

## บทที่ 9

### ตัวแปรตามที่เป็นตัวแปรเชิงคุณภาพหรือมีค่าจำกัด (Qualitative and Limited Dependent Variables)

#### 9.1 ความหมายและความสำคัญ

คำว่า Qualitative Dependent Variable หมายถึงกรณีที่ตัวแปรตาม  $Y$  เป็นตัวแปรเชิงคุณภาพ หรือ Categorical Variable ซึ่งนักวิจัยจำเป็นต้องกำหนดดัชนีสำหรับค่าตัวแปรดังกล่าวในรูปของตัวแปรตัวมี หรือกล่าวอย่างไม่เป็นทางการได้ว่า Qualitative Dependent Variable ก็คือตัวแปรตามที่เป็นตัวแปรตัวมีนั่นเอง เราเรียกสมการ  $Y = f(X's)$  ที่  $Y$  เป็น Qualitative Dependent Variable ว่า Quantal Response Model

คำว่า Limited Dependent Variable คือตัวแปรตาม  $Y$  ที่มีค่าจำกัดอยู่เฉพาะในบางช่วงของ Real Line เช่น  $Y > 0$  (เรียกว่าปัญหา Truncated Sample) หรือตัวแปร  $Y$  มีค่าได้ตลอดช่วงของ Real Line ยกเว้นบางช่วง (เรียกว่าปัญหา Censored Sample)

ทั้งกรณี Qualitative Dependent Variable และ Limited Dependent Variable เป็นสถานการณ์ที่มีผลให้ไม่อาจใช้ OLS -Estimator ตามปกติได้เพราะ OLS -Estimator ในกรณีเหล่านี้จะเป็น Biased Estimator และ Inconsistent สำหรับการศึกษานี้เป็นการศึกษาถึงกรณีของ Qualitative Dependent Variable เท่านั้น ผู้สนใจสามารถศึกษาเรื่อง Limited Dependent Variable ได้จากเอกสารอ้างอิง

พิจารณาสมการถดถอย  $Y_i = \alpha + \beta X_i + u_i$  ;  $i = 1, 2, \dots, M$  โดยที่  $Y$  เป็นตัวแปรเชิงสุ่มที่มีค่าเป็น 0 และ 1 เท่านั้น เช่น  $Y_i = 1$  ถ้าหน่วยสำรวจที่  $i$  คือ ผู้ป่วยที่เป็นมะเร็งบริเวณศีรษะและลำคอ และ  $Y_i = 0$  ถ้าหน่วยสำรวจคือผู้ป่วยที่เป็นมะเร็งบริเวณอื่น ๆ หรือ  $Y_i = 0$  ถ้าหน่วยสำรวจคือผู้มีเงินได้ที่มีการศึกษาระดับอุดมศึกษา และ  $Y_i = 1$  ถ้าหน่วยสำรวจคือผู้มีเงินได้ที่มีการศึกษาระดับอื่น ๆ

จากตัวอย่างนี้จะเห็นได้ว่า  $Y_i$  มีค่าได้ 2 ค่าคือ 0 และ 1 ผลสะท้อนที่ตามมาก็คือ  $u_i$  จะมีค่าได้ 2 ค่าคือ  $u_i = Y - (\alpha + \beta X_i) = 1 - (\alpha + \beta X_i)$  ถ้า  $Y_i = 1$  และ  $u_i = -(\alpha + \beta X_i)$  ถ้าปัญหาที่ปรากฏอย่างเด่นชัดก็คือ

1. เราไม่อาจควบคุมให้  $E(u_i) = 0$  ได้โดยง่าย กล่าวคือถ้าต้องการให้ข้อตกลงว่า  $E(u_i) = 0$  เป็นจริงเราจำเป็นต้องควบคุมให้  $u_i$  ปรากฏค่าด้วยความน่าจะเป็น  $(\alpha + \beta X_i) = p_i$  หรือ  $1 - (\alpha + \beta X_i) = q_i$  เท่านั้น<sup>1</sup> แต่การควบคุมให้  $0 < p_i < 1$  และ  $0 < q_i < 1$  เป็นเรื่องที่เป็นไปได้ยากเพราะโดยปกติ  $(\alpha + \beta X_i)$  อาจมีค่ามากกว่า 1 หรือแม้กระทั่งเป็นประมาณติดลบ คือ น้อยกว่า 0 ได้

2.  $V(u_i)$  จะไม่คงที่ เพราะในกรณีนี้ตัวแปรสุ่ม  $Y$  มีการแจกแจงแบบเบอร์นูลลี (Bernoulli Distribution) ดังนั้น  $V(u_i) = p_i q_i = (\alpha + \beta X_i)(1 - \alpha - \beta X_i)$  ซึ่งจะไม่คงที่แต่จะแปรเปลี่ยนไปตามค่าของตัวแปร  $X$  ปัญหา Heteroscedasticity จึงเกิดขึ้น

3. สมการประมาณค่า  $Y_i = \hat{\alpha} - \hat{\beta} X_i$  จะหวั่นไหว (Sensitive) ได้ง่ายเมื่อ  $X$  มีค่าเปลี่ยนแปลงไป

