

บทที่ 4

การวิเคราะห์การถดถอยโลจิสติกสำหรับข้อมูลทวิ

4.1 กล่าวนำ

ในบทที่ 3 ได้อธิบายถึงการประยุกต์ใช้การถดถอยเชิงเส้น กับตัวแปรตามที่มีข้อมูลเป็นข้อมูลทวิแล้ว พบว่าขึ้นไม่สามารถแก้ปัญหาเกี่ยวกับข้อจำกัดที่ว่า $0 \leq E(Y_i) = P_i \leq 1$ ได้ในทันที จะได้กล่าวถึงการพิจารณาโมเดลสำหรับ $E(Y_i) = P_i$ ที่เหมาะสม เพื่อให้ได้โมเดลการถดถอยที่มีคุณสมบัติเป็นไปตามข้อจำกัดดังกล่าว

โมเดลการถดถอยสำหรับข้อมูลทวิ ที่จะให้ค่าจากโมเดลมีค่าอยู่ระหว่าง 0 และ 1 นั้น พนว่า ควรจะเป็นโมเดลเส้นตรงที่มีค่าลู่เข้าสู่ 0 และ 1 โมเดลที่มีลักษณะดังกล่าวนี้ ได้แก่

1. โมเดลโลจิสติก (logistic model)

$$P_i = [1 + \exp(-U)]^{-1}$$

ซึ่งคือ ผังกั้นความน่าจะเป็นสะสม (cumulative probability density function) ของการแจกแจงแบบโลจิสติก (logistic distribution) ที่มีผังกั้นการแจกแจงความน่าจะเป็น (probability density function ; p.d.f.) คือ

$$f(U) = e^{-U} / [1 + e^{-U}]^2 \quad ; -\infty < U < \infty$$

เมื่อ $U = \sum_{j=0}^{q-1} \beta_j X_j \quad ; j = 0, 1, 2, \dots, q-1$
 $X_0 = 1$

และสามารถแปลงโมเดลโลจิสติก ได้เป็น $\text{logit}(P_i) = \ln(P_i/Q_i) =$

$$U = \sum_{j=0}^{q-1} \beta_j X_j \quad \text{เมื่อ } Q_i = 1 - P_i$$

2. โมเดลโลรบิก(probit model)

$$P_i = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^u \exp(-1/2z^2) dz ; -\infty < z < \infty$$

ซึ่งเป็นฟังก์ชันการแจกแจงความน่าจะเป็นสะสม ของการแจกแจงปกติมาตรฐาน
ที่มี p.d.f. เป็น

$$f(z) = \frac{1}{\sqrt{2\pi}} \exp(-1/2z^2)$$

การใช้โมเดลโลจิสติกและโมเดลโลรบิก เป็นโมเดลการถดถอยลำดับข้อมูลทวิ จะให้ผลคล้ายคลึงกัน แต่ความนิยมในโมเดลโลจิสติก มีมากกว่าโมเดลโลรบิก เนื่องจากสามารถนำไปประยุกต์ใช้ในงานทดลองวิจัยได้มากกว่า โดยเฉพาะอย่างยิ่งงานทดลองวิจัยเกี่ยวกับสิ่งมีชีวิต (bioassay) ดังนั้นในที่นี้ จึงจะกล่าวถึงเฉพาะการวิเคราะห์การถดถอยลำดับข้อมูลทวิ ที่ใช้โมเดลโลจิสติกเป็นโมเดลการถดถอย ซึ่งเรียกว่าการวิเคราะห์การถดถอยโลจิสติก (logistic regression analysis) โดยมีข้อกำหนดว่า ข้อมูลทวิของตัวแปรตามที่ใช้ต้องมาจากการแจกแจงแบบทวินามเท่านั้น

4.2 โมเดลโลจิสติกเชิงเส้น(The linear logistic model)

ถ้า P_i คือความน่าจะเป็นที่จะเกิดผลสำเร็จ ของการแจกแจงแบบทวินาม ขนาด n_i
ซึ่งมีความลัมพันธ์กับตัวแปรอิสระ $X_{j,i}$ และงดด้วยโมเดลโลจิสติก เป็น

$$P_i = [1 + e^{-u_i}]^{-1}$$

$$= e^{u_i} [1 + e^{u_i}]^{-1}$$

$$\text{เมื่อ } U_i = \sum_{j=0}^{q-1} \beta_j X_{j,i}; \quad j = 0, 1, \dots, q-1 \\ X_{0,i} = 1$$

