพความเป็นอยู่ของประชากร จะกำหนดน้ำหนักของตัวแปรต่าง ๆ ที่จะนำมาสร้างแบบจำลองด้วยการใช้วิธีวิเคราะห์การถดถอยเชิงซ้อน โดยให้รายได้เฉลี่ยต่อครัวเรือนต่อปีของประชากรในอำเภอเป็นตัวแปรตาม ซึ่งมีรูปแบบดังนี้

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \epsilon$$



โดยที่ β_i เป็นค่าสัมประสิทธิ์การถดถอยเชิงซ้อน ไม่ทราบค่าของตัวแบบ

$$i = 0, 1, 2, \dots, k$$

และถ้าหากว่ามีจำนวนข้อมูล n ตัว ต่อตัวแปรแต่ละตัวแล้วจะได้รูปของข้อมูลเป็นดังนี้

คือ

ข้อมูล	ตัวแปรตาม	ตัวแปรอิสระที่ 1	ที่ 2	...	ที่ k
1	Y_1	X_{11}	X_{21}		X_{k1}
2	Y_2	X_{12}	X_{22}		X_{k2}
3	Y_3	X_{13}	X_{23}		X_{k3}
.
.
.
n	Y_n	X_{1n}	X_{2n}		X_{kn}

จากข้อมูลข้างต้นนี้ หากนำมาเขียนในรูปของสมการเมตริกจะได้เป็น

$$\begin{pmatrix} Y_1 \\ Y_2 \\ Y_3 \\ \cdot \\ \cdot \\ \cdot \\ Y_n \end{pmatrix} = \begin{pmatrix} 1 & X_{11} & X_{21} & \cdot & \cdot & \cdot & X_{k1} \\ 1 & X_{12} & X_{22} & \cdot & \cdot & \cdot & X_{k2} \\ 1 & X_{13} & X_{23} & \cdot & \cdot & \cdot & X_{k3} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 1 & X_{1n} & X_{2n} & \cdot & \cdot & \cdot & X_{kn} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \cdot \\ \cdot \\ \cdot \\ \beta_k \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \cdot \\ \cdot \\ \cdot \\ \epsilon_n \end{pmatrix}$$

หรือ $Y = X\beta + \epsilon$

- โดยที่ Y คือข้อมูลของตัวแปรตาม โดยมีเวกเตอร์อันดับที่ $n \times 1$
- X คือข้อมูลของตัวแปรอิสระ โดยมีเมตริกอันดับที่ $n \times (k+1)$
- β คือสัมประสิทธิ์ของการถดถอย โดยมีเวกเตอร์อันดับที่ $(k+1) \times 1$
- ε คือค่าความคลาดเคลื่อนของค่าประมาณจากค่าจริง โดยมีเวกเตอร์อันดับที่ $n \times 1$

3.2.1 ข้อสมมติเบื้องต้น

ในการที่เราจะศึกษาถึงวิธีการประมาณและคุณสมบัติของรูปแบบ ซึ่งเราไม่ทราบค่า นั้น เราจะต้องตั้งข้อสมมติ ซึ่งจำเป็นสำหรับการประมาณเสียก่อน และข้อสมมติที่จะกล่าวถึงต่อไปนี้ก็มีความสำคัญมาก เพราะคุณสมบัติของตัวประมาณโดยวิธีกำลังสองน้อยที่สุด (Least Squares Method) ซึ่งเป็นวิธีที่เรานำมาใช้ จะขึ้นอยู่กับข้อสมมติเหล่านั้น ซึ่งจะต้องนำมาพิจารณาประกอบก่อนลงมือทำการวิเคราะห์คือ

$$1) \quad Y = X\beta + \varepsilon$$

นั่นคือ Y_i เป็นฟังก์ชันเชิงเส้นของ X_{ij} กับ ε_i โดยที่ $i = 1, 2, \dots, n$
 $j = 0, 1, \dots, k \quad X_{0i} = 1$

$$2) \quad E(\varepsilon) = 0$$

นั่นคือ $E(\varepsilon) = 0$ สำหรับทุกค่า i

$$3) \quad E(\varepsilon_i \varepsilon_j) = \sigma^2 I$$

นั่นคือ $E(\varepsilon_i^2) = \sigma^2$ สำหรับทุกค่า i นั่นคือ ε_i มีความแปรปรวนคงที่

$$E(\varepsilon_i \varepsilon_j) = 0 \quad \text{เมื่อ } i \neq j \quad \text{แสดงว่าแต่ละคู่ของ } \varepsilon \text{ ไม่เกี่ยวข้องกัน}$$

(pairwise uncorrelated)

4) ค่าความคลาดเคลื่อนของค่าประมาณจากค่าจริง ε_i เป็นตัวแปรเชิงสุ่มที่มีการกระจายแบบปกติ $N(0, \sigma^2)$ ซึ่งข้อนี้จำเป็นต้องใช้ในกรณีที่ต้องการสร้างช่วงความเชื่อมั่นของ $\beta_0, \beta_1, \beta_2, \dots, \beta_k$ หรือเมื่อต้องการทดสอบสมมติฐานใด ๆ เกี่ยวกับตัวประมาณ

5) X เป็นเมตริกอันดับที่ $n \times (k+1)$ ซึ่งมีค่าคงที่ สิ่งทำให้ X และ ϵ เป็นตัวแปรอิสระกัน เราจึงได้ว่า

$$E(\epsilon/X) = E(\epsilon) = 0$$

$$E(X'\epsilon) = X'E(\epsilon) = 0$$

$$E(\epsilon\epsilon'/X) = E(\epsilon\epsilon') = \sigma^2 I$$

6) X มี rank $k + 1 < n$ หมายความว่าจำนวนค่าสังเกต n จะมากกว่าจำนวนพารามิเตอร์ $(k+1)$ ที่จะต้องประมาณค่า

จากข้อสมมติดังกล่าวมาแล้วข้างต้นเราจะเห็นได้ว่า

$$1) E(Y/X) = X\beta + E(\epsilon/X) = X\beta \text{ เช่น}$$

$$E(Y_i/X_{1i}, X_{2i}, \dots, X_{ki}) = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki}$$

$$2) E[(Y - E(Y))(Y - E(Y))'/X] = E(\epsilon\epsilon'/X) = \sigma^2 I$$

$$\text{เช่น } E[(Y_i - E(Y_i))(Y_i - E(Y_i))'/X_{1i}, X_{2i}, \dots, X_{ki}] = \sigma^2 \text{ สำหรับทุกค่า } i$$

3.2.2 การประมาณค่าพารามิเตอร์

วิธีการในการประมาณค่าพารามิเตอร์การถดถอยนั้น กระทำได้ 2 วิธีคือ

- 1) วิธีการสังคมน้อยที่สุด
- 2) วิธีน่าจะเป็นสูงสุด (Maximum Likelihood Method)

ทั้ง 2 วิธีดังกล่าวเป็นวิธีในการคำนวณหาค่าเส้นถดถอยจากข้อมูลที่กำหนดให้หรือกล่าวอีกนัยหนึ่งก็คือ ถ้ากำหนดสมการ

$$Y = X\beta + \epsilon \text{ แล้วเทคนิคทั้งหมดข้างต้นก็จะเป็นเครื่องมือในการคำนวณ}$$

$$\hat{Y} = X\hat{\beta} \text{ นั้นเอง หากแต่ว่าแต่ละวิธีก็มีคุณสมบัติของตัวเอง}$$

แต่ในที่นี้จะนำมากล่าวถึงเพียงวิธีเดียวเท่านั้นคือ วิธีที่ 1 ซึ่งถือว่าเป็นวิธีที่สำคัญและนิยมใช้กันมาก โดยมีหลักการของการประมาณแบบกำลังสองน้อยที่สุด คือการทำให้ผลรวมกำลังสองของส่วนเบี่ยงเบนของค่าสังเกตต่างจากค่าเฉลี่ยมีค่าน้อยที่สุด และจะทำให้ได้ตัวประมาณที่ไม่เอนเอียง

$$\text{จากรูปแบบ } Y_i = X_i\beta + \epsilon_i$$

$$\text{และ } E(Y_i) = X_i\beta \quad \therefore \quad E(\epsilon_i) = 0$$

ซึ่ง $E(Y_i)$ คือค่าเฉลี่ยของ Y_i

$$\begin{aligned} \text{ดังนั้น } \epsilon_i &= Y_i - E(Y_i) \\ &= Y_i - X_i\beta \end{aligned}$$

โดยวิธีกำลังสองน้อยที่สุด จะต้องหาว่ากำลังสองของค่าความคลาดเคลื่อนเท่ากับ

$$\begin{aligned} \epsilon'\epsilon &= (Y - X\beta)'(Y - X\beta) \\ &= Y'Y - \beta'XY - Y'X\beta + \beta'XX\beta \\ &= Y'Y - 2\beta'XY + \beta'XX\beta \end{aligned}$$

(เนื่องจาก $Y'X\beta$ เป็นสเกลล่า ดังนั้น ทรานสโพสของมันจึงมีค่าเท่าเดิม
จึงได้ว่า $Y'X\beta = \beta'XY$)