4. เราไม่อาจใช้ t-test, F-test ได้รวมทั้ง  $R^2$  ก็ขาดความหมายที่แท้จริง และไม่ใช้เป็นดัชนีในการตัดสินใจได้อีกทั้ง  $\hat{\sigma}^2$  ยังเป็น Inconsistent Estimator ของ  $\sigma^2$  นอกจากนี้แล้วเรายังไม่อาจใช้วิธีการประมาณค่าต่าง ๆ ที่อาศัย Linear Function ของ  $Y_i$  มาประมาณค่าพารามิเตอร์ เช่น OLS, GLS ได้เพราะ  $Y$  มิได้แจกแจงแบบปกติ เว้นแต่จะได้รับการปรับโครงสร้างและสภาพทั่วไปเสียก่อน

## 9.2 ตัวอย่าง

เพื่อความเข้าใจปัญหาและติดตามการพัฒนาทางทฤษฎีได้โดยไม่ติดขัดผู้เขียนขอยกตัวอย่างงานวิจัยที่ดำเนินการโดยที่ตัวแปร  $Y$  เป็นตัวแปรเชิงคุณภาพ (Dichotomous หรือ Binary Choice Model) ดังนี้

จากข้อมูลเกี่ยวกับความพอใจของทหารว่าพึงพอใจที่จะประจำกองพลในภาคเหนือ หรือภาคใต้ ( $Y_i = 1$  ถ้าทหารนายที่  $i$  พึงพอใจจะประจำกองพลในภาคเหนือ  $Y_i = 0$  ถ้าทหารนายที่  $i$  พึงพอใจจะประจำกองพลในภาคใต้) โดยที่ตัวแปรอิสระคือ  $X_2 =$  ภูมิลำเนาเดิม ( $X_2 = 1$  ถ้ามีภูมิลำเนาเดิมอยู่ภาคเหนือ  $X_2 = 0$  ถ้าภูมิลำเนาเดิมอยู่ภาคใต้)  $X_3 =$  ที่ตั้งปัจจุบันของกองพลที่สังกัด ( $X_3 = 1$  ถ้าสังกัดปัจจุบันอยู่ในภาคเหนือ  $X_3 = 0$  ถ้าสังกัดปัจจุบันอยู่ในภาคใต้)

---

<sup>1</sup> ตัวแปร  $Y_i$  มีการแจกแจงแบบเบอร์นูลลี กล่าวคือ โดยทั่วไปถ้า  $X$  มี pdf เป็น  $f(x) = p^x q^{1-x}$ ;  $x = 0, 1$  เราเรียกว่า  $f(x)$  ว่า Bernoulli Distribution โดยที่ถ้า  $x = 0$  จะได้  $f(x) = q$  และถ้า  $x = 1$  จะได้  $f(x) = p$

ผลการบันทึกข้อมูลจากทหาร 8,036 นาย ปรากฏข้อมูลดังตาราง<sup>1</sup>

j	k	ℓ	X <sub>2</sub> ภูมิฐานะเดิม	X <sub>3</sub> ที่ตั้งสังกัดปัจจุบัน	Y ที่ตั้งสังกัดที่ต้องการ	ความถี่ (f <sub>jkℓ</sub> )
1	1	1	เหนือ	เหนือ	เหนือ	1342
1	1	0	เหนือ	เหนือ	ใต้	198
1	0	1	เหนือ	ใต้	เหนือ	1750
1	0	0	เหนือ	ใต้	ใต้	760
0	1	1	ใต้	เหนือ	เหนือ	487
0	1	0	ใต้	เหนือ	ใต้	446
0	0	1	ใต้	ใต้	เหนือ	472
0	0	0	ใต้	ใต้	ใต้	2581

และเมื่อจำแนกกลุ่มตัวอย่างจากตารางข้างต้นเพื่อพิจารณาความน่าจะเป็นจะพบว่าเราสามารถเสนอตารางได้ดังนี้

j	k	X <sub>2</sub> ภูมิฐานะเดิม	X <sub>3</sub> ที่ตั้งสังกัดปัจจุบัน	Y : ที่ตั้งสังกัดที่ต้องการ						รวม (n <sub>jk</sub> )
				ภาคเหนือ (ℓ=1)			ภาคใต้ (ℓ=0)			
				f <sub>jk1</sub>	p <sub>jk</sub>	ω <sub>jk</sub>	f <sub>jk0</sub>	q <sub>jk</sub>	$\frac{1}{\omega_{jk}}$	
1	1	เหนือ	เหนือ	1342	.871	<b>6.778</b>	198	.129	.148	1540
1	0	เหนือ	ใต้	1750	.697	2.303	760	.303	.434	2510
0	1	ใต้	เหนือ	487	.522	1.092	446	.478	.916	933
0	0	ใต้	ใต้	472	.155	0.183	2581	.845	5.468	3053
				4051	.504	1.020	3985	.496	.984	8036

<sup>1</sup>Goodman, Leo, A., Analyzing Qualitative Categorical Data • Log - Linear Models and Latent - Structure Analysis. (Addison • Wesley Publishing Company Advanced Book Program. London, 1978), p.27-56

