

บทที่ 2

ระเบียบวิธีที่ใช้ในการวิจัย

สิ่งที่สนใจศึกษาในการวิจัยครั้งนี้ คือ การเปรียบเทียบการจำแนกกลุ่มโดยวิธีการวิเคราะห์การถดถอยทวิและการวิเคราะห์จำแนกประเภท เมื่อตัวแปรอิสระมีการแจกแจงแบบเบ้ ในการศึกษาเปรียบเทียบจะใช้ค่าเฉลี่ยของร้อยละที่พยากรณ์ถูกต้อง เป็นเกณฑ์เปรียบเทียบ โดยจะกล่าวถึงรายละเอียดสำหรับแต่ละวิธีดังต่อไปนี้

- 2.1 ข้อสมมติเบื้องต้นของตัวแบบการวิเคราะห์การถดถอยเชิงเส้น
- 2.2 วิธีกำลังสองน้อยที่สุดแบบธรรมดา (Ordinary Least Square method)
- 2.3 การวิเคราะห์การถดถอยแบบทวิ (Binary Regression)
- 2.4 การวิเคราะห์จำแนกประเภท (Discriminant Analysis)

พิจารณาจาก 2 กรณี คือ

- 2.4.1 ประชากรทั้ง 2 กลุ่ม มีปัจจัยต่อไปนี้เท่ากัน

- Cost of misclassification

- Prior Probability

- 2.4.2 ประชากรทั้ง 2 กลุ่ม มีปัจจัยต่อไปนี้แตกต่างกัน

- Cost of misclassification ทั้ง 2 กลุ่มเท่ากัน

- Prior Probability ในแต่ละกลุ่มแตกต่างกัน

- 2.1 ข้อสมมติทั่วไปของตัวแบบการถดถอยเชิงเส้น

ในการสมมติความสัมพันธ์เชิงเส้นระหว่างตัวแปร y และ ตัวแปรอิสระ m ตัวของ x_1, x_2, \dots, x_m และค่าผิดพลาด ϵ นั่นคือ ถ้ามีตัวอย่าง n ค่าสังเกตของ y และค่า x เราสามารถเขียนสมการได้เป็น

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_m x_{im} + \epsilon_i ; i = 1, \dots, n \quad (2.1)$$

ค่าสัมประสิทธิ์ β และค่าพารามิเตอร์ของการแจกแจงของ ϵ ไม่ทราบค่า ปัญหาคือเราต้องการประมาณ β ที่ไม่ทราบเหล่านี้เมื่อมี n สมการของ (2.1) ซึ่งเราสามารถเขียนในรูปของเมตริกซ์ดังนี้

$$y = x\beta + \varepsilon \quad (2.2)$$

เมื่อ

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \cdot \\ \cdot \\ \cdot \\ y_n \end{bmatrix}, \quad x = \begin{bmatrix} 1 & x_{11} & \dots & x_{1m} \\ 1 & x_{21} & \dots & x_{2m} \\ \cdot & \cdot & & \cdot \\ \cdot & \cdot & & \cdot \\ \cdot & \cdot & & \cdot \\ 1 & x_{n1} & \dots & x_{nm} \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \cdot \\ \cdot \\ \cdot \\ \beta_m \end{bmatrix}, \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \cdot \\ \cdot \\ \cdot \\ \varepsilon_n \end{bmatrix}$$

โดยมีข้อสมมติคือ

1. $E(\varepsilon) = 0$
2. $E(\varepsilon'\varepsilon) = \sigma^2 I_n$
3. X เป็นชุดของค่าตัวเลขคงที่
4. และมีอันดับ $m < n$

จากข้อสมมติที่ 1 $E(\varepsilon) = 0$ นั่นคือ ตัวแปร ε_i จะมีค่าคาดหวังเป็นศูนย์สำหรับทุกค่าของ $i, i = 1, \dots, n$