และมี Q_i คือ ความน่าจะเป็นที่จะไม่เกิดผลสำเร็จ

$$\begin{aligned} \text{โดย } Q_i &= 1 - P_i \\ &= e^{-U_i} [1 + e^{-U_i}]^{-1} \\ &= [1 + e^{U_i}]^{-1} \end{aligned}$$

นิยาม odds of success หรือ Odds คืออัตราส่วนระหว่างความน่าจะเป็นที่จะเกิดผลสำเร็จ ต่อความน่าจะเป็นที่จะไม่เกิดผลสำเร็จ ; P_i / Q_i จะได้

$$\text{Odds} = P_i / Q_i = e^{U_i}$$

$$\text{และ } \ln \text{Odds} = \text{logit}(P_i) = U_i = \sum_{j=0}^{q-1} \beta_j X_{j,i} \quad (4.2.1)$$

เรียก $\text{Logit}(P_i)$ ว่าเป็น link function ซึ่งแสดงว่าการแปลงโมเดลโลจิสติกของ P ให้อยู่ในรูปของ $\ln \text{Odds}$ ได้เป็นสมการเส้นตรง และเรียกสมการ (4.2.1) ว่าเป็นโมเดลโลจิสติกเชิงเส้น สำหรับตัวแปรอิสระ หรือโมเดลโลจิต (logit model)

4.3 การประมาณค่าพารามิเตอร์

ในการประมาณค่าพารามิเตอร์ β_j สำหรับโมเดลโลจิสติกเชิงเส้น $\text{logit}(P_i) = \sum_{j=0}^{q-1} \beta_j X_{j,i}$ จะใช้วิธีประมาณค่าแบบภาวะน่าจะเป็นสูงสุด ซึ่งมีพังก์ชันภาวะน่าจะเป็นสูงสุด คือ

$$L(\beta) = \prod_{i=1}^n n_i C_{y_i} P_i^{y_i} Q^{n_i - y_i}$$

$$\ln L(\theta) = \sum_{i=1}^n [\ln n_i C_{y_i} + y_i \ln P_i + (n_i - y_i) \ln Q_i] \quad (4.3.1)$$

$$= \sum_{i=1}^n [\ln n_i C_{y_i} + y_i \ln P_i / Q_i + n_i \ln Q_i]$$

$$= \sum_{i=1}^n \{ \ln n_i C_{y_i} + y_i U_i + n_i \ln [1 + e^{U_i}]^{-1} \}$$

$$\frac{\partial \ln L(\theta)}{\partial \theta_j} = \sum_{i=1}^n y_i X_{j,i} - \sum_{i=1}^n n_i X_{j,i} e^{U_i} [1 + e^{U_i}]^{-1}$$

$$= \sum_{i=1}^n y_i X_{j,i} - \sum_{i=1}^n n_i X_{j,i} \hat{p}_i ; j = 0, 1, \dots, q-1 \quad (4.3.2)$$

$X_{0,1} = 1$

เมื่อกำสมการ (4.3.2) ให้เท่ากับ 0 จะให้สมการที่ไม่ใช้เลี้นตรงจำนวน q สมการ ที่มีค่าประมาณของ θ_j ซึ่งไม่ทราบค่าติดอยู่ในสมการ การหาค่าประมาณของ θ_j จาก สมการ q สมการนี้ จะใช้การประมาณค่าแบบภาวะน่าจะเป็นสูงสุดที่ทำซ้ำ (iterative maximum likelihood estimation โดย Newton-Raphson procedure และ Fisher's method of scoring และวิธีกำลังสองต่ำสุดถ่วงน้ำหนักที่ทำซ้ำ (iteratively weighted least squares ซึ่งขบวนการประมาณค่าพารามิเตอร์ตั้งกล่าวนี้จะต้องอาศัยคอมพิวเตอร์ช่วย ในการประมาณผล