### ข้อสังเกต

$$1. P_{jk} \text{ (หรือ } P_{jk}^{X_2 X_3}) = \text{Prob(ต้องการสังกัดภาคเหนือ} | x_2 = x_{2j} \text{ และ } x_3 = x_{3k})$$

$$= \frac{f_{jk1}}{n_{jk}} \text{ เช่น}$$

$$\hat{P}_{11} = \text{Pr(ต้องการสังกัดภาคเหนือ} | x_2 = 1, x_3 = 1) = \frac{1342}{1540} = .871$$

$$\hat{P}_{13} = \text{Pr(ต้องการสังกัดภาคเหนือ} | x_2 = 0, x_3 = 1) = \frac{487}{933} = .522$$

$$2. \omega_{jk} \text{ (หรือ } \omega_{jk}^{X_2 X_3}) = \text{Odds (ต้องการสังกัดภาคเหนือ} | x_2 = x_{2j}, x_3 = x_{3k})$$

$$= \frac{f_{jk1}}{f_{jk0}} = \text{อัตราส่วน (Ratio)}$$

$$\hat{\omega}_{11} = \text{Odds (ต้องการสังกัดภาคเหนือ} | x_2 = 1, x_3 = 1) = \frac{1342}{198} = 6.778$$

= 6.778 : 1 หรือทหารที่มีภูมิลำเนาอยู่ภาคเหนือที่สังกัดกองพลในภาคเหนือ แต่ต้องการสังกัดเหนือบ้างภาคใต้บ้างคิดเป็นอัตราส่วน 6.778 ต่อ 1

$$\hat{\omega}_{14} = \text{Odds (ต้องการสังกัดภาคเหนือ} | x_2 = 0, x_3 = 0) = \frac{472}{2581} = .183$$

3. Odds และ Prob สามารถใช้คำนวณย้อนสู่กันและกันได้ดังนี้

$$\text{Odds} = \frac{\text{Prob}}{(1-\text{Prob})}$$

$$\text{Prob} = \frac{\text{Odds}}{(1+\text{Odds})}$$

เช่น

$$\hat{\omega}_{11} = 6.778 \quad \hat{P}_{11} = \hat{\omega}_{11} / (1 + \hat{\omega}_{11}) = 6.778 / 7.778 = .871$$

$$\hat{\omega}_{14} = .183 \quad \hat{P}_{14} = \hat{\omega}_{14} / (1 + \hat{\omega}_{14}) = .183 / 1.183 = .155$$

$$\hat{P}_{11} = .871 \quad \hat{\omega}_{11} = \hat{P}_{11} / (1 - \hat{P}_{11}) = .871 / .129 = 6.751$$

$$\hat{P}_{13} = .522 \quad \hat{\omega}_{13} = \hat{P}_{13} / (1 - \hat{P}_{13}) = .522 / .478 = 1.155$$

หมายเหตุ ค่า  $\hat{\omega}_{11}$  และ  $\hat{\omega}_{13}$  ตามตัวอย่างนี้คลาดเคลื่อนไปเพราะการปัดเศษ ถ้าใช้ค่าที่ไม่ปัดเศษจะได้เท่ากับในตารางซึ่งคำนวณโดยตรงตามข้อ 1 และ 2

ตัวอย่างนี้เป็นตัวอย่างที่ชี้ให้เห็นว่าในการวิเคราะห์สมการถดถอยที่ตัวแปร Y เป็น Quantal เราจำเป็นต้องจำแนกค่าสังเกตออกเป็นกลุ่ม ๆ (Subsample) เรียกว่า Setting ของตัวแปรอิสระ จากนั้นจึงคำนวณหาความน่าจะเป็น p และ q หรือ Odds เพื่อนำไปใช้ประโยชน์ในการวิเคราะห์ต่อไป ในขั้นตอนนี้ผู้เขียนขอแสดงตัวอย่างนี้เพียงระดับการกำหนด Setting ของตัวแปรอิสระ<sup>1</sup> และการหาความน่าจะเป็นเท่านั้น ซึ่งคิดว่าน่าจะมีประโยชน์เพียงพอที่จะใช้เป็นพื้นความรู้ในการศึกษาเรื่องของ Probit และ Logit ต่อไป และจะเสนอการประมาณค่าพารามิเตอร์ต่อไปเมื่อศึกษาถึงเรื่องนั้น

### 9.3 การประมาณค่าสำหรับ Binary Choice Model

เมื่อ  $Y_i$  ในสมการ  $Y = f(X's)$  เป็นตัวแปรตมที่มีค่าได้ 2 ค่า (Binary Choice หรือ Dichotomous) คือ 0, 1 (หรือรหัสอื่น ๆ เช่น 1, 2 หรือ -1, 1 ตามที่ผู้วิจัยเห็นสมควร โดยที่ X's อาจเป็น Quantal หรือมิใช่ก็ได้<sup>2</sup> เรามีวิธีประมาณค่าได้หลายวิธีดังนี้