จากนั้นจึงเริ่มหาตัวประมาณโดยทำให้กำลังสองของค่าความคลาดเคลื่อนที่ได้มานั้นมีค่าน้อยที่สุดโดยใช้ การดิฟเฟอเรนเชียล

$$\frac{\partial (\epsilon'\epsilon)}{\partial \beta} = \frac{\partial (Y'Y - 2\beta'XY + \beta'XX\beta)}{\partial \beta} = 0$$

$$- 2X'Y + 2XX\beta = 0$$

$$X'Y = X'X\beta$$

จากนี้จะหาตัวประมาณของ β คือ b

ซึ่งเป็นชุดสมการปกติ (normal equations) และถ้าหากจะนำมาเขียนในรูปสเกลล่าจะได้ว่า

$$\begin{aligned}\Sigma Y_i &= b_0 n + b_1 \Sigma X_{1i} + b_2 \Sigma X_{2i} + \dots + b_k \Sigma X_{ki} \\ \Sigma X_{1i} Y_i &= b_0 \Sigma X_{1i} + b_1 \Sigma X_{1i}^2 + b_2 \Sigma X_{1i} X_{2i} + \dots + b_k \Sigma X_{1i} X_{ki} \\ \Sigma X_{2i} Y_i &= b_0 \Sigma X_{2i} + b_1 \Sigma X_{1i} X_{2i} + b_2 \Sigma X_{2i}^2 + \dots + b_k \Sigma X_{2i} X_{ki} \\ &\vdots \\ &\vdots \\ &\vdots \\ \Sigma X_{ki} Y_i &= b_0 \Sigma X_{ki} + b_1 \Sigma X_{1i} X_{ki} + b_2 \Sigma X_{2i} X_{ki} + \dots + b_k \Sigma X_{ki}^2\end{aligned}$$

การจะหา b_0 ได้จะต้องมีเงื่อนไขว่าชุดสมการปกติทั้งหมด $k+1$ สมการนี้ต้องเป็นอิสระซึ่งกันและกันโดยที่ $X'X$ จะต้องเป็นเมตริกที่ไม่เอกเทศ เพราะจะทำให้สามารถหา $(X'X)^{-1}$ ได้

$$b_0 = (X'X)^{-1} X'Y$$

และเมื่อทำการดิฟเฟอเรนเชียลอันดับที่ 2 จะได้

$$\frac{\partial (\Sigma \epsilon^2)}{\partial \beta} = 2X'X$$

ซึ่งเป็นเดฟนิททางบวกเมื่อ X มี rank เต็มตามจำนวนคอลัมน์คือ $k+1$

เพราะฉะนั้น $b_0 = (X'X)^{-1} X'Y$ จึงเป็นตัวประมาณที่ได้จากการทำให้กำลังสองของความคลาดเคลื่อนมีค่าน้อยที่สุด

$$\text{จึงได้ว่า } \hat{Y}_i = X_i b_0$$

ต่อไปจะพิจารณาค่าความแปรปรวนของ b_0

$$\begin{aligned}\text{จาก } b_0 &= (X'X)^{-1} X'Y \\ &= (X'X)^{-1} X'(X\beta + \epsilon) \\ &= \beta + (X'X)^{-1} X'\epsilon \\ b_0 - \beta &= (X'X)^{-1} X'\epsilon\end{aligned}$$

$$\begin{aligned}
E(\hat{\beta} - \beta)(\hat{\beta} - \beta)' &= E \left[(X'X)^{-1} X'\varepsilon \right] \left[(X'X)^{-1} X'\varepsilon \right]' \\
&= E \left[(X'X)^{-1} X'\varepsilon\varepsilon' X (X'X)^{-1} \right] \\
&= (X'X)^{-1} X' E(\varepsilon\varepsilon') X (X'X)^{-1} \\
&= \sigma^2 (X'X)^{-1} \quad (\because E(\varepsilon\varepsilon') = \sigma^2 I)
\end{aligned}$$

นั่นคือ $\text{Var}(\hat{\beta}) = \sigma^2 (X'X)^{-1}$

จากข้างต้นจึงได้ว่า $E(\hat{\beta} - \beta)(\hat{\beta} - \beta)'$ เป็นวาเรียนซ์-โควาเรียนซ์ (Variance-Covariance) ของ $\hat{\beta}$ แต่เนื่องจากไม่ทราบค่าของ σ^2 จึงต้องประมาณค่าของ σ^2 ก่อน

จาก $y = X\beta + \varepsilon$

และ $\hat{y} = X\hat{\beta}$

ให้ $e = y - \hat{y}$

$$= y - X\hat{\beta}$$

$$= y - X(X'X)^{-1} X'y$$

$$= (I - X(X'X)^{-1} X') y$$

$$= My \quad (\text{ให้ } M = I - X(X'X)^{-1} X')$$

$$= M(X\beta + \varepsilon)$$

$$= MX\beta + M\varepsilon$$

$$= M\varepsilon$$

ทั้งนี้เนื่องจาก $MX = 0$

$$\begin{aligned}
 \text{โดย } MX &= (I - X (X'X)^{-1} X') X \\
 &= X - X (X'X)^{-1} X'X \\
 &= 0
 \end{aligned}$$

$$\begin{aligned}
 \text{หา } \hat{\epsilon} &= \hat{\epsilon}' M M \epsilon \\
 &= \hat{\epsilon}' M \epsilon
 \end{aligned}$$

$$\begin{aligned}
 \text{ทั้งนี้เนื่องจาก } MM &= (I - X(X'X)^{-1} X') (I - X(X'X)^{-1} X') \\
 &= I - X(X'X)^{-1} X' - X(X'X)^{-1} X' + X(X'X)^{-1} X' X (X'X)^{-1} X' \\
 &= I - X(X'X)^{-1} X' \\
 &= M
 \end{aligned}$$

$$\therefore \hat{\epsilon} \hat{\epsilon}' = \text{tr} \hat{\epsilon}' M \epsilon$$

ทั้งนี้เนื่องจาก M เป็นสเกลาร์จึงเท่ากับ Trace ของตัวมันเอง

$$\hat{\epsilon} \hat{\epsilon}' = \text{tr} M \hat{\epsilon} \hat{\epsilon}'$$

$$\therefore \text{tr}(AB) = \text{tr}(BA)$$

$$\begin{aligned}
 \therefore E(\hat{\epsilon} \hat{\epsilon}') &= E(\text{tr} M \hat{\epsilon} \hat{\epsilon}') \\
 &= \text{tr} M \cdot E(\hat{\epsilon} \hat{\epsilon}')
 \end{aligned}$$

เพราะว่า trace เป็นฟังก์ชันเชิงเส้น

$$\begin{aligned}
 E(\hat{\epsilon} \hat{\epsilon}') &= \sigma^2 \text{tr} M \quad (\dots E(\hat{\epsilon} \hat{\epsilon}') = \sigma^2 I) \\
 &= \sigma^2 (n-k-1)
 \end{aligned}$$

$$\begin{aligned}
\text{ทั้งนี้เนื่องจาก } \text{tr } M &= \text{tr} (I - X(X'X)^{-1} X') \\
&= \text{tr} (I_n) - \text{tr} (X(X'X)^{-1} X') \\
&= n - \text{tr} X'X(X'X)^{-1} \quad (\text{tr} (AB) = \text{tr} (BA)) \\
&= n - \text{tr} (I_{k+1}) \\
&= n - k - 1
\end{aligned}$$

$$\begin{aligned}
\text{ดังนั้นจึงได้ว่า } s^2 &= \frac{e'e}{n - k - 1} \\
&= \frac{Y' M Y}{n - k - 1} \quad \text{ซึ่งเป็นตัวประมาณที่ไม่เอนเอียงของ } \sigma^2
\end{aligned}$$

และ $S^2(X'X)^{-1}$ ก็เป็นตัวประมาณที่ไม่เอนเอียงของ $\sigma^2(X'X)^{-1}$

3.2.3 คุณสมบัติของตัวประมาณโดยวิธีกำลังสองน้อยที่สุด

จากข้อสมมติที่กำหนดไว้ในหัวข้อข้างต้น ถ้าหากเป็นจริงแล้ว จะได้ว่า การหาค่าตัวประมาณโดยวิธีกำลังสองน้อยที่สุดจะมีคุณสมบัติดังนี้คือ

$$1. \quad b = (X'X)^{-1} X'Y \text{ เป็นตัวประมาณที่ไม่เอนเอียงของ } \beta$$

พิสูจน์

$$\begin{aligned}
b &= (X'X)^{-1} X'Y \\
&= (X'X)^{-1} X'(X\beta + \epsilon) \\
&= (X'X)^{-1} X'X\beta + (X'X)^{-1} X'\epsilon \\
&= \beta + (X'X)^{-1} X'\epsilon
\end{aligned}$$

$$E(b) = \beta$$

นั่นคือ จะได้ว่า $b = (X'X)^{-1} X'Y$ เป็นตัวประมาณที่ไม่เอนเอียงของ β

2. $\hat{\beta} = (X'X)^{-1} X'Y$ จะเป็นตัวประมาณที่ไม่เอนเอียงเชิงเส้นที่ดีที่สุด (Best Linear Unbiased Estimator) ในความหมายคือ ตัวประมาณอื่น ๆ ของ β ซึ่งเป็นเชิงเส้นกับ Y และไม่เอนเอียงเช่นกัน จะมีค่าความแปรปรวน-โควาเรียนซ์เมตริก ซึ่งมีค่าเกินกว่าค่าความแปรปรวน-โควาเรียนซ์ของตัวประมาณ $\hat{\beta}$ หรืออาจกล่าวได้โดยง่าย ๆ ก็คือในบรรดาตัวประมาณในจำพวกของตัวประมาณที่ไม่เอนเอียงเชิงเส้นด้วยกันแล้วตัวประมาณโดยวิธีกำลังสองน้อยที่สุดนั้นจะมีค่าความแปรปรวนต่ำสุด โดยทฤษฎีของกอลด์-มาคอฟฟ์และเป็นที่ยอมรับกันอย่างแพร่หลาย