4.3.1 การประมาณค่าแบบภาวะน่าจะเป็นสูงสุดที่ทำซ้ำ

1. โดยวิธี Newton-Raphson

ให้ $B(\theta)$ คือ คอลัมน์เวคเตอร์ขนาด $q \times 1$ ของสมการ (4.3.2) จำนวน q

สมการ โดยค่าในแคลวิล์ j คือ $\partial \ln L(\beta) / \partial \beta_j$ ซึ่งเรียกว่า coefficient score และให้ $H(\beta)$ คือแมทริกซ์ขนาด $q \times q$ ที่ได้จากการหาอนุพันธ์ครั้งที่สองของสมการ (4.3.1) เทียบกับค่า β_k โดยค่าในแคลวิล์ j คอลัมน์ที่ k คือ $\partial^2 \ln L(\beta) / \partial \beta_j \partial \beta_k$ สำหรับ $j, k = 0, 1, \dots, q-1$ แมทริกซ์ $H(\beta)$ บางครั้งเรียกว่า Hessian matrix

ให้ $U(\beta)$ คือ คอลัมน์เวคเตอร์ขนาด $q \times 1$ ของ coefficient scores ที่ได้จากการแทนค่าประมาณของพารามิเตอร์ β , $\hat{\beta}$ และโดย Taylor series กระจาย $U(\hat{\beta})$ รอบค่า β^* เมื่อ β^* คือค่าที่ใกล้เคียง $\hat{\beta}$ จะได้

$$U(\hat{\beta}) \approx U(\beta^*) + H(\beta^*) (\hat{\beta} - \beta^*) \quad (4.3.1.1)$$

โดยนิยามของการประมาณค่าพารามิเตอร์ β_j โดยวิธีภาวะน้ำจะเป็นสูงสุด จะต้องได้ว่า

$$\frac{\partial \ln L(\beta)}{\partial \beta_j} \Bigg|_{\hat{\beta}} = 0$$

ซึ่งจะได้ $U(\hat{\beta}) = 0$ ดังนั้น จะได้ (4.3.1.1) เป็น

$$0 \approx U(\beta^*) + H(\beta^*) (\hat{\beta} - \beta^*) \quad (4.3.1.2)$$

$$\hat{\beta} \approx \beta^* - [H(\beta^*)]^{-1} U(\beta^*)$$

จาก (4.3.1.2) เมื่อมีการทำซ้ำ ๆ กัน จนถึงรอบที่ $r+1$ จะได้ค่าประมาณของ β ในรอบที่ $r+1$ เป็น

$$\hat{\beta}_{r+1} = \hat{\beta}_r - [H(\hat{\beta}_r)]^{-1} U(\hat{\beta}_r) \quad (4.3.1.3)$$

เมื่อ $r = 0, 1, \dots$ และ $\hat{\beta}_0$ คือคอลัมน์เวคเตอร์ของค่าประมาณเริ่มต้นของพารามิเตอร์ β

2. โดย Fisher's method of scoring

ให้ $I(\theta)$ คือค่าคาดหวังของ $H(\theta)$ คูณด้วย -1 และเรียก $I(\theta)$ ว่า information matrix โดย $I(\theta)$ เป็นแมทริกซ์ขนาด $q \times q$ ซึ่งมีค่าในแถวที่ j คอลัมน์ที่ k เป็น $-E \left[\frac{\partial^2 \ln L(\theta)}{\partial \theta_j \partial \theta_k} \right]$

inverse ของ $I(\theta)$ คือ $[I(\theta)]^{-1}$ เรียกว่า asymptotic variance-covariance matrix ของค่าประมาณ โดยวิธีภาวะน่าจะเป็นสูงสุด

และโดยขั้นตอนการทำซ้ำ ๆ กัน เช่นเดียวกับวิธี Newton-Raphson จะได้ (4.3.1.3) เป็น

$$\hat{\theta}_{r+1} = \hat{\theta}_r + [I(\hat{\theta}_r)]^{-1} U(\hat{\theta}_r) \quad (4.3.1.4)$$

การประมาณค่าพารามิเตอร์ θ โดยวิธี Newton-Raphson และ Fisher's method of scoring จะให้ค่าประมาณลู่เข้าสู่ค่าประมาณแบบภาวะน่าจะเป็นสูงสุดของ θ ความคลาดเคลื่อนมาตรฐานของค่าประมาณ คือรากที่สองของค่าในแนววงแยงของ $-[H(\theta)]^{-1}$ หรือ $[I(\theta)]^{-1}$ ซึ่งค่าประมาณและความคลาดเคลื่อนมาตรฐานของค่าประมาณโดยวิธีทั้งสองนี้ จะให้ค่าเหมือนกัน เมื่อการวิเคราะห์การถดถอยนี้ไม่เดลเป็นโมเดลโลจิสติกเชิงเส้นซึ่งใช้กับข้อมูลที่มีการแจกแจงแบบทวินาม และค่าประมาณที่ได้นี้จะมีการแจกแจงใกล้เคียงการแจกแจงปกติ