#### 9.3.1 Linear Probability Model

วิธีนี้เสนอโดยโกลด์เบอร์เกอร์, เซลล์เนอร์ และ ลี (Goldberger, A.S., 1964; Zellner, A and Lee, T.H., 1965)<sup>3</sup> วิธีปฏิบัติคือจากสมการ  $Y = X\beta + U$

ให้กำหนด Setting ของตัวแปรอิสระ  $X_2, X_3, \dots, X_k$  ออกเป็นส่วน ๆ รวม M ส่วน ๆ ละ

$n_i$  หน่วย ให้  $Y_i = \sum_j^{n_i} Y_{ij}$  เมื่อ  $Y_{ij} = 0, 1; i = 1, 2, \dots, M; j = 1, 2, \dots, n_i$  และ

$p_i = \frac{Y_i}{n_i}; i = 1, 2, \dots, M$  คือความน่าจะเป็นที่  $Y_i = 1$  ใน Setting ที่  $i; i = 1, 2, \dots, M$  หรือ  $p_i = \frac{Y_i}{n_i}$  ก็คือ Sample Proportion ใน Setting ที่  $i$  ของ X's

<sup>1</sup> การกำหนด Setting ของตัวแปรอิสระก็คือการจัดกลุ่ม (Combination) ของรหัสของตัวแปรตามตัวอย่างนี้  $X_2 = 0, 1$  ขณะที่  $X_3 = 0, 1$  เราจึงสามารถจัดกลุ่มค่าสังเกตของ  $X_2$  และ  $X_3$  ได้  $2^2$  กลุ่ม คือ (0, 0), (0, 1) (1, 0), (1, 1) ด้วยเหตุนี้ถ้ามีตัวแปรอิสระที่เป็นตัวแปรตมมีทั้งสิ้น  $m$  ตัว เราสามารถจัดกลุ่มค่าสังเกตของตัวแปรอิสระได้ทั้งสิ้น  $2^m$  กลุ่ม

<sup>2</sup> คำว่า Quantal มีความหมายเดียวกับคำว่า Discrete ในทางปฏิบัติเราใช้ในความหมายของ Dummy Variable

<sup>3</sup> Judge, G.G., et. al., Op. Cit., p.587



ดังนั้นเมื่อพิจารณาเมตริกซ์  $\Phi$  เราจึงพบว่า  $\Phi \neq \sigma^2 I_M$  หรือเกิดปัญหา Heteroscedasticity ทางออกก็คือประมาณค่า  $\beta$  โดยอาศัย GLS เมื่อทราบค่าของ  $\Phi$  กล่าวคือ

$$\hat{\beta} = (X'\Phi^{-1}X)^{-1}(X'\Phi^{-1}p)$$

และในกรณีที่ไมทราบค่าของ  $\Phi$  หรือนัยหนึ่งก็คือไมทราบค่าของ  $p_i$  ให้ดำเนินการประมาณค่าของ  $p_i$  เสียก่อนแล้วจึงประมาณค่าเวกเตอร์  $\beta$  โดยวิธี EGLS ตามขั้นตอนต่อไปนี้

ขั้นที่ 1 จากสมการ  $p = X\beta + u$  โดยที่  $p = [p_1, p_2, \dots, p_M]$  เมื่อ  $p_i = \frac{y_i}{n_i}$   $i = 1, 2, \dots, M$  ให้ประมาณค่าเวกเตอร์  $\beta$  ตามวิธี OLS ได้

$$\hat{\beta} = (X'X)^{-1}X'p$$

ขั้นที่ 2 แทนที่  $\beta$  ด้วย  $\hat{\beta}$  จะได้สมการประมาณค่า  $p = X\hat{\beta}$  หรือ  $\hat{p} = X\hat{\beta}$  สมาชิกของเวกเตอร์  $\hat{p}$  คือ  $\hat{p}_i$  ;  $i = 1, 2, \dots, M$

ขั้นที่ 3 ประมาณค่า  $v(u_i)$  โดยใช้  $\hat{p}_i$  แทน  $p_i$  ได้  $v(u_i) = \frac{\hat{p}_i \hat{q}_i}{n_i}$ ;  $i = 1, 2, \dots, M$  เมื่อ  $\hat{q}_i = 1 - \hat{p}_i$  และพบว่า

$$\hat{\Phi} = \begin{bmatrix} \frac{\hat{p}_1 \hat{q}_1}{n_1} & & & \\ & \frac{\hat{p}_2 \hat{q}_2}{n_2} & & \\ & & \ddots & \\ & & & \frac{\hat{p}_M \hat{q}_M}{n_M} \end{bmatrix} \quad \hat{\Phi}^{-1} = \begin{bmatrix} \frac{n_1}{\hat{p}_1 \hat{q}_1} & & & \\ & \frac{n_2}{\hat{p}_2 \hat{q}_2} & & \\ & & \ddots & \\ & & & \frac{n_M}{\hat{p}_M \hat{q}_M} \end{bmatrix}$$