จากการที่ตัวประมาณของ β เป็นเชิงเส้นกับ Y เราจึงสามารถเขียนตัวประมาณนั้นได้เป็น $A^* Y$ โดย A^* เป็นเมตริก $(k+1) \times n$ ซึ่งไม่ขึ้นอยู่กับค่าของ Y

$$\begin{aligned} \text{ให้ } A^* &= (X'X)^{-1} X' + A \\ \therefore A^* Y &= \left[(X'X)^{-1} X' + A \right] Y \\ &= \left[(X'X)^{-1} X' + A \right] (X\beta + \epsilon) \\ &= [I + AX] \beta + \left[(X'X)^{-1} X' + A \right] \epsilon \\ E(A^* Y) &= (I + AX) \beta + 0 \\ &= \beta + AX \beta \end{aligned}$$

ตัวประมาณ $A^* Y$ จะเป็นตัวประมาณที่ไม่เอนเอียง ถ้า $AX = 0$

$$\begin{aligned} \therefore \text{จาก } A^* Y &= \beta + \left[(X'X)^{-1} X' + A \right] \epsilon \\ A^* Y - \beta &= \left[(X'X)^{-1} X' + A \right] \epsilon \\ (A^* Y - \beta) (A^* Y - \beta)' &= \left[(X'X)^{-1} X' + A \right] \epsilon \epsilon' \left[(X'X)^{-1} X' + A \right]' \\ E(A^* Y - \beta) (A^* Y - \beta)' &= \left[(X'X)^{-1} X' + A \right] E(\epsilon \epsilon') \left[X(X'X)^{-1} + A' \right] \\ &= \sigma^2 \left[(X'X)^{-1} X' + A \right] \left[X(X'X)^{-1} + A' \right] \\ &= \sigma^2 \left[(X'X)^{-1} + AX(X'X)^{-1} + (X'X)^{-1} X' A' + A A' \right] \\ &= \sigma^2 \left[(X'X)^{-1} + A A' \right] \end{aligned}$$



เนื่องจากเรากำหนดให้ $AX = 0 \rightarrow XA' = 0$ ด้วย

$$\begin{aligned} \text{จึงได้ว่า } \text{Var}(A^*X) &= \sigma^2 (X'X)^{-1} + \sigma^2 AA' \\ &= \text{Var}(b) + \sigma^2 AA' \end{aligned}$$

นั่นคือ วาเรียนซ์-โควาเรียนซ์ของตัวประมาณจะมีค่าเกิน วาเรียนซ์-โควาเรียนซ์ของตัวประมาณ ซึ่งประมาณโดยวิธีกำลังสองน้อยที่สุด

$$\rightarrow b = (X'X)^{-1} X'Y \quad \text{จะเป็นตัวประมาณที่ไม่เอนเอียงเชิงเส้นที่ดีที่สุด}$$

3.3 วิเคราะห์องค์ประกอบหลัก

เป็นวิธีที่ให้แบบแผนความสัมพันธ์ระหว่างปัจจัยต่าง ๆ และแยกปัจจัยเหล่านี้เป็นกลุ่ม ๆ ตามความสัมพันธ์ที่เกี่ยวข้อง ตัวประกอบที่จะให้ค่าถ่วงน้ำหนัก (Loading) ซึ่งสามารถแปลงคะแนนดิบของปัจจัยนั้น ๆ ให้อยู่ในรูปค่าถ่วงน้ำหนัก เพื่อจัดลำดับความสำคัญให้ได้ชัดเจนถูกต้องยิ่งขึ้น

หลักการของวิเคราะห์องค์ประกอบโดยทั่วไป คือหาตัวแปรใหม่ที่แสดงความสัมพันธ์เชิงเส้นตรง (Linear Combination) ของตัวแปรเดิมหลาย ๆ ตัวแปร

ขั้นตอนของการคำนวณหารูปแบบของตัวแปรใหม่ในองค์ประกอบ (component) ต่าง ๆ

1) กำหนดตัวแปรที่ใช้ในการศึกษาวัดความสำคัญของพื้นที่ ให้ X เป็นตัวแปรต่าง ๆ ที่ใช้วัด

$$X = \begin{pmatrix} X_{11} & X_{12} & \dots & X_{1k} \\ X_{21} & X_{22} & \dots & X_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n1} & X_{n2} & \dots & X_{nk} \end{pmatrix}$$

มีตัวแปร k ตัวเช่นเดียวกับการวิเคราะห์การถดถอยเชิงซ้อนและมีพื้นที่ n พื้นที่

2) หาค่าความแปรปรวนร่วม (Covariance) ของตัวแปรทั้งหมด

$$S = E \left[(X - E(X)) (X - E(X)) \right]$$

S = Covariance matrix ของ X

3) หาค่าค่าแรงแคเตอร์ลิสต์ติก รุท (Characteristic roots ; λ)

$$| S - \lambda I | = 0$$

I = Identity matrix

$$\lambda = \left[\lambda_1 \lambda_2 \dots \lambda_k \right]$$

λ_i = ค่าแรงแคเตอร์ลิสต์ติก รุทที่ i ของ S

$$\lambda_1 \geq \lambda_2 \geq \lambda_3 \dots \geq \lambda_k$$

4) หาค่าค่าแรงแคเตอร์ลิสต์ติก เวกเตอร์ (Characteristic vectors; a_i)

$$S a_i = \lambda_i a_i$$

a_i = ค่าแรงแคเตอร์ลิสต์ติก เวกเตอร์ที่ i ที่สอดคล้องกับค่าแรงแคเตอร์

ลิสต์ติก รุท λ_i

$$a_i = \begin{bmatrix} a_1 \\ a_2 \\ \cdot \\ \cdot \\ a_k \end{bmatrix}$$

5) หา normalized ของค่าแรงแคเตอร์ลิสต์ติก เวกเตอร์

$$\text{normalized} = \begin{pmatrix} a_1 / \sqrt{\text{lenght}} \\ a_2 / \sqrt{\text{lenght}} \\ \vdots \\ a_k / \sqrt{\text{lenght}} \end{pmatrix}$$

$$\text{lenght} = a_1^2 + a_2^2 + \dots + a_k^2$$

นำค่า normalized ของค่าแรงแคเตอร์สถิติเวคเตอร์ มาเรียงกันจะได้
พริ้นลิปอลคอมโพเนนท์เวคเตอร์ (Principal Component Vector)

6) หมุนแกนองค์ประกอบให้แกนนตั้งฉากซึ่งกันและกัน (Orthogonal) จะได้
องค์ประกอบที่เป็นอิสระต่อกัน ค่าที่ได้จากการหมุนแกนคือค่าองค์ประกอบ (Principal
Component) ที่จะนำไปใช้ในการวิเคราะห์เพื่อหาตัวแปรที่สำคัญ โดยพิจารณาจากองค์ประกอบ
ที่มีค่าความแปรปรวนสูงสุด

องค์ประกอบที่ i (Principal Component ที่ i)

$$Y_i = \sum_{j=1}^k a'_{ij} X_j = \sum_{j=1}^k a_{ij} X_j$$

$$Y_i = a_{i1} X_1 + a_{i2} X_2 + \dots + a_{ik} X_k$$

จะเลือกองค์ประกอบที่ให้ความแปรปรวนสูงสุด โดยพิจารณา

$$a'_{ii} a_{ii} = 1$$

$$a'_{ii} a_{jj} = 0 \quad i \neq j$$

Y_1 จะมีค่าวาเรียนซ์สูงที่สุด
 Y_2 จะมีค่าวาเรียนซ์รองลงมา
 \vdots
 \vdots
 \vdots
 Y_k จะมีค่าวาเรียนซ์ต่ำที่สุด

ดังนั้นสิ่งไข้องค์ประกอบที่ 1 (Principal Component ที่ 1)

$$\text{รูปแบบที่ใช้ } Y_1 = a_{11}X_1 + a_{12}X_2 + \dots + a_{1k}X_k$$

$$\text{หรือ } Y = a_1X_1 + a_2X_2 + \dots + a_kX_k$$

3.4 วิธีวิเคราะห์หาลหสัมพันธ์คาโนนิกอล

ตามปกติเมื่อมีตัวแปรลุ่มอยู่ 2 ตัวคือ X และ Y จะสามารถหาความสัมพันธ์ของตัวแปรทั้งสองได้ โดยพิจารณาจากสัมประสิทธิ์ลหสัมพันธ์ (Correlation Coefficient) หรือ ρ เมื่อ

$$\rho = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$$

ถ้าค่า Y ขึ้นอยู่กับค่า X หลาย ๆ ตัว วิธีการที่จะสามารถหาความสัมพันธ์ระหว่าง Y และ X เหล่านั้นจะใช้วิธีวิเคราะห์หาลหสัมพันธ์พหุคูณ (Multiple Correlation)

โดยให้ X เป็นปัจจัยหรือตัวแปรที่ต้องการศึกษาซึ่งมีอยู่ P ตัว

$$\text{ดังนั้น } X' = [X_1 \ X_2 \ \dots \ X_p]$$

และ \hat{Y} เป็นค่าประมาณของผลรวมเชิงเส้น

$$\hat{Y}_j = \hat{\beta}_1 X_{1j} + \hat{\beta}_2 X_{2j} + \dots + \hat{\beta}_p X_{pj}, \quad j = 1, 2, \dots, n$$

เมื่อ $\hat{\beta}_S = (S = 1, 2, \dots, P)$ เป็นค่าสัมประสิทธิ์การถดถอยเชิงซ้อน (Multiple Regression Coefficient) ซึ่งหาค่าได้จากกฎประมาณโดยวิธี Least-Square อันจะทำให้สามารถหาค่าสัมประสิทธิ์สหสัมพันธ์โดยการหาค่าความสัมพันธ์อย่างง่าย (Simple Correlation) ระหว่าง Y และ X นั้นเอง ยังมีอีกกรณีหนึ่ง ซึ่งนอกจากจะมีตัวแปร X หลาย ๆ ตัวแล้วยังมี Y หลาย ๆ ตัวคือ

$$Y' = [Y_1 \ Y_2 \ \dots \ Y_q]$$

วิธีการที่จะหาความสัมพันธ์ระหว่าง X หลาย ๆ ตัว และ Y หลาย ๆ ตัวก็คือ วิธีวิเคราะห์หลักสัมพันธ์คาโนนิกอล จะได้ค่าความสัมพันธ์คาโนนิกอล (Canonical Correlation) ซึ่งจะได้จากความสัมพันธ์ระหว่างผลรวมเชิงเส้นของกลุ่มตัวแปร X และผลรวมเชิงเส้นของกลุ่มตัวแปร Y