จากข้อสมมติที่ 2 จะได้ว่า

$$E(\varepsilon'\varepsilon) = \begin{bmatrix} E(\varepsilon_1^2) & E(\varepsilon_1\varepsilon_2) & \dots & E(\varepsilon_1\varepsilon_n) \\ E(\varepsilon_2\varepsilon_1) & E(\varepsilon_2^2) & \dots & E(\varepsilon_2\varepsilon_n) \\ \cdot & \cdot & & \cdot \\ \cdot & \cdot & & \cdot \\ E(\varepsilon_n\varepsilon_1) & E(\varepsilon_n\varepsilon_2) & \dots & E(\varepsilon_n^2) \end{bmatrix} \quad n \times n$$

$$E(\epsilon'\epsilon) = \begin{bmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{bmatrix} \quad n \times n$$

$$E(\epsilon'\epsilon) = \sigma^2 I_n$$

ตามข้อสมมติที่ 3 เมทริกซ์ X ต้องเป็นตัวเลขคงที่ หมายถึง การสุ่มตัวอย่างที่ซ้ำกันจะเกิดค่าหลายค่าในเวกเตอร์ y ซึ่งแปรผันในเวกเตอร์ ε และจำนวนค่าสังเกตของ X จะต้องมีจำนวนเกินจำนวนพารามิเตอร์ที่จะประมาณ ตามข้อสมมติที่ 4 ทั้งนี้ จะต้องไม่มีความสัมพันธ์เชิงเส้นระหว่างตัวแปรอิสระด้วย เราเรียกความสัมพันธ์เชิงเส้นระหว่างตัวแปรอิสระว่า multicollinearity

2.2 วิธีกำลังสองน้อยที่สุดแบบธรรมดา (Ordinary least square method)

วิธีการประมาณค่าสัมประสิทธิ์การถดถอยเชิงเส้นวิธีนี้ มีหลักในการประมาณค่าสัมประสิทธิ์คือ ทำให้ผลบวกกำลังสองของค่าผิดพลาดมีค่าน้อยที่สุด ซึ่งแสดงรายละเอียดดังนี้

นิยาม จากสมการ $y = x\beta + \epsilon$ เมื่อ $\epsilon \sim N(0, \sigma^2 I_n)$ ตัวประมาณกำลังสองน้อยที่สุดของ β คือ $\hat{\beta}$ ที่ทำให้ผลบวกกำลังสองของค่าผิดพลาด (Sum square errors) หรือ SSE มีค่าน้อยที่สุด

จากนิยาม จะทำการหาตัวประมาณค่ากำลังสองน้อยที่สุดได้ดังนี้

กำหนด $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k)'$
 จะได้ว่า $y = X\hat{\beta}$ (2.3)

ดังนั้น ϵ เป็นเวกเตอร์แนวตั้งของค่าผิดพลาด n ค่าจาก (2.3)

ผลบวกกำลังสองของค่าผิดพลาดคือ

$$SSE = \epsilon'\epsilon$$

$$SSE = (y - X\hat{\beta})'(y - X\hat{\beta})$$

$$SSE = y'y - 2\hat{\beta}'X'y + \hat{\beta}'X'X\hat{\beta}$$

การหาค่าน้อยที่สุดของผลบวกกำลังสองของค่าผิดพลาดทำได้โดยการหาอนุพันธ์ (differentiate) เทียบกับ β แล้วกำหนดให้เท่ากับ 0 ดังนี้

$$\begin{aligned} \frac{\partial}{\partial \beta} (y'y - 2\beta x'y + \beta^2 x'x) &= 0 \\ -2x'y + 2\beta x'x &= 0 \\ \beta &= (x'x)^{-1} x'y \end{aligned} \quad (2.4)$$

เมื่อ $(x'x)^{-1}x'$ เป็นเมตริกซ์ของค่าคงที่ สมาชิกใน β เป็นฟังก์ชันเชิงเส้นของ y นั่นคือ β เป็นตัวประมาณเชิงเส้น และแทนสมการ (2.2) ในสมการ (2.4) จะได้สมการในรูปแบบใหม่ คือ