ขั้นที่ 4 ใช้  $\hat{\Phi}$  ในการประมาณค่าเวกเตอร์  $\beta$  ตามวิธี EGLS กล่าวคือ

$$\hat{\beta} = (X'\hat{\Phi}^{-1}X)^{-1}(X'\hat{\Phi}^{-1}p), \quad \hat{V}(\hat{\beta}) = (X'\hat{\Phi}^{-1}X)^{-1}, \quad p = X\hat{\beta}$$

วิธีนี้เป็นวิธีที่ง่ายที่สุดแต่มีปัญหาซึ่งนับเป็นจุดอ่อน 2 ประการคือ

1. เราไม่อาจรับประกันได้ว่าค่า  $p_i$  จากขั้นที่ 2 จะมีค่าปรากฏในช่วง  $[0, 1]$  เพราะอาจเป็นไปได้ที่  $\hat{p}_i > 1$  หรือ  $\hat{p}_i < 0$

2. เราไม่อาจรับประกันได้ว่าค่าพยากรณ์ของ  $p_i$  จากสมการประมาณค่า  $p = x\hat{\beta}$  จะมีค่าปรากฏในช่วง  $[0, 1]$

วิธีแก้ปัญหาคือเป็นจุดอ่อนก็คือ ให้เปลี่ยนไปใช้เทคนิคการประมาณค่าแบบอื่น ขอให้สังเกตว่าเมื่อ  $Y$  เป็นตัวแปรตมมีเราจะไม่พยากรณ์ค่าของ  $Y$  แต่จะพยากรณ์สัดส่วนของ  $Y$  ในแต่ละ Setting ของ  $X$ 's

### 9.3.2 Transformation Approach

ทางออกสำหรับปัญหาค่า  $\hat{p}$  ไม่ตกอยู่ใน Unit Interval  $[0, 1]$  นั้นเราสามารถแก้ได้ 3 วิธีคือ วิธี Inequality - Restricted Least Square (IRLS) วิธี Normit analysis (หรือ Probit Analysis) และวิธี Logit Analysis

วิธี IRLS ก็คือวิธีกำหนดเงื่อนไขว่า  $p_i \leq 1$  และ  $p_i \geq 0$  ซึ่งก็คือ  $x_i'\beta \leq 1$  และ  $x_i'\beta \geq 0$  เมื่อ  $x_i'$  คือแถวที่  $i$  ของเมตริกซ์  $X$  ในสมการ  $p = X\beta + u$  นั้นเอง ค่าประมาณของ  $\beta$  ตามวิธี IRLS ก็คือค่าประมาณที่เกิดขึ้นจาก

minimize  $e'e$  ภายใต้ข้อจำกัดว่า

$$(1) \quad X\beta \leq 1_M \quad \text{เมื่อ} \quad 1_M = [1, 1, \dots, 1]'$$

$$(2) \quad X\beta \geq 0_M \quad \text{เมื่อ} \quad 0_M = [0, 0, \dots, 0]'$$

แต่ค่าประมาณของ  $\beta$  ตามวิธี IRLS ยังมีปัญหาคือ เรายังไม่ได้พัฒนา Sampling Property ของ  $\hat{\beta}$  ivo อย่างกว้างขวางเพียงพอแก่งานวิเคราะห์<sup>1</sup>

วิธีประมาณค่า  $\beta$  อีกวิธีหนึ่งที่สามารถแก้ปัญหาเรื่องค่าประมาณ  $\hat{p}$  ไม่ปรากฏใน Unit Interval  $[0, 1]$  ก็คือ Normit Analysis และ Logit Analysis

---

<sup>1</sup> อ่าน Judge, G.G., et al., Ibid., p.78-94, 588-591



หลักการโดยสรุปของทั้งวิธี Normit Analysis และ Logit Analysis ก็คือการแปลงค่า  $P_i$  และ  $p_i$  ให้เป็น Quantile<sup>1</sup> ของ Distribution ใดๆ ที่พิจารณาเห็นว่าเหมาะสมแล้วนำค่า Quantile ดังกล่าวไปใช้วิเคราะห์สมการถดถอย จากนั้นจึงแปลงค่า Quantile กลับสู่รูปความน่าจะเป็น  $P_i$  และ  $p_i$  ซึ่งมีผลให้ค่าประมาณและค่าพยากรณ์ของ  $P_i$  และ  $p_i$  ปรากฏอยู่ใน Unit Interval

### 9.3.2.1 Normit Analysis<sup>2</sup>

จากสมการความสัมพันธ์  $p_i = P_i + u_i$  ;  $i = 1, 2, \dots, M$  .... (1)