เมื่อกำหนดค่า

$$X_j^* = \alpha_1 X_1 + \alpha_2 X_2 + \dots + \alpha_p X_p \quad \dots \dots \dots (1)$$

$$Y_j^* = \beta_1 Y_1 + \beta_2 Y_2 + \dots + \beta_q Y_q$$

สมการนี้ยังไม่ทราบค่าของ α_i และ β_i ที่จะทำให้ความสัมพันธ์ระหว่าง Y^* และ X^* มีค่ามากที่สุดที่จะเป็นไปได้

ถ้าให้ ρ_c Canonical Correlation Coefficient ระหว่าง X^* และ Y^*

$$\text{ซึ่ง } \rho_c = \frac{\text{Cov}(X^*, Y^*)}{\sqrt{\text{Var}(X^*) \text{Var}(Y^*)}} \quad \dots \dots \dots (2)$$

และ $\rho_c^{(1)}$ เป็น Canonical Correlation Coefficient ที่มีค่ามากที่สุด เรียกว่า First Canonical Correlation ซึ่งเกิดจาก X_1^* และ Y_1^*

$\rho_c^{(2)}$ เป็น Canonical Correlation Coefficient ที่มีค่ามากที่สุดเป็นอันดับสอง เรียกว่า Second Canonical Correlation ซึ่งเกิดจาก X_2^* และ Y_2^*

$\rho_c^{(3)}$ เป็น Canonical Correlation Coefficient ที่มีค่ามากที่สุดเป็นอันดับ
 สามเรียกว่า Third Canonical Correlation ซึ่งเกิดจาก X_3^* และ Y_3^*

· · · ·
 · · · ·
 · · · ·

X_1^* และ Y_1^* เรียกว่า first canonical variables

X_2^* และ Y_2^* เรียกว่า second canonical variables

X_3^* และ Y_3^* เรียกว่า third canonical variables

· · · ·
 · · · ·
 · · · ·

คู่ของ canonical variable ต่าง ๆ นั้น จะมีหลักการเหมือนกัน

กรณี $q < p$ ก็จะมี Canonical Correlation ทั้งหมด q ค่า

และ Canonical Variable ทั้งหมด q คู่

ผลลัพธ์ต่าง ๆ เหล่านี้ได้มาจากหลักการต่อไปนี้คือ

ถ้า X มีค่าเฉลี่ยเป็น μ_1 และความแปรปรวนเท่ากับ Σ_{11}

Y มีค่าเฉลี่ยเป็น μ_2 และความแปรปรวนเท่ากับ Σ_{22}

ให้ $\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$ (p+q) x (p+q)

โดยที่ Σ_{11} เป็นเมตริกของความแปรปรวน (Variable matrix) ของ X
 ซึ่งมีขนาด $p \times p$

Σ_{22} เป็นเมตริกของความแปรปรวน (Variable matrix) ของ Y
 ซึ่งมีขนาด $q \times q$

Σ_{12} เป็นเมตริกของความแปรปรวนร่วม (Covariance matrix) ของ X
และ Y ซึ่งมีขนาด $p \times q$

$$\Sigma_{21} = \Sigma_{12}'$$

$$\Sigma = \begin{array}{c} \left[\begin{array}{ccc|ccc} \sigma_{X_1 X_1} & \sigma_{X_1 X_2} & \dots & \sigma_{X_1 X_p} & \sigma_{X_1 Y_1} & \sigma_{X_1 Y_2} & \dots & \sigma_{X_1 Y_q} \\ \sigma_{X_2 X_1} & \sigma_{X_2 X_2} & \dots & \sigma_{X_2 X_p} & \sigma_{X_2 Y_1} & \sigma_{X_2 Y_2} & \dots & \sigma_{X_2 Y_q} \\ \vdots & \vdots & & \vdots & \vdots & \vdots & & \vdots \\ \sigma_{X_p X_1} & \sigma_{X_p X_2} & & \sigma_{X_p X_p} & \sigma_{X_p Y_1} & \sigma_{X_p Y_2} & & \sigma_{X_p Y_q} \end{array} \right. \\ \hline \left. \begin{array}{ccc|ccc} \sigma_{X_1 Y_1} & \sigma_{X_2 Y_1} & \dots & \sigma_{X_p Y_1} & \sigma_{Y_1 Y_1} & \sigma_{Y_1 Y_2} & \dots & \sigma_{Y_1 Y_q} \\ \vdots & \vdots & & \vdots & \vdots & \vdots & & \vdots \\ \sigma_{X_1 Y_q} & \sigma_{X_2 Y_q} & \dots & \sigma_{X_p Y_q} & \sigma_{Y_q Y_1} & \sigma_{Y_q Y_2} & \dots & \sigma_{Y_q Y_q} \end{array} \right] \end{array}$$

Canonical variables ในรูปของ matrix คือ

$$X_{\alpha}^* = X_{\alpha}$$

$$Y_{\beta}^* = Y_{\beta}$$

เมื่อ X_{α}^* และ Y_{β}^* เป็นเวกเตอร์ของ Canonical variable

ขนาด $p \times 1$ และ $q \times 1$ ตามลำดับ

X_{α} และ Y_{β} เป็นเมตริกของตัวแปรสุ่มแรกเริ่ม (Original random variables)

ขนาด $n \times p$ และ $n \times q$ ตามลำดับ

α เป็น Coefficient vectors ขนาด $p \times 1$

β เป็น Coefficient vectors ขนาด $q \times 1$

ดังนั้นจากสมการ (2) สามารถเขียน Canonical Correlation Coefficient

ได้เป็น

$$\rho_c = \frac{\text{Cov}(X_\alpha, Y_\beta)}{\sqrt{\text{Var}(X_\alpha) \text{Var}(Y_\beta)}}$$



หรือ

$$\rho_c = \frac{\alpha' \Sigma_{12} \beta}{\sqrt{\alpha' \Sigma_{11} \alpha \quad \beta' \Sigma_{22} \beta}} \dots\dots\dots (2a)$$

หาค่า ρ_c จากสมการ (2a) โดยที่ต้องการได้ ρ_c ที่มีค่าสูงสุด ดังนั้นจะ Maximize ρ_c เทียบกับ α และ β ซึ่งวิธีดังกล่าวนี้จะง่ายขึ้น เมื่อ Canonical Variable แต่ละตัวมีคุณสมบัติที่ว่าความแปรปรวนเป็น 1 นั่นคือ

$$\text{Var}(X_j^*) = \alpha' \Sigma_{11} \alpha = 1 \dots\dots\dots (3)$$

$$\text{Var}(Y_j^*) = \beta' \Sigma_{22} \beta = 1 \dots\dots\dots (4)$$

ซึ่งเงื่อนไขนี้จะทำให้ได้ $\rho_c = \alpha' \Sigma_{12} \beta$ เมื่อ Maximize ρ_c โดยวิธี Least Square ซึ่งจะได้ดังนี้

$$\Sigma_{12} \beta - \lambda \Sigma_{11} \alpha = 0 \dots\dots\dots (5)$$

$$\Sigma_{21} \alpha - \lambda \Sigma_{22} \beta = 0 \dots\dots\dots (6)$$

เมื่อ λ เป็น Lagrange multiplier

คูณสมการ (5) ด้วย λ ได้

$$\Sigma_{12} \lambda \beta = \lambda^2 \Sigma_{11} \alpha \dots\dots\dots (7)$$

คูณสมการ (6) ด้วย Σ_{22}^{-1}

$$\Sigma_{22}^{-1} \Sigma_{21} \alpha = \lambda \beta \dots\dots\dots (8)$$

นำสมการ (8) แทนในสมการ (7)

$$\Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \alpha = \lambda^2 \Sigma_{11} \alpha$$

คูณทั้งสองข้างด้วย Σ_{11}^{-1} และจัดเทอมใหม่ จะได้

$$(\Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} - \lambda^2 I) \alpha = 0 \dots\dots\dots(9)$$

กำหนดเดียวกันจะได้

$$(\Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} - \lambda^2 I) \beta = 0 \dots\dots\dots(10)$$

ทั้งสมการ (9) และ (10) เป็น homogeneous equations ซึ่งสามารถหาคำตอบได้ แต่จากการที่ได้กำหนดไว้ก่อนแล้วว่า $q \leq p$ ดังนั้น ถ้าหากทราบค่า λ^2 และ β แล้วก็จะสามารถหา α ได้จากสมการ (5)

เนื่องจาก

$$\alpha = \frac{\Sigma_{11}^{-1} \Sigma_{12} \beta}{\lambda} \dots\dots\dots(11)$$

จากสมการ (5) เมื่อคูณด้วย α' จะได้

$$\alpha' \Sigma_{12} \beta = \lambda \alpha' \Sigma_{11} \alpha$$

จากสมการ (6) เมื่อคูณด้วย β' จะได้

$$\beta' \Sigma_{21} \alpha = \lambda \beta' \Sigma_{22} \beta$$

เนื่องจาก $\alpha' \Sigma_{11} \alpha = \beta' \Sigma_{22} \beta = 1$

ดังนั้น $\alpha' \Sigma_{12} \beta = \beta' \Sigma_{21} \alpha = \lambda \dots\dots\dots(12)$