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

จะได้ว่า $E(\beta) = \beta + (x'x)^{-1} x' E(\epsilon) = \beta$ ถ้า $E(\epsilon) = 0$ นั่นคือ β เป็นค่าที่ไม่เอนเอียง และ $V(\beta) = E(\beta - \beta)(\beta - \beta)'$

$$\begin{aligned} &= (x'x)^{-1} x' E(\epsilon\epsilon') x'(x'x)^{-1} \\ &= \sigma^2 (x'x)^{-1} \quad \text{โดยใช้ } E(\epsilon\epsilon') = \sigma^2 I_n \end{aligned}$$

เนื่องจากข้อมูลของตัวแปรตามที่จะใช้ศึกษาในครั้งนี้มีเพียง 2 ค่า คือ 1 หรือ 0 จึงมีความจำเป็นจะต้องมีการปรับค่าตัวแปรตามให้เป็น 1 หรือ 0 ตามสัดส่วนที่ต้องการ สำหรับเกณฑ์ที่เหมาะสมในการจำแนกค่าพยากรณ์ในสมการที่ (2.3) ออกเป็น 2 กลุ่ม โดยวิธีนี้ คือ

$$y_i = \begin{cases} 1 & \text{เมื่อ } y_i \geq y^* \\ 0 & \text{เมื่อ } y_i < y^* \end{cases}$$

โดยที่

$$\begin{aligned} y^* &= \frac{1}{2} \left[\frac{\sum_{i=1}^{n_1} y_{i1}}{n_1} + \frac{\sum_{i=1}^{n_0} y_{i0}}{n_0} \right] \\ &= \frac{1}{2} (\hat{y}_1 + \hat{y}_0) \end{aligned}$$

- \hat{y}_1 = ค่าเฉลี่ยของค่าพยากรณ์ที่ได้จากสมการถดถอยพหุ เมื่อตัวแปรตามมีค่าเป็น 1
 \hat{y}_0 = ค่าเฉลี่ยของค่าพยากรณ์ที่ได้จากสมการถดถอยพหุ เมื่อตัวแปรตามมีค่าเป็น 0
 n_1 = จำนวนตัวอย่าง เมื่อตัวแปรตามมีค่าเป็น 1
 n_0 = จำนวนตัวอย่าง เมื่อตัวแปรตามมีค่าเป็น 0
 $n_1 + n_0 = n$

2.3 การวิเคราะห์การถดถอยแบบทวิ (Binary regression)

การวิเคราะห์การถดถอยทวิเป็นเทคนิคที่พัฒนาขึ้นมาเพื่อแก้ไขปัญหาและข้อขัดแย้งที่เกิดจากการวิเคราะห์การถดถอย โดยวิธีกำลังสองน้อยที่สุดแบบธรรมดา สำหรับการวิเคราะห์การถดถอยทวิจะมีลักษณะคล้ายคลึงกับการวิเคราะห์ถดถอย โดยใช้วิธีกำลังสองน้อยที่สุดทั่ว ๆ ไป แต่เพิ่มเติมขั้นตอนการแปลงข้อมูลซึ่งได้จากวิเคราะห์ถดถอย โดยวิธีกำลังสองน้อยที่สุด เพื่อให้ค่าพยากรณ์ของตัวแปรตามที่ประมาณได้มีค่าอยู่ในช่วง $[0, 1]$

Arnold Zellner และ Tong Hun Lee (1965) ได้เสนอวิธีการประมาณค่าสัมประสิทธิ์การถดถอยแบบทวิ และการแปลงค่าพยากรณ์ให้อยู่ในช่วง $[0, 1]$ วิธีการที่ Zellner นำเสนอเป็นวิธีที่สามารถนำมาประยุกต์ใช้ เพื่อแก้ไขปัญหาดังกล่าวที่กำลังประสบอยู่ขณะนี้ สำหรับการศึกษาวิจัยนี้จะเลือกรูปแบบการแปลงข้อมูลให้สอดคล้องกับการแจกแจงที่กำหนด ด้วยเส้นโค้งของการแจกแจงต่าง ๆ ที่สำคัญ 3 รูปแบบ คือ