ให้ Normit ของ  $P_i$  คือ

Normit  $P_i = I_i$  เมื่อ  $I_i = X_i' \beta = \beta_1 + \beta_2 X_{2i} + \dots + \beta_k X_{ki}$  โดยถือว่า  $P_i$  มีความสัมพันธ์กับ  $I_i$  ดังนี้

$$P_i = F(I_i) = \int_{-\infty}^{I_i} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} t^2\right) dt^3$$

$$\text{หรือ } I_i = F^{-1}(P_i) \quad \dots (2)$$

จาก (1) คือ  $p_i = P_i + u_i$  ;  $i = 1, 2, \dots, M$  จะพบว่า

$$F^{-1}(p_i) = F^{-1}(P_i + u_i) \quad \dots (3)$$

โดยอาศัย Taylor's Series ของ  $f(x)$  ที่กระจายรอบ  $a$  คือ

$$f(x) = f(a) + (x-a) f'(a) + \frac{(x-a)^2}{2!} f''(a) + R$$

ดังนั้นเมื่อกระจาย  $F^{-1}(P_i + u_i)$  รอบ  $P_i$  จะพบว่า

<sup>1</sup> ถ้า  $P_i = \int_{-\infty}^{I_i} f_x(x) dx = \Pr(X < I_i) = F(I_i)$  เราเรียก  $I_i$  ว่าเป็น Quantile of Order  $P_i$  แสดงว่าการแปลงรูป  $P_i$  เป็น  $I_i$  เป็นกระบวนการที่อาศัย CDF ของตัวแปรสุ่มเป็นตัวกลาง และเนื่องจาก CDF คือฟังก์ชันของตัวแปรสุ่มที่มี Domain (ในที่นี้คือค่าของ  $I_i$ ) ปรากฏอยู่ใน Real Line ขณะที่ Counter Domain (ในที่นี้คือ  $P_i$  หรือ  $p_i$ ) มีค่าปรากฏอยู่ช่วง  $[0, 1]$  ดังนั้นปัญหา  $P_i$  และ  $p_i$  ไม่ปรากฏค่าใน Unit Interval จึงหมดไป

<sup>2</sup> Zellner, A. and Lee, T.H., "Joint Estimation of Relationships Involving Discrete Random Variables", (Econometrica, 33, April, 1965) p. 382-394

<sup>3</sup> Probit  $P_i = I_i + 5$  และขอให้สังเกตว่า  $I_i$  คือ Quantile of Order  $P_i$  ของ Normal CDF

$$F^{-1}(P_i + u_i) = F^{-1}(P_i) + u_i \frac{d}{dP_i} F^{-1}(P_i) + R_i \quad \dots (4)$$

นั่นคือ  $F^{-1}(P_i) = F^{-1}(P_i) + u_i \frac{d}{dP_i} F^{-1}(P_i) + R_i$

$$I_i^0 = I_i + u_i \frac{d}{dP_i} F^{-1}(P_i) + R_i \quad \dots (5)$$

โดยที่  $I_i^0 = F^{-1}(P_i)$  เรียกว่า Observed Normit และ  $I_i = F^{-1}(P_i)$  เรียกว่า True Normit แต่  $R_i$  ( $R_i$  คือ Remainder ของ Taylor's Series) มีค่าต่ำและสามารถตัดทิ้งไปได้ ขณะเดียวกันเราสามารถพิสูจน์ได้ว่า  $u_i \frac{d}{dP_i} F^{-1}(P_i) = \frac{u_i}{Z(P_i)}$  เมื่อ  $Z(P_i)$  คือค่าของ  $f_Z(P_i)$  กล่าวคือ

$$Z(P_i) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2} P_i^2\right\}$$

ดังนั้นสมการที่ 5 จึงเปลี่ยนรูปเป็น

$$I_i^0 = I_i + \frac{u_i}{Z(P_i)} \quad ; \quad i = 1, 2, \dots, M$$

หรือสมการถดถอยที่แปลงรูปแล้วก็คือ

$$I_i^0 = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \dots + \beta_k X_{ki} + \frac{u_i}{Z(P_i)} \quad ; \quad i = 1, 2, \dots, M \quad \dots (6)$$

พิจารณาสมการที่ 6 จะพบว่า

$$E\left(\frac{u_i}{Z(P_i)}\right) = \frac{E(u_i)}{Z(P_i)} = 0 \quad , \quad i = 1, 2, \dots, M$$

$$V\left(\frac{u_i}{Z(P_i)}\right) = \frac{1}{[Z(P_i)]^2} V(u_i) = \frac{1}{[Z(P_i)]^2} \frac{P_i^2 \sigma_i^2}{n_i} \quad ; \quad i = 1, 2, \dots, M$$

$$V(U) = \phi = \begin{bmatrix} p_1 q_1 / n_1 [Z(p_1)]^2 & & & \\ & p_2 q_2 / n_2 [Z(p_2)]^2 & & \\ & & \ddots & \\ & & & p_M q_M / n_M [Z(p_M)]^2 \end{bmatrix}_{M \times M}$$

$$= \sigma^2 I_M$$

และเนื่องจาก  $p_i$  เป็นค่าประมาณของ  $P_i$  ดังนั้น

$$\hat{\phi} = \begin{bmatrix} p_1 q_1 / n_1 [Z(p_1)]^2 & & & \\ & p_2 q_2 / n_2 [Z(p_2)]^2 & & \\ & & \ddots & \\ & & & p_M q_M / n_M [Z(p_M)]^2 \end{bmatrix}_{M \times M}$$

ดังนั้นโดยอาศัย EGLS เราสามารถประมาณค่า  $\beta$  ได้ดังนี้

$$\hat{\beta} = (X' \hat{\phi}^{-1} X)^{-1} X' \phi^{-1} I^0, \quad \hat{V}(\hat{\beta}) \cong (X' \hat{\phi}^{-1} X)^{-1}$$