นั่นคือ จากสมการ (2a) จะได้ว่า $\alpha' \Sigma_{12} \beta = \lambda$ เป็น Canonical Correlation Coefficient ภายใต้เงื่อนไขของสมการ (3) และสมการ (4)

$$\rho_c = \lambda = \alpha' \Sigma_{12} \beta$$

ถ้าหากถอดรากที่สองของ Characteristic root coefficient ที่มีค่ามากที่สุด
ของสมการ (9) หรือ (10) จะได้ first canonical correlation

จากสมการ (10) หาค่า λ^2 และ β ได้ และจาก (11) หาค่า α ได้ ให้

$\lambda_1^2 > \lambda_2^2 > \dots > \lambda_q^2$ เป็น Characteristic root ของสมการ (10)

$\beta_1, \beta_2, \dots, \beta_q$ เป็น Characteristic vectors ที่สัมพันธ์กับ

Characteristic root λ_i^2 ดังกล่าว

$\alpha_1, \alpha_2, \dots, \alpha_q$ เป็น characteristic vectors ซึ่งคำนวณได้

จากการแทนค่า λ_i และ β_i ลงใน (11)

ดังนั้น Canonical Correlation Coefficient ที่ i คือ λ_i

และเซตของ Canonical variables ก็คือ $X_i^* = X \alpha_i$

$$Y_i^* = Y \beta_i$$

ค่าสัมประสิทธิ์สหสัมพันธ์คาโนคอลล (Canonical Correlation Coefficient) ที่ได้มีคุณสมบัติ

เช่นเดียวกับค่าสัมประสิทธิ์สหสัมพันธ์อย่างง่าย (Simple Correlation Coefficient)

กล่าวคือค่าสัมประสิทธิ์สหสัมพันธ์คาโนคอลลจะมีค่าอยู่ระหว่าง 0 และ 1 และจะไม่เปลี่ยนค่า

แม้ว่าจะเปลี่ยนหน่วยการวัดก็ตาม

การคำนวณ เนื่องจากคุณสมบัติดังกล่าวข้างต้นของ Canonical Correlation
Coefficient ดังนั้นเมื่อไม่ทราบค่า Σ เราจะประมาณ Σ ด้วย Covariance Matrix (S)

$$S = \begin{pmatrix} s_{11} & s_{12} \\ s_{21} & s_{22} \end{pmatrix}$$

เมื่อ S เป็น Covariance Matrix ของตัวอย่างหรืออาจจะคำนวณได้จาก
Correlation matrix (R)

$$R = \begin{pmatrix} R_{11} & R_{12} \\ R_{21} & R_{22} \end{pmatrix}$$

เมื่อ R เป็น Correlation Matrix ของตัวอย่าง

R_{11} เป็น matrix ขนาด $p \times p$ ซึ่งแสดงความสัมพันธ์ระหว่างตัวแปร X
ซึ่งมี p ตัว

R_{22} เป็น matrix ขนาด $q \times q$ ซึ่งแสดงความสัมพันธ์ระหว่างตัวแปร Y
ซึ่งมี q ตัว

R_{12} เป็น matrix ขนาด $p \times q$ ซึ่งแสดงความสัมพันธ์ระหว่างตัวแปร X
และตัวแปร Y

$$R_{21} = R'_{12}$$

ค่าของ Canonical Correlation Coefficients ที่คำนวณจาก S หรือ R
นั้น จะได้ผลเช่นเดียวกัน แต่ค่า α_i และ β_i นั้นจะได้ไม่เหมือนกัน กล่าวคือ

ถ้าหากคำนวณ α_i และ β_i จาก S matrix แล้ว ค่าที่ได้นั้นจะเป็นค่าตัวแปรแรก
เริ่ม (Original Variables) ถ้าหากคำนวณจาก R matrix ค่าที่ได้จะเปลี่ยนอยู่ในรูปตัว
แปรแรกเริ่มมาตรฐาน (Standardized Original Variables) อย่างไรก็ตามก็
สามารถที่จะเปลี่ยนค่าสัมประสิทธิ์ที่คำนวณได้จากการใช้ S นั้นไปเป็นค่าสัมประสิทธิ์ที่ได้จาก R
ดังจะเห็นได้จาก

$$\begin{bmatrix} \frac{1}{\sqrt{S_{11}}} & 0 \\ 0 & \frac{1}{\sqrt{S_{22}}} \end{bmatrix} \begin{bmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{S_{11}}} & 0 \\ 0 & \frac{1}{\sqrt{S_{22}}} \end{bmatrix} = \begin{bmatrix} 1 & R_{12} \\ R_{21} & 1 \end{bmatrix}$$

ดังนั้น เมื่อนำ diagonal matrix ของส่วนกลับของค่าส่วนเบี่ยงเบนมาตรฐาน มาคูณทั้งหน้าและหลังของ Covariance Matrix แล้ว ก็จะสามารถหาค่าของ Correlation Matrix α , β ได้ กล่าวคือ

$$\text{ถ้า } \alpha_{i1}^* , \alpha_{i2}^* , \dots , \alpha_{ip}^* \text{ และ } \beta_{i1}^* , \beta_{i2}^* , \dots , \beta_{iq}^* ;$$

$$(i = 1, 2, \dots, q)$$

เป็นค่าที่คำนวณได้จาก matrix S แล้วจะสามารถหาค่า α β ซึ่งคำนวณได้จาก R ได้ โดยการหารแต่ละค่าสัมประสิทธิ์เหล่านั้นด้วยค่าส่วนเบี่ยงเบนมาตรฐานของตัวแปรแรกเริ่ม (Original Variable) นั่นคือ

$$\frac{\alpha_{i1}^*}{S_{x_{i1}}} , \frac{\alpha_{i2}^*}{S_{x_{i2}}} , \frac{\alpha_{i3}^*}{S_{x_{i3}}} , \dots , \frac{\alpha_{ip}^*}{S_{x_{ip}}}$$

และ

$$\frac{\beta_{i1}^*}{S_{y_{11}}} , \frac{\beta_{i2}^*}{S_{y_{22}}} , \frac{\beta_{i3}^*}{S_{y_{33}}} , \dots , \frac{\beta_{iq}^*}{S_{y_{qq}}}$$

การทดสอบนัยสำคัญของการวิเคราะห์หลักสัมพันธภาพแคนอนิคอล เมื่อทำการวิเคราะห์หลักสัมพันธภาพแคนอนิคอลแล้ว สิ่งที่จะได้คือ

$$(1) \rho_c$$

$$(2) \alpha \quad \text{และ} \quad \beta$$

ค่าของ α และ β ซึ่งสอดคล้องกับ λ^2 ของเมตริก $\Sigma_{11}^{-1} \quad \Sigma_{12} \quad \Sigma_{22}^{-1} \quad \Sigma_{21}$ โดยที่ λ^2 ก็คือความแปรปรวนที่ร่วมกันของตัวแปรทั้งสองชุดนั่นเอง ด้วยเหตุนี้ก็จะมี α และ β อยู่หลายชุดด้วยกัน ดังนั้นสิ่งที่จะต้องคำนึงถึงคือ Canonical Variables ทั้ง q คู่ นั้นมีความสัมพันธ์กันหรือไม่ คือต้องมีการทดสอบสมมติฐานที่ว่า ตัวแปรทั้งสองชุดนั้นมีความสัมพันธ์กันหรือไม่ นั่นคือทดสอบว่า Σ_{12} หรือ $\Sigma_{21} = 0$ หรือไม่

$$H_0 : \Sigma_{12} = 0 \qquad H_0 : \Sigma_{21} = 0$$

$$H_a : \Sigma_{12} \neq 0 \qquad H_a : \Sigma_{21} \neq 0$$

ใช้วิธีการทดสอบของ Bartlett ดังนี้

วิธีทดสอบของ Bartlett ใช้ตัวสถิติ V_0 โดยที่

$$V_0 = - \left[n-1 - \frac{1}{2} (p+q+1) \right] \ln \Lambda_0$$

เมื่อ

$$\Lambda_0 = \frac{1}{\prod_{i=1}^q \left(1 + \frac{\lambda_i^2}{(1-\lambda_i^2)} \right)}$$

$$= \frac{1}{\prod_{i=1}^q \left(\frac{1}{(1-\lambda_i^2)} \right)}$$

$$= \prod_{i=1}^q (1-\lambda_i^2)$$

V_0 จะมีการกระจายแบบไคลส์แควร์ที่องศาความเป็นอิสระเท่ากับ pq ที่ระดับนัยสำคัญที่ต้องการ ถ้าค่า V_0 ที่คำนวณได้มากกว่าค่าไคลส์แควร์จากตารางจะสรุปได้ว่า $\rho_c^{(1)}$ ซึ่งเป็น first canonical correlation coefficient นั้นไม่เท่ากับ 0 และจะทดสอบความสัมพันธ์ของ ρ_c ของคู่ที่เหลือต่อไปคือ

$$\Lambda_1 = \prod_{i=2}^q (1 - \lambda_i^2)$$

ซึ่งจะนำไปสู่ Bartlett V_1 ดังนี้

$$V_1 = - \left[n-2 - \frac{1}{2} (p+q+1) \right] \ln \Lambda_1$$

Λ_1 มีการกระจายแบบไคลส์แควร์ที่องศาความเป็นอิสระเท่ากับ $(p-1)(q-1)$ การสรุปผลจะทำเหมือน V_0 ดังนั้น ถ้าหากมี k canonical correlation coefficients ที่ทดสอบแล้วปรากฏว่ามีนัยสำคัญส่วนที่เหลือก็就会被ทดสอบว่า $(q-k)$ Canonical Correlation Coefficient เหล่านั้นจะเท่ากับ 0 หรือไม่กล่าวคือ

$$\Lambda_k = \prod_{i=(k+1)}^q (1 - \lambda_i^2)$$

$$V_k = - \left[n-1-k - \frac{1}{2} (p+q+1) \right] \ln \Lambda_k$$

V_k มีการกระจายแบบไคลส์แควร์ที่องศาความเป็นอิสระเท่ากับ $(p-k)(q-k)$

ในการประมวลผลด้วยโปรแกรมสำเร็จรูป SPSS นั้น ผลลัพธ์ที่ได้จะมีค่าหรือค่าของความน่าจะเป็นที่จะไม่ยอมรับ H_0 เมื่อ H_0 จริง ในกรณีที่จะทำทดสอบสมมติฐาน ณ ระดับนัยสำคัญ $\alpha = 0.05$ นั้น อาจใช้ค่าของความน่าจะเป็นที่จะไม่ยอมรับ H_0 เมื่อ H_0 จริง มาเปรียบเทียบกับ $\alpha = 0.05$ เพื่อใช้เป็นเกณฑ์ในการตัดสินใจว่าจะยอมรับสมมติฐาน H_0 หรือไม่โดยที่