1. Normit Model (Normal Curve)
2. Logit Model (Logistic Curve)
3. Gompit Model (Gompertz Curve)

ในการประมาณค่าสัมประสิทธิ์การถดถอยแบบทวินี้ มีความจำเป็นที่จะต้องอาศัยวิธีกำลังสองน้อยที่สุดแบบทั่วไป (Generalized least square) มาช่วยในการประมาณค่าพารามิเตอร์ ดังนี้

วิธีกำลังสองน้อยที่สุดแบบทั่วไป (Generalized least square method: GLS)

ปัญหาหนึ่งที่เราพบในการศึกษาวิจัยครั้งนี้พบว่า เกิดปัญหา Heteroscedasticity เราพบว่า $E(e_i e_j) \neq 0$ สำหรับทุกค่าของ $i \neq j$ และ $E(e_i e_j) \neq \sigma^2$ สำหรับทุกค่าของ $i = j$

นั่นคือ $E(\epsilon' \epsilon) = \sigma^2 V$ โดยที่ V เป็นเมตริกซ์บวกแน่นอน (Positive definite) จึงสามารถที่จะหาเมตริกซ์ P ที่ทำให้ $PP' = V$ ได้

จากสมการ $y = x\beta + \epsilon$ สามารถแปลงให้อยู่ในรูปแบบใหม่ได้ดังนี้

$$y^* = x^* \beta + \varepsilon^* \quad (2.5)$$

โดยที่ $y^* = P^{-1}y$

$$x^* = P^{-1}x$$

และ $\varepsilon^* = P^{-1}\varepsilon$

จากสมการ (2.5) จะได้ว่า

$$E(\varepsilon^{*'} \varepsilon^*) = \sigma^2$$

$$E(\varepsilon^{*'} \varepsilon^*) = \sigma^2 I_n$$

ดังนั้น เราจึงสามารถประมาณค่าพารามิเตอร์ β ได้โดยวิธีกำลังสองน้อยที่สุดตั้งที่กล่าวมาแล้วในหัวข้อ 2.2 และ ค่าประมาณของ β สามารถเขียนได้เป็น

$$\begin{aligned} \hat{\beta}_{GLS} &= (x^{*'} x^*)^{-1} x^{*'} y^* \\ &= (x' V^{-1} x)^{-1} (x' V^{-1} y) \end{aligned}$$

$$\begin{aligned} V(\hat{\beta}_{GLS}) &= \sigma^2 (x^{*'} x^*)^{-1} \\ &= \sigma^2 (x' V^{-1} x)^{-1} \end{aligned}$$

โดยปกติใน GLS เราถือว่า σ^2 เป็นพารามิเตอร์ที่ไม่ทราบค่า ในขณะที่ V เป็นเมตริกซ์ที่เราต้องทราบค่าของสมาชิกทุกตัวเพื่อประมาณค่า β และความแปรปรวนของ β ในทางปฏิบัติเรามักไม่ทราบค่าของ V จะทราบได้ก็โดยการประมาณค่า V ด้วยวิธีที่เหมาะสมสำหรับการศึกษาคั้งนี้จะใช้ค่าประมาณที่เสนอโดยวิธีของ Arnold Zellner และ Tong Hun Lee (1985) การประมาณค่า V ด้วย \hat{V} นี้เรียกว่า Estimated Generalized Least Square (EGLS)

สำหรับการประมาณค่าโดยวิธีกำลังสองน้อยที่สุดในหัวข้อ 2.2 นั้น ทำให้เกิดปัญหาหลายประการ เพื่อแก้ไขปัญหาดังกล่าว Zellner และ Lee จึงได้เสนอวิธีการแปลงค่าพหุคูณของตัวแปรตามในรูปของความน่าจะเป็นด้วยเส้นโค้งของการแจกแจงแบบต่าง ๆ 3 รูปแบบ โดยมีขั้นตอนในการประมาณค่าดังนี้