โดยที่  $I^0 = [F^{-1}(p_1), F^{-1}(p_2), \dots, F^{-1}(p_M)]'$  คือเวกเตอร์ของ Observed Normit และ

$$\hat{\beta} = [\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_M]'$$

### 9.3.2.2 Logit Analysis

Logit Analysis เป็นเทคนิคการแปลงค่า  $p_i$  สู่ดัชนี  $I_i$  โดยผ่าน Logistic Distribution คือ กำหนดให้

---

<sup>1</sup>อ่าน Zellner, A. and Lee, T.H., Ibid.. p.384

$$P_i = \frac{1}{1 + \exp[-I_i]} = \frac{1}{1 + \exp[-X_i' \beta]} = \frac{1}{1 + \exp[-(\beta_1 + \beta_2 X_{2i} + \dots + \beta_k X_{ki})]};$$

$i = 1, 2, \dots, M \quad \dots (1)$

การพัฒนาวิธีประมาณค่า  $\beta$  ดำเนินการดังนี้

จากสมการ  $p_i = P_i + u_i$  และ  $1-p_i = 1 - (P_i + u_i)$  จะพบว่า

$$\frac{P_i}{1-P_i} = \frac{P_i + u_i}{(1-P_i) - u_i}$$

$$= \frac{P_i (1 + u_i/P_i)}{1 - P_i (1 - u_i/Q_i)} \text{ เมื่อ } Q_i = 1 - P_i$$

ดังนั้น  $\ln\left(\frac{P_i}{1-P_i}\right) = \ln\left(\frac{P_i}{1-P_i}\right) + \ln\left(1 + \frac{u_i}{P_i}\right) - \ln\left(1 - \frac{u_i}{Q_i}\right); i = 1, 2, \dots, M$

<sup>1</sup>  $P_i$  คือ Probability ที่จะได้ Setting ที่  $i$  ของ X's และ  $\frac{P_i}{1-P_i}$  เรียกว่า Observed Odd ของ Setting ที่  $i$  ของ X's หรือ Observed Logit ขณะที่  $\frac{P_i}{1-p_i}$  เรียกว่า True Odd หรือ True Logit

สำหรับ Logistic Distribution มี CDF และคุณลักษณะต่าง ๆ ดังนี้

1.  $F(x) = [1 + e^{-(x-\alpha)/\beta}]^{-1}; -\infty < \alpha < \infty, \beta > 0$  ขอให้สังเกต Parameter Space

ของ  $\alpha$  และ  $\beta$  จะเห็นว่า มีลักษณะเดียวกับ  $\mu$  และ  $\sigma^2$  ของ  $N(\mu, \sigma^2)$

2.  $E(X) = \alpha$

3.  $V(X) = \frac{\beta^2 \pi^2}{3}$

4.  $M_x(t) = e^{\alpha t} \pi \beta t \operatorname{Cosec}(\pi \beta t)$

ถ้าให้  $\alpha = 0, \beta = 1$  จะพบว่า

$$F(x) = [1 + e^{-x}]^{-1}, E(X) = 0, V(X) = 3.29, M_x(t) = \pi t \operatorname{Cosec}(\pi t)$$

ขอให้สังเกตว่าในกรณีนี้ Logistic Distribution จะมีธรรมชาติใกล้เคียง  $N(0, 1)$  มาก

โดยอาศัย Taylor's Series กระจายรอบ  $a = 0$  เราพบว่า

$$\ln(1+x) = x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \dots$$

$$\ln(1-x) = -x - \frac{x^2}{2} - \frac{x^3}{3} - \frac{x^4}{4} + \dots$$

$$\ln\left(\frac{p_i}{1-p_i}\right) = \ln\left(\frac{p_i}{1-p_i}\right) + \frac{u_i}{p_i} + R_{1i} - \left(-\frac{u_i}{Q_i} + R_{2i}\right); i = 1, 2, \dots, M$$

เมื่อ  $R_{1i}$  และ  $R_{2i}$  คือ Remainder ของ  $\ln\left(1 + \frac{u_i}{P_i}\right)$  และ  $\ln\left(1 - \frac{u_i}{Q_i}\right)$  ตามลำดับ

นั่นคือ

$$\begin{aligned} \ln\left(\frac{p_i}{1-p_i}\right) &= \ln\left(\frac{p_i}{1-p_i}\right) + \frac{u_i}{p_i} + \frac{u_i}{Q_i} + R_i; i = 1, 2, \dots, M \\ &= \ln\left(\frac{p_i}{1-p_i}\right) + \frac{u_i(p_i + Q_i)}{p_i Q_i} + R_i; i = 1, 2, \dots, M \end{aligned}$$

$$\text{นั่นคือ } \ln\left(\frac{p_i}{1-p_i}\right) = \ln\left(\frac{p_i}{1-p_i}\right) + \frac{u_i}{p_i Q_i} + R_i; i = 1, 2, \dots, M \dots (2)$$