4.3.2 การประมาณค่าตัวยั่งวิธีกำลังสองต่อสุ่ลต่อวั้นน้ำหนัก กระทำจาก

ไมเดลโลจิก ; $\text{logit}(P_i) = \ln P_i/Q_i = U_i = \sum_{j=0}^{q-1} \beta_j X_{ji}, \quad \text{ได้}$

$$\hat{\theta} = (X'WX)^{-1} X'WU$$

เพื่อให้ค่า $\hat{\beta}$ ลู่เข้าสู่ค่าคงที่ จะทำการประมาณค่าด้วยวิธีกำลังสองต่ำสุดถ่วงน้ำหนัก
ที่กำช้ำ ซึ่งต้องมีการปรับค่า P และ U ให้เป็น P^* และ U^* ตามลำดับ โดย McCullagh
and Nelder (1983) แสดงวิธีการปรับค่า ดังนี้

$$U_{ir}^* = U_{ir} + \frac{(y_i - n_i \hat{p}_{ir})}{n_i} \left[\frac{\partial \ln P_i / Q_i}{\partial P_i} \right]_r$$

$$\begin{aligned} \text{แต่ } \frac{\partial \ln P_i / Q_i}{\partial P_i} &= \frac{\partial \ln P_i}{\partial P_i} - \frac{\partial \ln (1 - P_i)}{\partial P_i} \\ &= \frac{1}{P_i} + \frac{1}{1 - P_i} = \frac{1}{\hat{p}_i(1 - \hat{p}_i)} = \frac{1}{\hat{p}_i \hat{q}_i} \end{aligned}$$

ดังนั้น เมื่อกำช้ำ r รอบจะได้

$$U_{ir}^* = X \hat{\theta}_r + \frac{(y_i - n_i \hat{p}_{ir})}{n_i \hat{p}_{ir} \hat{q}_{ir}}$$

$$\text{และ } W_{ir}^{*-1} = V(\hat{p}_{ir}) \left[\frac{\partial \ln P_i / Q_i}{\partial P_i} \right]_r^2$$

$$= \frac{\hat{p}_{ir} \hat{q}_{ir}}{n_i} \cdot \left[\frac{1}{\hat{p}_{ir} \hat{q}_{ir}} \right]^2$$

$$\begin{aligned} &= \frac{1}{n_i \hat{p}_{ir} \hat{q}_{ir}} \\ W_{ir}^* &= n_i \hat{p}_{ir} \hat{q}_{ir} \end{aligned}$$

ดังนั้น ค่าประมาณของ θ ในรอบที่ $r+1$ เป็น

$$\hat{\theta}_{r+1} = (X' W_r^* X)^{-1} X' W_r^* U_r^*$$

ซึ่งจะให้ค่าประมาณเป็นเช่นเดียวกับการประมาณค่าด้วยวิธีภาวะน่าจะเป็นสูงสุดที่ทำซ้ำโดย Fisher's method of scoring และสามารถแสดงให้เห็นได้ดังนี้

จากสมการ (4.3.1)

$$\begin{aligned} \ln L(\beta) &= \sum_{i=1}^n \{ \ln n_i C_{y_i} + y_i \ln P_i + (n_i - y_i) \ln Q_i \} \\ \frac{\partial \ln L(\beta)}{\partial \beta_j} &= \sum_{i=1}^n \left[\frac{\partial \ln L(\beta)}{\partial P_i} \cdot \frac{\partial P_i}{\partial U_i} \cdot \frac{\partial U_i}{\partial \beta_j} \right] \quad (4.3.2.1) \\ \frac{\partial \ln L(\beta)}{\partial P_i} &= \frac{y_i}{\hat{p}_i} - \frac{n_i - y_i}{1 - \hat{p}_i} = \frac{y_i - n_i \hat{p}_i}{\hat{p}_i \hat{q}_i} \\ \frac{\partial U_i}{\partial P_i} &= \frac{\partial \ln(P_i/Q_i)}{\partial P_i} = \frac{1}{\hat{p}_i \hat{q}_i} \\ \frac{\partial U_i}{\partial \beta_j} &= \frac{\partial (\sum_{j=0}^{q-1} \beta_j X_{j,i})}{\partial \beta_j} = X_{j,i} \end{aligned}$$