ทั้ง 3 รูปแบบมีขั้นตอนในการประมาณค่าต่าง ๆ ดังนี้

1.) ประมาณค่า β ด้วยวิธี OLS จะได้ $\hat{\beta} = (x'x)^{-1}x'y$

2.) กำหนดให้ ξ เป็นค่าตัวกลาง (Intermediate value) ที่ได้จากสมการประมาณค่า : $y = x\beta$ โดยแทนค่า y ด้วย ξ ในขั้นตอนนี้จะได้ค่า ξ บางค่าอยู่นอกช่วง $[0, 1]$ ซึ่งไม่สอดคล้องกับค่าที่แท้จริงของตัวแปรตาม

3.) แปลงค่าพหุคูณ ξ จากขั้นตอนที่ 2 ให้อยู่ในช่วง $[0, 1]$ ในรูปแบบของความน่าจะเป็นด้วยเส้นโค้งของการแจกแจง ต่อไปนี้

$$3.1 \text{ Normit Model : } \hat{p}_i = 1/\sqrt{2\pi} \int^{t_i} e^{-1/2z^2} dz$$

$$3.2 \text{ Logit Model : } \hat{p}_i = 1/(1+e^{-t_i})$$

$$3.3 \text{ Gompit Model : } \hat{p}_i = e^{t_i}$$

4). การนำค่า $\hat{\beta}$ ที่ได้จากขั้นตอนที่ 1 มาใช้ในการแปลงค่าพยากรณ์ จะทำให้ค่าพยากรณ์ \hat{y} และ $\hat{\beta}$ ที่ได้ยังไม่มีความเหมาะสม ดังนั้นจึงต้องประมาณค่า $\hat{\beta}$ ค่าใหม่ขึ้นมาโดยอาศัยค่า t_i จากขั้นตอนที่ 1 และค่าประมาณของ V ที่เสนอโดย Zellner และ Lee

ดังนั้น จะได้ค่าประมาณสัมประสิทธิ์การถดถอยพหุคูณใหม่ ในแต่ละวิธี ดังนี้

$$\hat{\beta}_{EGLE}^* = (X' \hat{V}^{-1} X)^{-1} X' \hat{V}^{-1} t_i$$

เนื่องจาก $E(e_i e_j) = \sigma^2 V$ และ V เป็น Unknown Matrix จำเป็นต้องประมาณค่า V ด้วย \hat{V} ก่อนการวิเคราะห์

สำหรับเมตริกซ์ X และเวกเตอร์ t_i จะเหมือนกันสำหรับแต่ละรูปแบบ แต่จะต่างกันที่เมตริกซ์ \hat{V} (Diagonal Matrix) ค่าของ \hat{V}_i ในแต่ละวิธีคำนวณได้จาก

Normit Model:

$$\hat{V}_i = (\hat{p}_i \hat{q}_i) / n_i [Z(\hat{p}_i)]^2 \dots, Z(\hat{p}_i) = 1/\sqrt{2\pi} e^{-1/2z_i^2}$$

Logit Model:

$$\hat{V}_i = 1/n_i \hat{p}_i \hat{q}_i$$

Gompit Model:

$$\hat{V}_i = \hat{p}_i \hat{q}_i / (\hat{p}_i \ln \hat{p}_i)^2$$

สำหรับค่า t_i , \hat{p}_i และ \hat{q}_i ได้จากขั้นตอนที่ 2 และ 3

5). สมการในการประมาณค่าชุดใหม่ คือ $t_i^* = X \hat{\beta}^*$ และนำค่า t_i^* ที่ได้ขึ้นมาหาค่าประมาณของ \hat{p}_i โดยใช้การแปลงข้อมูลตามขั้นตอนที่ 3 กำหนดให้เป็นค่าเดิมคือ \hat{p}_i

๖). นำค่า \hat{P}_i ที่ได้มาใช้เพื่อพยากรณ์ตัวแปรตามว่ามีค่าเป็น 1 หรือ 0 โดยใช้เกณฑ์เดียวกับการวิเคราะห์การถดถอยโดยวิธีกำลังสองน้อยที่สุดแบบธรรมดา แต่เปลี่ยนค่าพยากรณ์จาก \hat{y} เป็น \hat{P} ดังนี้

$$y_i = \begin{cases} 1 & \text{เมื่อ } P_i \geq P^* \\ 0 & \text{เมื่อ } P_i < P^* \end{cases}$$