ถ้าค่าความน่าจะเป็นที่จะไม่ยอมรับ H_0 หรือ H_0 จริง < 0.05 จะไม่ยอมรับ H_0

ถ้าค่าความน่าจะเป็นที่จะไม่ยอมรับ H_0 เมื่อ H_0 จริง > 0.005 จะยอมรับ H_0

3.5 วิธีวิเคราะห์จำแนกประเภท

เป็นวิธีวิเคราะห์ทางสถิติโดยมีจุดมุ่งหมายที่จะคัดเลือกตัวแปรชุดหนึ่ง ซึ่งนักวิจัยคิดว่าตัวแปรชุดนี้มีความสัมพันธ์กับสิ่งที่ต้องการศึกษา จนถึงขั้นที่ตัวแปรชุดนี้เป็นตัวแบ่งแยกประชากรออกเป็นกลุ่มต่าง ๆ ได้อย่างชัดเจน

3.5.1 ขั้นตอนในการวิเคราะห์จำแนกประเภท มี 2 ขั้นตอนคือ

ขั้นที่ 1 การคัดเลือกตัวแปรชุดหนึ่ง เพื่อสร้างสมการที่ใช้ในการจำแนกประชากรออกเป็นกลุ่มต่าง ๆ กัน ได้อย่างชัดเจน สมการนี้คือ สมการจำแนกประเภท (Discriminant equation)

ขั้นที่ 2 การจำแนกตัวอย่างที่ได้ศึกษามานั้นเข้าเป็นสมาชิกของประชากรแต่ละกลุ่มโดยอาศัยสมการจำแนกประเภท

ในกรณีที่มีประชากร 2 กลุ่ม จะใช้วิธีการของ Fisher มาใช้ในการจำแนกประเภท กล่าวคือ ถ้ากำหนดให้

Π_1 เป็นประชากรกลุ่มที่ 1

Π_2 เป็นประชากรกลุ่มที่ 2

ตัวแปรที่ศึกษา คือ ตัวแปร X ซึ่งมีทั้งหมด p ตัวคือ

$$X = \begin{bmatrix} X_1 & X_2 & \dots & X_p \end{bmatrix}$$

วิธีการของ Fisher นั้น จะแปลงตัวแปร X เหล่านี้ไปเป็นค่าของตัวแปรเพียงตัวเดียวคือ Y โดยที่ Y_1 และ Y_2 เป็นค่าสังเกตที่ได้จากประชากร Π_1 และ Π_2 ตามลำดับ

ถ้า μ_{1y} = ค่าเฉลี่ยของค่า Y ซึ่งได้มาจากค่า X ของประชากร Π_1

μ_{2y} = ค่าเฉลี่ยของค่า Y ซึ่งได้มาจากค่า X ของประชากร Π_2

μ_1 = $E(X/\Pi_1)$ = ค่าคาดหวังของตัวแปร X ที่มาจากประชากร Π_1

μ_2 = $E(X/\Pi_2)$ = ค่าคาดหวังของประชากร X ที่มาจากประชากร Π_2 } ..(1)

และโควาเรียนซ์เมตริก

$$\Sigma = E(X_i - \mu_i) (X_i - \mu_i)'; i = 1, 2, \dots \dots \dots (2)$$

พิจารณาผลรวมเชิงเส้น

$$Y = \beta' X \dots \dots \dots (3)$$

(1x1) (1xp) (px1)

เมื่อ Y คือ Fisher's Linear Discriminant Function (ผลรวมเชิงเส้น)

β คือ ค่าที่แสดงถึงความสำคัญของตัวแปร X_1, X_2, \dots, X_p

$$\beta' = [\beta_1 \quad \beta_2 \quad \dots, \quad \beta_p]$$

X คือตัวแปรที่ต้องการศึกษา ซึ่งมีทั้งหมด p ตัว

และ $X' = [X_1 \quad X_2 \quad \dots, \quad X_p]$

ดังนั้น จะได้
$$\left. \begin{aligned} \mu_{1y} &= E(Y/\Pi_1) = E(\beta' X/\Pi_1) = \beta' \mu_1 \\ \mu_{2y} &= E(Y/\Pi_2) = E(\beta' X/\Pi_2) = \beta' \mu_2 \end{aligned} \right\} \dots \dots \dots (4)$$

และค่าความแปรปรวนของ Y จากทั้ง 2 ประชากรนั้นคือ

$$\begin{aligned} \sigma_y^2 &= \text{Var}(\beta' X) = \beta' \text{Cov}(X) \beta \\ &= \beta' \Sigma \beta \dots \dots \dots (5) \end{aligned}$$

วิธีการของ Fisher คือพยายามหาผลรวมเชิงเส้นของค่า X ซึ่งจะทำให้ระยะทางระหว่าง μ_{1y} และ μ_{2y} มีค่ามากที่สุดที่จะเป็นไปได้

ค่าผลรวมเชิงเส้นซึ่งดีที่สุดจะต้องมีคุณสมบัติในการแบ่งแยกประชากรทั้ง 2 กลุ่มออกจากกัน ได้มากที่สุด ซึ่งวิธีการที่จะได้ผลรวมเชิงเส้นที่ดีนั้นก็โดยการหาค่า β ที่ทำให้อัตราส่วน

(ระยะทางระหว่างค่าเฉลี่ยของ y)² มีค่ามากที่สุด

ความแปรปรวนของ Y

$$\begin{aligned}
 \text{ให้ } \lambda &= \frac{(\text{ระยะทางระหว่างค่าเฉลี่ยของ } Y)^2}{\text{ความแปรปรวนของ } Y} \\
 &= \frac{(\mu_{1Y} - \mu_{2Y})^2}{\sigma_Y^2} \\
 &= \frac{(\beta' \mu_1 - \beta' \mu_2)^2}{\beta' \Sigma \beta} \\
 &= \frac{\beta' (\mu_1 - \mu_2) (\mu_1 - \mu_2)' \beta}{\beta' \Sigma \beta} \dots\dots\dots (6)
 \end{aligned}$$

ค่า β ที่ได้จะมีค่ามากที่สุดเมื่อ

$$\frac{\partial \lambda}{\partial \beta} = 0$$

นั่นคือ $C (\mu_1 - \mu_2)' - \Sigma \beta = 0$

$$\beta = C \Sigma^{-1} (\mu_1 - \mu_2)'$$

กำหนดให้ $C = 1$

จะได้ผลรวมเชิงเส้นเป็น

$$Y = \beta' X = (\mu_1 - \mu_2)' \Sigma^{-1} X \dots\dots\dots (7)$$

สมการ $Y = (\mu_1 - \mu_2)' \Sigma^{-1} X$ สามารถใช้เป็นตัวจำแนกค่า

สังเกตที่ได้มานั้นว่าจะอยู่ในประชากรกลุ่ม Π_1 หรือ Π_2 ได้โดยให้

$$Y_0 = (\mu_1 - \mu_2)' \Sigma^{-1} X_0$$

เป็นค่าของ discriminant function ของค่าสังเกต X_0 และให้

$$\begin{aligned}
 m &= \frac{1}{2} (\mu_{1y} + \mu_{2y}) \\
 &= \frac{1}{2} (\mu'_1 + \mu'_2) \\
 &= \frac{1}{2} (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 + \mu_2) \dots\dots\dots (8)
 \end{aligned}$$



เป็นจุดกึ่งกลาง (centroid) ระหว่างค่าเฉลี่ยของค่า Y จากทั้ง 2 ประชากร

สามารถเขียนกฎการจำแนกประเภทได้ดังนี้คือ

$$\left. \begin{aligned}
 \text{จัด } X_0 \text{ ให้อยู่ } \Pi_1 \text{ ถ้าหากว่า } Y_0 &= (\mu_1 - \mu_2)' \Sigma^{-1} X_0 < m \\
 \text{จัด } X_0 \text{ ให้อยู่ } \Pi_2 \text{ ถ้าหากว่า } Y_0 &= (\mu_1 - \mu_2)' \Sigma^{-1} X_0 > m
 \end{aligned} \right\} \dots\dots\dots (9)$$

ในทางปฏิบัติจะไม่ทราบค่า μ_1, μ_2 และ Σ ดังนั้นถ้ากลุ่มตัวอย่างจาก Π_1 และ Π_2 มาเป็นจำนวน n_1 และ n_2 แล้วโดยที่วัดค่าสังเกต

$$\begin{aligned}
 \bar{X}' &= [x_1 \quad x_2 \quad \dots \quad x_p] \quad \text{จะได้ว่า} \\
 \bar{X}_1 &= [x_{11} \quad x_{12} \quad \dots \quad x_{1n_1}] \\
 (p \times n_1) & \\
 \bar{X}_2 &= [x_{21} \quad x_{22} \quad \dots \quad x_{2n_2}] \\
 (p \times n_2) &
 \end{aligned}$$