ดังนั้นจะได้ (4.3.2.1) เป็น

$$\begin{aligned} \frac{\partial \ln L(\beta)}{\partial \beta_j} &= \sum_{i=1}^n \left[\frac{(y_i - n_i \hat{p}_i)}{\hat{p}_i \hat{q}_i} \cdot \hat{p}_i \hat{q}_i \cdot X_{j,i} \right] \\ &= \sum_{i=1}^n (y_i - n_i \hat{p}_i) X_{j,i} \quad (4.3.2.2) \end{aligned}$$

$$\text{ให้ } y_i^* = \frac{y_i - n_i \hat{p}_i}{n_i \hat{p}_i \hat{q}_i}$$

และ $W_i^* = n_i \hat{p}_i \hat{q}_i$

จะได้ $\frac{\partial \ln L(\beta)}{\partial \beta_j} = \sum_{i=1}^n y_i^* W_i^* X_{j,i}$

ให้ X คือ แมทริกซ์ขนาด $n \times q$ ของตัวแปรอิสระ X จำนวน $q - 1$ ตัว

W^* คือ แมทริกซ์ขนาด $n \times n$ ที่มีค่าในแนวเทղะแยก คือ W_i^*

และ Y^* คือ แมทริกซ์ขนาด $n \times 1$ ที่มีค่าในแนวแก้วที่ i คือ y_i^*

ดังนั้นจะได้ $\frac{\partial \ln L(\beta)}{\partial \beta_j} = U(\beta) = X' W^* Y^*$

พิจารณา $- E[\frac{\partial^2 \ln L(\beta)}{\partial \beta_j \partial \beta_k}] = E[\frac{\partial \ln L(\beta)}{\partial \beta_j} \cdot \frac{\partial \ln L(\beta)}{\partial \beta_k}] = I(\beta)$

จาก (4.3.2.2)

$$\begin{aligned} \text{จะได้ } - E[\frac{\partial^2 \ln L(\beta)}{\partial \beta_j \partial \beta_k}] &= E[\sum_{i=1}^n (y_i - n_i \hat{p}_i)(y_i - n_i \hat{p}_i) X_{j,i} X_{k,i}] \\ &= \sum_{i=1}^n [E(y_i - n_i \hat{p}_i)(y_i - n_i \hat{p}_i)] X_{j,i} X_{k,i} \\ &= \sum_{i=1}^n [\text{Cov}(Y_i, Y_i)] X_{j,i} X_{k,i} \\ &= \sum_{i=1}^n [\text{V}(Y_i)] X_{j,i} X_{k,i} ; \quad y_i, y_i \text{ เป็นอิสระกัน} \\ &= \sum_{i=1}^n n_i \hat{p}_i \hat{q}_i X_{j,i} X_{k,i} \end{aligned}$$

$$\begin{aligned}
 &= \sum_{i=1}^n W_i^* X_{j,i} X_{k,i} \\
 &= X' W_r^* X = I(\theta)
 \end{aligned}$$

จาก (4.3.1.4) จะได้ $\hat{\theta}_{r+1} = \hat{\theta}_r + (X' W_r^* X)^{-1} X' W_r^* Y_r^*$

$$= (X' W_r^* X)^{-1} [X' W_r^* (X\hat{\theta}_r + Y_r^*)]$$

เมื่อให้ $U_r^* = X\hat{\theta}_r + Y_r^*$

จะได้ $\hat{\theta}_{r+1} = (X' W_r^* X)^{-1} X' W_r^* U_r^*$

4.3.3 ขั้นตอนของขบวนการประมาณค่าด้วยวิธีกำลังสองต่อสุ่ลต่อวงน้ำหนักที่ทำซ้ำ

Collett (1991) ได้แสดงขั้นตอนของขบวนการประมาณค่าด้วยวิธีกำลังสองต่อสุ่ลต่อวงน้ำหนักที่ทำซ้ำ ดังนี้

$$1. \text{ ให้ค่าเริ่มต้นของ } \hat{p}_1 \text{ กับ } \hat{p}_{10} = \frac{y_1 + 0.5}{n_1 + 1}$$

คำนวณ $W_{10}^* = n_1 \hat{p}_{10} \hat{q}_{10}$

และ $U_{10}^* = \text{logit}(\hat{p}_{10}) = \ln \hat{p}_{10}/\hat{q}_{10}$

$$2. \text{ คำนวณค่า } \hat{\theta}_1 \text{ จาก } \hat{\theta}_1 = (X' W_o^* X)^{-1} X' W_o^* U_o^*$$