โดยที่
$$P^* = \frac{1}{2} (\hat{P}_1 + \hat{P}_0)$$

2.4 การวิเคราะห์จำแนกประเภท (Discriminant analysis)

การจำแนกกลุ่มและจัดเข้ากลุ่มกรณี 2 ประชากร เป็นวิธีการที่มุ่งจำแนกวัตถุออกเป็น 2 กลุ่ม ซึ่งเราทราบจำนวนกลุ่มชัดเจนแล้ว ในการแบ่งกลุ่มจะต้องมีตัวแปรอิสระที่จะมาช่วยในการสร้างเกณฑ์ เพื่อจำแนกกลุ่มและจัดสมาชิกใหม่ที่เพิ่มเข้ากลุ่มในกรณีนี้ถือว่ามี 2 ประชากร ดังนั้น จะเก็บข้อมูลจากประชากร π_1 มา n_1 หน่วยและจากประชากร π_0 มา n_0 หน่วย ($n_1 + n_0 = n$) ในการเก็บข้อมูลจะบันทึกคุณลักษณะประจำตัวของวัตถุแต่ละหน่วย คุณลักษณะดังกล่าวอาจจะมีมากกว่า 1 คุณลักษณะ ซึ่งเรียกว่า จำนวนตัวแปรอิสระ (P) ในการศึกษาวิจัยนี้จะถือว่าตัวแปรอิสระมาจากประชากรที่มีการแจกแจงแบบเบ้ และการจัดเข้ากลุ่มจะนำเอาความสูญเสียจากการจัดเข้ากลุ่มผิด (Cost of Misclassification) และ ความน่าจะเป็นโดยหลักเกณฑ์ (Prior Probability) มาเป็นเกณฑ์ร่วมกับ Fisher's Discriminant function (เป็นกรณีเฉพาะของเกณฑ์ใหม่) การจัดเกณฑ์ในลักษณะนี้เรียกว่า "Expected Cost of Misclassification (ECM)"

การวิเคราะห์จำแนกกลุ่มในกรณีที่มีประชากร 2 กลุ่ม

ถ้าการจำแนกกลุ่มเรามีตัวแปรที่ต้องการศึกษาอยู่ P ตัว มีประชากร 2 กลุ่ม คือ π_1 และ π_0 สามารถเขียนเวกเตอร์ของตัวแปรอิสระได้ดังนี้

$$X' = [x_1, x_2, x_3, \dots, x_p]$$

สำหรับเกณฑ์การจัดกลุ่มนี้ฟิชเชอร์ (Fisher, R.A., 1938) แนะนำให้ใช้ข้อมูลจากตัวแปรทุกตัว โดยถ่วงน้ำหนักด้วยวิธีเหมาะสม นั่นคือให้ใช้การประกอบกันของ

$x_1, x_2, x_3, \dots, x_p$ คือ

$$y = l_1 x_1 + l_2 x_2 + l_3 x_3 + \dots + l_p x_p = l' x$$

โดยที่ l คือ เวกเตอร์ของน้ำหนักที่มีผลให้อัตราส่วนต่อไปนี้มีค่าสูงสุด

$$\sigma = \frac{(\text{ระยะทางระหว่างค่าเฉลี่ยของ } Y)^2}{\text{ความแปรปรวนของ } Y}$$

$$\sigma = \frac{(\mu_{1y} - \mu_{0y})^2}{l' \Sigma l}$$

$$= \frac{l' (\mu_1 - \mu_0) (\mu_1 - \mu_0)' l}{l' \Sigma l} \quad (2.6)$$

กำหนดให้ π_1 เป็นประชากรกลุ่มที่ตัวแปรตามมีค่าเป็น 1

π_0 เป็นประชากรกลุ่มที่ตัวแปรตามมีค่าเป็น 0

$\mu_1 = E(x/\pi_1)$ ค่าเฉลี่ยของเวกเตอร์ x ที่บันทึกจากประชากร π_1

$\mu_0 = E(x/\pi_0)$ ค่าเฉลี่ยของเวกเตอร์ x ที่บันทึกจากประชากร π_0

ดังนั้น ค่าเฉลี่ยของ Y ที่เกิดจาก linear combination ของ X 's จากประชากร π_1 และ π_0 จึงมีค่าเป็น μ_{1y} และ μ_{0y} ตามลำดับ ดังนี้