จากสมการที่ 1 เราทราบว่า

$$p_i = \frac{1}{1 + \exp[-(\beta_1 + \beta_2 X_{2i} + \dots + \beta_k X_{ki})]}$$

และ

$$1 - p_i = \frac{\exp[-(\beta_1 + \beta_2 X_{2i} + \dots + \beta_k X_{ki})]}{1 + \exp[-(\beta_1 + \beta_2 X_{2i} + \dots + \beta_k X_{ki})]}$$

$$\text{ดังนั้น } \frac{p_i}{1-p_i} = \frac{1}{\exp[-(\beta_1 + \beta_2 X_{2i} + \dots + \beta_k X_{ki})]}; i = 1, 2, \dots, M$$

$$\text{และ } \ln\left(\frac{p_i}{1-p_i}\right) = \beta_1 + \beta_2 X_{2i} + \dots + \beta_k X_{ki} \dots (3)$$

แทนที่สมการ 3 ในสมการที่ 2 จะพบว่าสมการถดถอยที่แปลงรูปแล้วคือ

$$\ln\left(\frac{p_i}{1-p_i}\right) = \beta_1 + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + \frac{u_i}{p_i Q_i} + R_i; i = 1, 2, \dots, M \dots (4)$$

พิจารณาสมการที่ 4 จะพบว่า

$$E\left(\frac{u_i}{P_i Q_i}\right) = 0 \text{ และ } V\left(\frac{u_i}{P_i Q_i}\right) = \frac{1}{P_i^2 Q_i^2} \cdot \frac{P_i Q_i}{n_i} = \frac{1}{n_i P_i Q_i} \text{ ดังนั้น}$$

$$V(U) = \phi = \begin{bmatrix} 1/n_1 P_1 Q_1 & & & \\ & 1/n_2 P_2 Q_2 & & \\ & & \ddots & \\ & & & 1/n_M P_M Q_M \end{bmatrix}_{M \times M}$$

$$\text{และ } \hat{V}(U) = \hat{\phi} = \begin{bmatrix} 1/n_1 P_1 Q_1 & & & \\ & 1/n_2 P_2 Q_2 & & \\ & & \ddots & \\ & & & 1/n_M P_M Q_M \end{bmatrix}_{M \times M}$$

ดังนั้นโดยอาศัย EGLS เราสามารถประมาณค่าเวกเตอร์  $\beta$  ได้ดังนี้คือ

$$\hat{\beta} = (X' \hat{\phi}^{-1} X)^{-1} X' \hat{\phi}^{-1} I^0, \quad \hat{V}(\hat{\beta}) \cong (X' \hat{\phi}^{-1} X)^{-1}$$

โดยที่  $I^0$  คือเวกเตอร์ของ Observed Logit กล่าวคือ  $I^0 = [\ln \frac{P_1}{1-P_1}, \ln \frac{P_2}{1-P_2}, \dots, \ln \frac{P_M}{1-P_M}]'$

สมการประมาณค่าของทั้ง Normit Model และ Logit Model คือ

$$I_i = \hat{\beta}_1 + \hat{\beta}_2 X_{2i} + \dots + \hat{\beta}_k X_{ki}$$

แต่สิ่งที่เราต้องการพยากรณ์คือ  $p_i$  มิใช่  $I_i$  ดังนั้น เมื่อแทนค่าในขนาดของ X's ลงในสมการประมาณค่าซึ่งจะทำให้ได้ค่า  $\hat{I}_i$  ให้แปลงค่า  $\hat{I}_i$  สู่  $\hat{p}_i$  ดังนี้

$$\hat{p}_i = F(I_i)$$

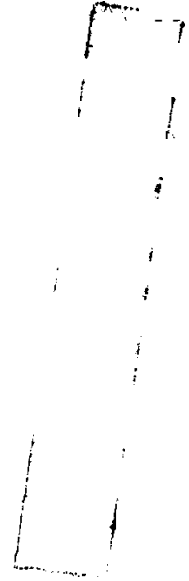
นั่นคือ

$$\hat{p}_i = \begin{cases} F(\hat{I}_i) = \Pr(Z < \hat{I}_i) & \text{สำหรับ Normit Model} \\ F(\hat{I}_i) = 1/[1+e^{-\hat{I}_i}] & \text{สำหรับ Logit Model} \end{cases}$$

จากนั้นให้ย้อนค่า  $\hat{p}_i$  สู  $\hat{Y}_i$  ดังนี้<sup>1</sup>

$$\text{Normit Model : } \hat{Y}_i = \begin{cases} 1 & \text{ถ้า } \hat{p}_i = \Pr(Z \leq \hat{I}_i) \\ 0 & \text{ถ้า } \hat{p}_i = \Pr(Z > \hat{I}_i) \end{cases}$$

$$\text{Logit Model : } \hat{Y}_i = \begin{cases} 1 & \text{ถ้า } \hat{p}_i = 1/[1+e^{-\hat{I}_i}] \\ 0 & \text{ถ้า } \hat{p}_i = 1 - 1/[1+e^{-\hat{I}_i}] \end{cases}$$



---

<sup>1</sup> Zellner, A and Lee, T.H., Ibid., p.384