3. นำค่า $\hat{\theta}_1$ จาก 2 เนื่อประมาณค่า $\hat{p}_{11} = [1 + \exp(-\sum_{j=0}^{q-1} \hat{\theta}_{j1} X_{j1})]^{-1}$

$$\text{คำนวณค่า } W_{11}^* = n_1 \hat{p}_{11} \hat{q}_{11}$$

$$\text{และ } U_{11}^* = \sum_{j=0}^{q-1} \hat{\theta}_{j1} X_{j1} + \frac{(y_1 - n_1 \hat{p}_{11})}{n_1 \hat{p}_{11} \hat{q}_{11}}$$

4. คำนวณค่า $\hat{\theta}_2$ จาก $\hat{\theta}_2 = (X'W_1^*X)^{-1}X'W_1^*U_1^*$

5. ดำเนินการตามขั้นตอน 3-4 ข้อ ๆ จนได้ค่า $\hat{\theta}$ ล้วนเข้าสู่ค่าคงที่ค่าหนึ่ง ซึ่งจะได้เป็น $\hat{\theta}$ ตามที่ต้องการ

4.4 การทดสอบภาวะสารูปดีของโมเดลโลจิสติกเชิงเส้น

เมื่อกำกับประมาณค่าพารามิเตอร์ $\hat{\theta}$ จนได้โมเดลโลจิสติกเชิงเส้นตามที่ต้องการแล้ว ควรจะทำการทดสอบดูว่า โมเดลที่ได้นี้มีภาวะสารูปดีกับค่าลังเกต (test of goodness of fit) หรือไม่ ซึ่งถ้าพบว่า ค่าที่ประมาณได้จากโมเดลมีภาวะสารูปดีกับค่าลังเกต ก็แสดงว่า โมเดลนั้นใช้ได้ แต่ถ้าไม่เป็นตั้ง เช่นที่กล่าว ก็จำเป็นที่จะต้องหาโมเดลที่เหมาะสมต่อไป

4.4.1 การทดสอบแบบอัตราลั่นภาวะน่าจะเป็น

การทดสอบแบบอัตราลั่นภาวะน่าจะเป็น (likelihood ratio test) เป็นการเปรียบเทียบระหว่างฟังก์ชันภาวะน่าจะเป็นสูงสุด 2 ฟังก์ชัน คือ

1. ฟังก์ชันภาวะน่าจะเป็นสูงสุดของโมเดล ที่ใช้ค่าพารามิเตอร์ซึ่งได้มาจากการประมาณค่า โดยวิธีภาวะน่าจะเป็นสูงสุด เรียกว่าฟังก์ชันภาวะน่าจะเป็นสูงสุดของโมเดลที่ประมาณได้ โดยเรียกโมเดลที่ประมาณได้ว่า current model ให้ฟังก์ชันภาวะน่าจะเป็นของโมเดลนี้ แทนด้วย L_c

2. ฟังก์ชันภาวะน่าจะเป็นสูงสุดของ โมเดล ที่พอดีกับค่าสังเกต หรือ full model ให้ฟังก์ชันภาวะน่าจะเป็นสูงสุดของ โมเดลนี้ แทนด้วย L_f

และให้

$$\begin{aligned} D &= -2 \ln [L_c / L_f] \\ &= -2 [\ln L_c - \ln L_f] \end{aligned}$$

จากสมการ (4.3.1) จะได้ฟังก์ชันภาวะน่าจะเป็นสูงสุดของ current model เป็น

$$\ln L_c = \sum_{i=1}^n [\ln n_i C_{y_i} + y_i \ln \hat{p}_i + (n_i - y_i) \ln \hat{q}_i]$$

และฟังก์ชันภาวะน่าจะเป็นสูงสุดของ full model เป็น

$$\ln L_f = \sum_{i=1}^n [\ln n_i C_{y_i} + y_i \ln \tilde{p}_i + (n_i - y_i) \ln \tilde{q}_i]$$

เมื่อ $\tilde{p}_i = y_i/n_i$

ดังนั้น $D = 2 \sum_{i=1}^n [y_i \ln (\tilde{p}_i/\hat{p}_i) + (n_i - y_i) \ln (\tilde{q}_i/\hat{q}_i)]$

$$= 2 \sum_{i=1}^n [y_i \ln \frac{y_i}{\hat{y}_i} + (n_i - y_i) \ln \frac{n_i - y_i}{n_i - \hat{y}_i}]$$