$$\mu_{1y} = E(Y/\pi_1) = E(l' x/\pi_1) = l' \mu_1$$

$$\mu_{0y} = E(Y/\pi_0) = E(l' x/\pi_0) = l' \mu_0$$

ส่วนความแปรปรวนของ Y จากทั้ง 2 ประชากร คือ

$$\begin{aligned} \sigma_y^2 &= \text{var}(l' x) = l' \text{Cov}(x) l \\ &= l' \Sigma l \end{aligned}$$

โดยที่

$$\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2p} \\ \cdot & & & \cdot \\ \cdot & & & \cdot \\ \sigma_{p1} & \sigma_{p2} & \dots & \sigma_{pp} \end{bmatrix} ; l' = [l_1, l_2, \dots, l_p]$$

การหาสมการที่เหมาะสมในการแบ่งกลุ่มประชากรทั้งสองออกจากกันให้ได้มากที่สุดคือ การหาค่า t ที่ทำให้ δ มีค่ามากที่สุด ในทางปฏิบัติค่า μ_1 , μ_0 และ Σ เป็นค่าที่เรามักจะไม่มีทราบ ดังนั้น ถ้าสุ่มตัวอย่างจาก τ_1 และ τ_0 มาจำนวน n_1 และ n_0 ตามลำดับแล้ววัดค่าสังเกตเพื่อประมาณค่า μ_1 , μ_0 และ Σ

ถ้าแบ่งเมตริกซ์ X ออกเป็น 2 ส่วน ดังนี้

$$X' = \begin{bmatrix} x_{111} & x_{112} & \dots & x_{11n} & : & x_{211} & x_{212} & \dots & x_{21n} & : \\ & & & & : & & & & & : \\ x_{121} & x_{122} & \dots & x_{12n} & : & x_{221} & x_{222} & \dots & x_{22n} & : \\ & & & & : & & & & & : \\ \cdot & \cdot & & & : & \cdot & \cdot & & & : \\ \cdot & \cdot & & & : & \cdot & \cdot & & & : \\ \cdot & \cdot & & & : & \cdot & \cdot & & & : \\ x_{1p1} & x_{1p2} & \dots & x_{1pn} & : & x_{2p1} & x_{2p2} & \dots & x_{2pn} & : \\ & & & & : & & & & & : \end{bmatrix}$$

$$\text{หรือ } X = [x_1' \mid x_0']$$

x_1' และ x_0' เป็นเมตริกซ์ของค่าสังเกต ซึ่งแต่ละเมตริกซ์จะได้เวกเตอร์ค่าเฉลี่ยดังนี้

$$x_1' = [x_{11} \quad x_{12} \quad \dots \quad x_{1p}]$$

$$x_0' = [x_{01} \quad x_{02} \quad \dots \quad x_{0p}]$$

และประมาณค่า Covariance Matrix Σ ด้วย S_p คือ Pooled Covariance Matrix

$$S_p = \frac{(n_1 - 1) S_1 + (n_0 - 1) S_0}{n_1 + n_0 - 2}$$

เมื่อ

$$S_1 = \frac{1}{n_1 - 1} \sum_{j=1}^{n_1} (x_{1j} - x_1)(x_{1j} - x_1)'$$

$$S_0 = \frac{1}{n_0 - 1} \sum_{j=1}^{n_0} (x_{0j} - x_0)(x_{0j} - x_0)'$$

จากสมการ (2.6) จะได้

$$= \frac{1'(\bar{x}_1 - \bar{x}_0)(\bar{x}_1 - \bar{x}_0)'1}{1'S_1 1}$$

จะมีค่ามากที่สุด เมื่อ $\frac{\partial \hat{\delta}}{\partial 1} = 0$

\therefore จะได้ $\hat{1} = c S_p^{-1} (\bar{x}_1 - \bar{x}_0)'$; $c = 0$

และเมื่อ $c = 1$ จะได้เวกเตอร์ 1 เรียกว่า Fisher's linear discriminant function ดังนี้

$$\begin{aligned} g &= 1'x \\ &= (\bar{x}_1 - \bar{x}_0)' S_p^{-1} x \end{aligned}$$