จะได้ค่าประมาณของ μ_1, μ_2 และ Σ^{-1} เป็น \bar{X}_1, \bar{X}_2 และ

S^{-1} ตามลำดับโดยที่
pooled

$$\bar{x}_1 = \frac{1}{n_1} \sum_{j=1}^{n_1} x_{1j} \dots\dots(10) ; \bar{x}_2 = \frac{1}{n_2} \sum_{j=1}^{n_2} x_{2j} \dots\dots(11)$$

$$S = \left[\frac{n_1 - 1}{(n_1-1)+(n_2-1)} \right] S_1 + \left[\frac{n_2 - 1}{(n_1-1)+(n_2-1)} \right] S_2$$

$$= \frac{(n_1 - 1) S_1 + (n_2-1) S_2}{n_1+n_2-2} \dots\dots\dots(12)$$

$$S_1 \text{ (pxp)} = \frac{1}{n_1-1} \sum_{j=1}^{n_1} (x_{1j} - \bar{x}_1) (x_{1j} - \bar{x}_1)'$$

$$S_2 \text{ pxp} = \frac{1}{n_2-1} \sum_{j=1}^{n_2} (x_{2j} - \bar{x}_2) (x_{2j} - \bar{x}_2)'$$

ดังนั้น เมื่อแทนค่า \bar{x}_1 , \bar{x}_2 และ S_{pooled}^{-1} ลงใน (7) จะได้

Fisher's Sample Linear Discriminant Function เป็น

$$y = \hat{\beta}' x$$

$$= (\bar{x}_1 - \bar{x}_2)' S_{pooled}^{-1} x \dots\dots\dots(13)$$

และค่าจุดกึ่งกลางระหว่างตัวแปร $\bar{x}_1 = \hat{\beta}' \bar{x}_1$ $\bar{x}_2 = \hat{\beta}' \bar{x}_2$

$$\hat{m} = \frac{1}{2} (\bar{x}_1 + \bar{x}_2)$$

$$= \frac{1}{2} (\bar{x}_1 - \bar{x}_2)' S_{pooled}^{-1} (\bar{x}_1 + \bar{x}_2) \dots\dots\dots(14)$$

กฎการจำแนกประเภทจะเขียนดังนี้

$$\text{จัดให้ } x_0 \text{ อยู่ใน } \Pi_1 \text{ ถ้า } y_0 < m$$

$$\text{จัดให้ } x_0 \text{ อยู่ใน } \Pi_2 \text{ ถ้า } y_0 > m$$

$$\text{เมื่อ } y_0 = (\bar{x}_1 - \bar{x}_2)' s_{\text{pooled}}^{-1} x_0$$

จากสมการสามารถเขียนได้เป็น

$$y = \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

ค่า $\beta_1, \beta_2, \dots, \beta_p$ จะเป็นค่าแสดงถึงความสำคัญของตัวแปร x_1, x_2, \dots, x_p ซึ่งโดยปกติ x_i แต่ละค่าจะมีหน่วยไม่เหมือนกัน การเปรียบเทียบค่า β_i เหล่านี้ทำได้ โดยการปรับค่า β_i เหล่านี้ให้เป็นมาตรฐาน (Standardized) เสียก่อน

การแปลงค่า β_i เหล่านี้สามารถทำได้โดยเอาสมาชิกในแนวเส้นทแยงมุม (diagonal) ของเมตริก W มาถอดรากที่สอง แล้วนำไปคูณกับค่า β_i ทุก ๆ ค่าตามสูตรต่อไปนี้

$$\beta_i^* = (\sqrt{w_{ii}}) (\beta_i)$$

เมื่อ W คือ $(x - \bar{x}_{g0}) (x - \bar{x}_{g0})'$

$$x = \begin{pmatrix} x_{111} & x_{121} & \dots & x_{1n1} & | & x_{112} & x_{122} & \dots & x_{1n2} \\ x_{211} & x_{221} & \dots & x_{2n1} & | & x_{212} & x_{222} & \dots & x_{2n2} \\ \vdots & \vdots & & \vdots & | & \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots & | & \vdots & \vdots & & \vdots \\ x_{p11} & x_{p21} & \dots & x_{pn1} & | & x_{p12} & x_{p22} & \dots & x_{pn2} \end{pmatrix}$$

กลุ่มที่ 1

กลุ่มที่ 2

$$\bar{X}_{go} = \begin{bmatrix} \bar{X}_{11} & \cdots & \bar{X}_{11} & \cdots & \bar{X}_{12} & \cdots & \bar{X}_{12} \\ \bar{X}_{21} & \cdots & \bar{X}_{21} & \cdots & \bar{X}_{22} & \cdots & \bar{X}_{22} \\ \vdots & & \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots & & \vdots \\ \bar{X}_{p1} & \cdots & \bar{X}_{p1} & \cdots & \bar{X}_{p2} & \cdots & \bar{X}_{p2} \end{bmatrix}$$

β_i^* คือ Standardized discriminant weight ของ discriminant function

β_i คือ Raw discriminant weight ของ discriminant function

W_{ii} คือ Diagonal element ของ W เมตริก โดยที่ $i = 1, 2, \dots, p$

ล้มการจำแนกประเภทที่ได้มานั้นจะมีความล้มการในการแบ่งแยกกลุ่มได้ดีหรือไม่นั้นสามารถดูได้จากค่าความผิดพลาด เนื่องจากการจำแนกผิด (miss classification error) ซึ่งทราบได้จากการที่เปรียบเทียบว่า สิ่งที่ศึกษาอยู่นั้นเป็นสมาชิกของประชากรกลุ่มหนึ่งซึ่งทำให้ได้ผลที่ผิดจากความเป็นจริงไป

3.5.2 วิธีการคัดเลือกตัวแปรเข้าไปในล้มการจำแนกประเภท

จากการที่กล่าวมาแล้วข้างต้นนั้นเป็นหลักการโดยทั่วไปของวิธีการสร้างล้มการจำแนกประเภทในกรณีที่ว่าในการวิจัยนั้น มีตัวแปรเป็นจำนวนมาก การคัดเลือกตัวแปรให้เหลือจำนวนน้อยที่สุดแต่มีความล้มการในการใช้เป็นตัวจำแนกมากที่สุดนั้น สามารถทำได้โดยใช้การคัดเลือกตัวแปรทีละตัว โดยที่จะหาตัวแปรที่ดีที่สุดตัวแรกและตัวแปรที่ดีที่สุดที่สองที่จะช่วยให้ล้มการจำแนกประเภทมีความล้มการจำแนกประเภทได้ดีที่สุด จากนั้นก็จะเลือกตัวที่สามและตัวต่อไปที่จะช่วยการจำแนกให้ดีขึ้นตามลำดับ ในแต่ละขั้นตอนตัวแปรที่ได้รับการคัดเลือกมาก่อนแล้วนั้นอาจถูกตัดทิ้งไปหากพบว่าเมื่อนำมารวมกับตัวแปรอื่น ๆ แล้วไม่ช่วยให้ล้มการจำแนกประเภทดีขึ้น วิธีการนี้เรียกว่า วิธีการสร้างล้มการจำแนกประเภทแบบขั้นตอน (Stepwise discriminant analysis) ซึ่งในการวิจัยครั้งนี้ไม่ได้ใช้วิธีนี้ แต่จะกำหนดตัวแปรเข้าไปในล้มการโดยใช้ตัวแปรชุดเดียวกับตัวแปรที่ใช้ในการวิเคราะห์การถดถอยเชิงซ้อน เนื่องจากต้องการจะเปรียบเทียบความเหมาะสมของแต่ละวิธีการวิเคราะห์ด้วยตัวแปรชุดเดียวกัน

จากวิธีการดังกล่าวข้างต้น จะทำให้ได้สมการจำแนกประเภทและค่า
ค่าแครคเตอร์ลิสติก รุท ซึ่งเป็นค่าที่ได้จากขั้นตอนการหาสมการวิเคราะห์จำแนกประเภทโดยเป็นค่า
ความแปรปรวนของค่า Y ซึ่งได้จากการแปลงรูป (Transform) จากค่า X ต่าง ๆ
 $[X_1 \ X_2 \ \dots \ X_p]$ เป็นค่าที่ใช้วัดความสำคัญเชิงเปรียบเทียบของสมการได้ โดยที่สมการ
วิเคราะห์จำแนกประเภทที่ได้มานั้นได้ตามลำดับความสำคัญของค่า ค่าแครคเตอร์ลิสติก รุท ต่าง ๆ
จากมากที่สุดและรองลงไปตามลำดับ ค่าผลรวมของ ค่าแครคเตอร์ลิสติก รุท ทั้งหมดนั้นเป็นค่าความ
แปรปรวนทั้งหมดของตัวแปรจำแนกประเภท (ตัวแปร X) ค่าค่าแครคเตอร์ลิสติก รุท แต่ละค่าจึง
คิดเป็นอัตราส่วนร้อยละของค่ารวมของ ค่าแครคเตอร์ลิสติก รุท ทั้งหมด ทำให้สามารถนำค่านี้ไปอ้างอิง
ถึงความสำคัญเชิงเปรียบเทียบของสมการจำแนกประเภทได้

3.5.3 การทดสอบนัยสำคัญของสมการจำแนกประเภท

เมื่อได้สมการจำแนกประเภทมาแล้วและอยากทราบว่าสมการนี้สามารถมีอำนาจ
ในการจำแนกกลุ่มได้อย่างมีนัยสำคัญทางสถิติหรือไม่ เราสามารถทดสอบได้จาก

$$V_m = [N - 1 - (p+k)/2] \ln(1 + \lambda_m)$$

โดยที่ V_m คือ Bartlett V Statistic ซึ่งมีการกระจายเป็นแบบ
Chi-Square ที่มี d.f. = $p+k-2m$