เมื่อ $\hat{Y}_i = n_i \hat{p}_i$ และ $\hat{p}_i = [1 + \exp(-\sum_{j=0}^{q-1} b_j X_{j,i})]^{-1}$

เรียก D ว่า deviance ซึ่งจะเป็นตัวสถิติที่ใช้ในการทดสอบ (test statistic)

เมื่อ D มีค่ามาก แสดงว่า โมเดลที่ประมาณได้จะไม่พอดีกับ โมเดลของค่าสังเกต แต่ถ้า โมเดลทั้งสองมีความใกล้เคียงกัน ก็จะได้ D มีค่าน้อย โดยทฤษฎีจะได้ว่า D มีการแจก

แจงลู่เข้าสู่การแจกแจงแบบไคสแควร์ (chi-square distribution) ด้วย $d.f. = n - q$ ดังนั้น เมื่อพบว่า D มีค่ามากกว่าค่า χ^2 จากตาราง แสดงว่ามีความแตกต่างระหว่างโมเดลประมาณ กับโมเดลของค่าสังเกตอย่างมีนัยสำคัญทางสถิติ ซึ่งหมายความว่าโมเดลที่ประมาณได้ ยังไม่พอดีกับโมเดลของค่าสังเกต จะต้องมีการพิจารณาหาโมเดลใหม่ที่เหมาะสมต่อไป

ในการทดสอบภาวะสรุปดี โดยการทดสอบแบบอัตราส่วนภาวะน่าจะเป็นสูงสุดนี้ Collett (1991) พิจารณาว่า ถ้า deviance ที่คำนวณได้ มีค่าใกล้เคียงกับ $d.f. = n - q$ สามารถสรุปได้ว่า โมเดลที่ประมาณได้มีความพอดีกับค่าสังเกต ทั้งนี้เนื่องจากข้อสมมติที่ว่า ข้อมูลที่มาจาก การแจกแจงแบบทวินาม จะให้ค่าความแปรปรวนของความคลาดเคลื่อน ซึ่งในที่นี้คือ deviance หากด้วย $d.f. = n - q$ มีค่าเท่ากับ 1^1

4.4.2 Pearson's X^2 -statistic

การทดสอบภาวะสรุปดี ระหว่างโมเดลประมาณกับโมเดลของค่าสังเกต สามารถทำได้อีกวิธีหนึ่ง คือ ใช้ Pearson's X^2 -statistic ซึ่งนิยามค่าไว้ ดังนี้

$$X^2 = \sum_{i=1}^n \frac{(y_i - n_i \hat{p}_i)^2}{n_i \hat{p}_i \hat{q}_i}$$

โดย X^2 จะมีการแจกแจงลู่เข้าสู่การแจกแจงแบบไคสแควร์ด้วย $d.f. = n-q$

เช่นเดียวกับการแจกแจงของ deviance

แม้ว่าการทดสอบภาวะสรุปดี ระหว่างโมเดลประมาณกับโมเดลของค่าสังเกต จะสามารถใช้ได้ทั้งค่า deviance และ Pearson's X^2 -statistic เป็นตัวทดสอบสถิติ

¹ McCullagh and Nelder. Generalized liner models (1983) pp.83

ก็ตาม แต่โดยทั่วไปแล้ว การใช้ deviance จะมีความเหมาะสมกว่าการใช้ Pearson's χ^2 -statistic เพราะสามารถนำไปประยุกต์ใช้ในการเปรียบเทียบระหว่างโมเดลประมาณต่าง ๆ ได้ด้วย ซึ่งจะได้กล่าวต่อไป

4.5 การเปรียบเทียบโมเดลโลจิสติกเชิงเส้น(Comparing linear logistic model)