ซึ่งจะมีผลให้อัตราส่วน δ มีค่าสูงสุด

ในการจำแนกกลุ่มโดยวิธีการวิเคราะห์จำแนกประเภทสามารถสรุปเป็นขั้นตอนในการจำแนกกลุ่มได้ดังต่อไปนี้ โดยจะพิจารณาค่า Prior Probability และ Cost of Misclassification ประกอบการจำแนกกลุ่มด้วย

1. คำนวณสมการจำแนกกลุ่ม

Fisher's Linear Discriminant function

$$g = 1'x = (\bar{x}_1 - \bar{x}_0)' S_p^{-1} x$$

2. คำนวณค่า y_{new} คือ ค่าของ discriminant function สำหรับค่าสังเกต x_{new} หน่วยใหม่ที่ปรากฏภายหลัง หรือค่าสังเกตที่ต้องการพยากรณ์ โดยที่แทนค่า x ด้วย x_{new} ในสมการขั้นตอนที่ 1

3. คำนวณค่า Midpoint (\hat{m}) คือ $\hat{m} = \frac{1}{2} (\bar{x}_1 - \bar{x}_0)' S_p^{-1} (\bar{x}_1 + \bar{x}_0)$

4. จัด x_{new} เข้ากลุ่ม $r1$ หรือ $r0$ โดยใช้เกณฑ์ร่วมใหม่ที่พยายามทำให้ค่า EMC มีค่าต่ำที่สุด

ถ้า $y_{new} - \hat{m} \geq \ln \frac{c(1/0) \cdot p_0}{c(0/1) \cdot p_1}$ จะจัด x_{new} เข้ากลุ่ม $r1$

ถ้า $y_{n+1} - m < \ln \frac{c(1/0)}{c(0/1)} \cdot P_0$ จะจัด x_{n+1} เข้ากลุ่ม π_0

อสมการนี้ เรียกว่า Anderson's classification rule

$c(1/0)$ ความสูญเสียที่เกิดจากการจัดวัตถุเข้ากลุ่ม π_1 ทั้งที่ความจริงวัตถุ
นั้นมาจากกลุ่ม π_0

$c(0/1)$ ความสูญเสียที่เกิดจากการจัดวัตถุเข้ากลุ่ม π_0 ทั้งที่ความจริงวัตถุ
นั้นมาจากกลุ่ม π_1

การวิจัยครั้งนี้จะไม่นำค่า $c(1/0)$ และ $c(0/1)$ มาพิจารณา เพราะในทาง
ปฏิบัติการที่จะกำหนดเท่าได้นั้น จะขึ้นอยู่กับผู้วิจัยและลักษณะของงานวิจัย ดังนั้นจึงถือว่า
 $c(1/0)$ และ $c(0/1)$ มีค่าเท่ากัน แต่ในทางปฏิบัติควรจะคำนึงความสูญเสีย ดังกล่าว
เสมอ

Prior Probability : P_1, P_0 หมายถึง ความน่าจะเป็นโดยหลักเกณฑ์
โอกาสที่จะเป็นสมาชิกของกลุ่ม π_1 และ π_0 ตามลำดับ จะนำมาพิจารณาในการวิจัยนี้
เนื่องจากทราบความแตกต่างในขนาดของประชากร เพราะว่าถ้าประชากรกลุ่มใดมีขนาด
ใหญ่กว่าโอกาสที่วัตถุนั้นจะเป็นสมาชิกของกลุ่มนั้นย่อมมีมากกว่า