N คือ จำนวนตัวอย่างทั้งหมดจากประชากรทั้ง 2 กลุ่ม

p คือ จำนวนตัวแปรทั้งหมด

k คือ จำนวนกลุ่ม

λ_m คือ ค่าแครคเตอร์ลิสติก รุท ต่าง ๆ เช่น ค่าตัวที่ 1, 2, 3, ..., m

ค่า λ_m นี้คำนวณได้จาก

$$|A - \lambda_m I| = 0$$

เมื่อ A คือ โควาเรียนซ์เมตริกของ X

ค่า λ_m นี้จะเลือกเอาเฉพาะค่าที่ไม่เป็นศูนย์

ถ้าค่า V_m ที่คำนวณได้มากกว่าค่าไคลเควร์ที่ระดับนัยสำคัญที่กำหนดไว้แล้ว แสดงว่าสมการนี้สามารถใช้จ่ายแยกกลุ่มได้อย่างมีนัยสำคัญทางสถิติ

การคำนวณค่าอำนาจการจำแนกของกลุ่มของตัวแปร (Total Discriminatory Power)

$$\hat{W}^2 = 1 - \frac{N}{(m-k)(1+\lambda_1) + (1+\lambda_2) \dots (1+\lambda_m) + 1}$$

เมื่อ \hat{W}^2 = ค่าอำนาจในการแยกตัวแปรได้จากการประมาณค่า

N = จำนวนตัวอย่างทั้งหมดจากประชากรทั้ง 2 กลุ่ม

k = จำนวนกลุ่ม

λ_m = ค่าค่าแรกเตอร์สถิติ รุท ต่าง ๆ

m = จำนวนค่าค่าแรกเตอร์สถิติ รุท

3.6 การคัดเลือกแบบจำลองทางคณิตศาสตร์ที่เหมาะสมสำหรับการคัดเลือกพื้นที่เพื่อประกาศเขตปฏิรูปที่ดิน

จากแบบจำลองทั้ง 4 ชุดที่ได้จากการวิเคราะห์การถดถอยเชิงซ้อน การวิเคราะห์องค์ประกอบหลัก การวิเคราะห์ค่าโนดอล และการวิเคราะห์จำแนกประเภท ใช้ข้อมูลจาก 37 อำเภอ ใน 5 จังหวัดที่คัดเลือกมาวิจัยครั้งนี้ นำค่าตัวแปรทั้งหมดของแต่ละอำเภอมาทนลงในแบบจำลองที่ได้จากการวิเคราะห์ทั้ง 4 วิธีนี้จะได้คะแนนความสำคัญของแต่ละอำเภอ ทั้ง 37 อำเภอ ต่อจากนั้นทำการจัดลำดับของอำเภอที่ลุ่มตัวอย่างมา 37 อำเภอนี้ตามคะแนนความสำคัญของแต่ละอำเภอ ก็จะได้ลำดับที่ความสำคัญของพื้นที่ (อำเภอ) ที่แสดงถึงระดับการพัฒนาแล้วทั้ง 4 ชุด

จากการศึกษาเรื่องเส้นวัดความยากจนของนักเศรษฐศาสตร์ และการศึกษาเกี่ยวกับความยากจนของประชากรในชนบทของประเทศไทยของธนาคารโลก ปรากฏว่าต่างก็ใช้รายได้ของประชากรมาเป็นตัวพิจารณา แสดงว่ารายได้เป็นตัวบ่งชี้ที่ดีในการวัดสภาพความเป็นอยู่ของประชากร แต่สิ่งที่ไม่สามารถบอกได้คือความสำคัญของปัญหาในแต่ละตัวที่ประชากรในชนบทกำลังเผชิญอยู่ ดังนั้นจึงต้องดำเนินการวิเคราะห์การถดถอยเชิงซ้อน การวิเคราะห์องค์ประกอบหลัก การวิเคราะห์ค่าโนดอล และการวิเคราะห์ค่าแยกประเภท ซึ่งได้กำหนดน้ำหนักความสำคัญของตัวแปรแต่ละตัวขึ้นมา ก็จะสามารรถชี้ให้เห็นถึงความสัมพันธ์ของปัญหาในแต่ละด้าน ซึ่งจะมีประโยชน์ในการคัดเลือกพื้นที่เพื่อประกาศเขตปฏิรูปที่ดินต่อไป อย่างไรก็ตามแบบจำลองที่จะใช้ได้ก็ควรจะเป็นแบบจำลองที่เมื่อนำมาใช้แล้วให้อันดับของอำเภอที่ล่อคคล้องกับลำดับของอำเภอที่ใช้รายได้เฉลี่ยต่อครัวเรือนต่อปีของประชากรในแต่ละอำเภอ

3.6.1 การวัดความล่อคคล้องของตัวแปรประเภทจัดอันดับ

ตัวสถิติที่ใช้สำหรับการหาความสัมพันธ์ระหว่างตัวแปรประเภทจัดอันดับมีหลายตัวด้วยกันเช่น สัมประสิทธิ์สหสัมพันธ์เชิงอันดับของสเปียร์แมน (γ_s) และสัมประสิทธิ์ความสัมพันธ์ของครัสคาลและกูดแมน (G) และสัมประสิทธิ์สหสัมพันธ์เชิงอันดับของเคนดัลลี (τ) ในที่นี้จะใช้สัมประสิทธิ์สหสัมพันธ์เชิงอันดับของสเปียร์แมน (γ_s) เพราะมีขั้นตอนของการคำนวณที่ง่ายกว่า ดังนี้

สัมประสิทธิ์สหสัมพันธ์เชิงอันดับของสเปียร์แมน (Spearman's Coefficient of Rank Correlation ; γ_s) เป็นวิธีการที่ใช้วัดความสัมพันธ์ระหว่างตัวแปรประเภทจัดอันดับ วิธีการก็คือ การแบ่งข้อมูลออกเป็นสองอนุกรม แล้วหาผลต่างระหว่างข้อมูลทั้งสองชุดนั้น แต่ละคู่ของอันดับเดียวกันแล้วยกกำลังสอง จากนั้นนำค่าที่ได้มารวมกัน เพื่อหาค่าสัมประสิทธิ์ที่ต้องการ ในทางปฏิบัติ การหาความสัมพันธ์อาจเป็นความสัมพันธ์ระหว่างตัวแปรประเภทเดียวกันของประชากรกลุ่มเดียวกัน แต่มีการจัดอันดับสองครั้ง หรือเป็นความสัมพันธ์ระหว่างอันดับของบุคคลเดียวกันบนตัวแปรสองตัว เป็นต้น

สัมประสิทธิ์สหสัมพันธ์เชิงอันดับของสเปียร์แมนคำนวณได้จากสูตร*

$$\gamma_s = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2 - 1)}$$



โดยที่ γ_s = ค่าสัมประสิทธิ์สหสัมพันธ์เชิงอันดับของสเปียร์แมน

d_i = ค่าของผลต่างระหว่างอันดับ

n = จำนวนตัวอย่าง

เป็นการจัดอันดับโดยไม่มีตำแหน่งซ้ำกัน

ในกรณีที่เป็นการจัดอันดับโดยมีตำแหน่งซ้ำกัน ก็ใช้วิธีเฉลี่ย เช่น ถ้ามีอันดับที่ 5 ซ้ำกัน 2 ตัว แต่ละตัวจะเป็นอันดับ 4.5 แต่ถ้ามีซ้ำกัน 3 ตัว ทุกตัวจะได้อันดับเดียวกันคือตัวกลาง เช่น ถ้าตัวที่ 6 7 8 มีค่าเท่ากัน ทุกตัวจะได้อันดับ 7 หมา และการใช้สูตรข้างต้นจำเป็นต้องมีการปรับปรุงเพื่อให้ได้ค่าที่ถูกต้องยิ่งขึ้นดังนี้

$$\gamma_s = \frac{\sum X_i^2 - \sum Y_i^2 - \sum d_i^2}{2 \sqrt{\sum X_i^2 \sum Y_i^2}}$$

* โดยทั่วไป ถ้าค่า γ_s เข้าใกล้ +1 หรือ -1 ถือว่ามีความสอดคล้องหรือมีความผกผันกันอย่างมาก แต่ถ้าค่า γ_s เข้าใกล้ 0 แสดงว่าการจัดลำดับไม่มีความสอดคล้องหรือผกผันกันแต่อย่างใด และถ้าค่า γ_s มีค่าต่ำกว่า +0.5 แสดงว่ามีความสอดคล้องกันบ้างหรือค่า γ_s มีค่าน้อยกว่า -0.5 ก็แสดงว่ามีความผกผันกันบ้าง

$$\text{โดยที่ } \Sigma X_i^2 = \frac{n^3 - n}{12} - \Sigma T_x$$

$$\Sigma Y_i^2 = \frac{n^3 - n}{12} = \Sigma T_y$$

$$T = \frac{t^3 - t}{12}$$

X_i คือ ตัวแปรตัวหนึ่ง

Y_i คือ ตัวแปรตัวหนึ่ง

T_x คือ จำนวนครั้งที่ซ้ำกันในตัวแปร X

T_y คือ จำนวนครั้งที่ซ้ำกันในตัวแปร Y

n คือ จำนวนข้อมูล

d_i คือ ผลต่างระหว่างอันดับ

t คือ จำนวนครั้งที่ซ้ำกันในอันดับเดียวกัน

การใช้สูตรปรับปรุงควรจะใช้ในกรณีที่เป็นการจัดอันดับโดยมีตำแหน่งที่ซ้ำกัน

มาก ๆ เนื่องจากจะทำให้ค่าของ γ_s ที่ได้มีความถูกต้องยิ่งขึ้น