การวิเคราะห์การทดสอบโดยโลจิสติกโดยใช้โมเดลโลจิสติกเชิงเส้น เพื่อศึกษาถึงอิทธิพลของตัวแปรอิสระ จำนวนหลาย ๆ ตัวนั้น มีปัญหาเกิดขึ้นว่า ตัวแปรอิสระตัวใดบ้างที่จะมีอิทธิพลและตัวแปรอิสระใดบ้างที่ไม่มีอิทธิพล หรือกล่าวอีกนัยหนึ่งก็คือ โมเดลที่จะถูกนำมาใช้นั้นควรจะให้มีตัวแปรอิสระใดบ้างที่ปรากฏอยู่ในโมเดล การพิจารณาเกี่ยวกับปัญหานี้ ก็คือการทดสอบสมมติฐานเกี่ยวกับ β_j ; $j = 0, 1, \dots, q-1$ นั้นเอง ทำได้โดยการเปรียบเทียบโมเดลโลจิสติกเชิงเส้นที่เป็นร่างแท (nested) กัน ทีละคู่ แล้วพิจารณาถึงผลกระทบของการนำตัวแปรเข้าหรือออกจากโมเดลว่า ทำให้ค่า deviance เปลี่ยนไปอย่างไร มีนัยสำคัญหรือไม่ วิธีการนี้เรียกว่าการวิเคราะห์ค่า deviance (analysis of deviance)

พิจารณาโมเดลโลจิสติกเชิงเส้น ที่เป็นร่างแทกัน 2 โมเดล ต่อไปนี้ คือ

$$\text{โมเดล (1)} : \text{logit}(P_i) = \beta_0 + \beta_1 X_1 + \dots + \beta_h X_h$$

$$\text{โมเดล (2)} : \text{logit}(P_i) = \beta_0 + \beta_1 X_1 + \dots + \beta_h X_h + \beta_{h+1} X_{h+1} + \dots + \beta_k X_k$$

ให้ D_1 คือ deviance ของโมเดล (1) ซึ่งมี d.f. = $n-(h+1)$

และ D_2 คือ deviance ของโมเดล (2) ซึ่งมี d.f. = $n-(k+1)$

โมเดล (1) เป็นร่างแทของโมเดล (2) ซึ่งแน่นอนว่าค่า D_2 ย่อมมีค่าน้อยกว่า D_1 ผลต่างระหว่าง deviance ของ 2 โมเดล เป็น $D_1 - D_2$ คือการเปลี่ยนแปลงของ deviance ซึ่งเนื่องมาจากการนำตัวแปร $X_{h+1}, X_{h+2}, \dots, X_k$ ใส่ในโมเดล ภายหลังจากในโมเดลมีตัวแปร X_1, X_2, \dots, X_h อญญาแล้ว

เนื่องจากได้กล่าวแล้วว่า deviance ของโมเดล จะมีการแจกแจงลุ้นเข้าสู่การแจกแจงแบบไคสแควร์ ดังนี้จะได้ว่า การแจกแจงของ $D_1 - D_2$ ก็จะลุ้นเข้าสู่การแจกแจงแบบไคสแควร์ด้วย ดังนี้

$$D_1 = -2[\ln L_{c_1} - \ln L_f] \sim \chi^2 ; \text{ d.f.} = n-(h+1)$$

$$D_2 = -2[\ln L_{c_2} - \ln L_f] \sim \chi^2 ; \text{ d.f.} = n-(k+1)$$

$$D_1 - D_2 = 2[\ln L_{c_2} - \ln L_{c_1}] \sim \chi^2 ; \text{ d.f.} = k-h$$

ดังนั้น ในการพิจารณาถึงผลกรากทบทวนตัวแปรอิสระที่จะนำเข้าในโมเดล หรือนำออกจากโมเดล หรือการทดสอบ $H_0 : \beta_{h+1} = \beta_{h+2} = \dots = \beta_k = 0$ โดยมี H_1 : มีค่า β_j ; $j = h+1, h+2, \dots, k$ อ่างน้อย 1 ค่าที่ไม่เท่ากับ 0 จะพิจารณาจากค่า $D_1 - D_2$ ว่า มีนัยสำคัญหรือไม่ ซึ่งถ้าพบว่าไม่มีนัยสำคัญ คือ D_1 กับ D_2 มีค่าไม่แตกต่างกันทางสถิติ โมเดล (1) จะถูกนำไปใช้ ซึ่งหมายความว่า $\beta_{h+1} = \beta_{h+2} = \dots = \beta_k = 0$ หรือตัวแปร $X_{h+1}, X_{h+2}, \dots, X_k$ จะไม่มีอิทธิพลต่อโมเดล แต่ถ้าพบว่า D_1 กับ D_2 แตกต่างกันอย่างมีนัยสำคัญ ก็แสดงว่า มีค่า β_j ; $j = h+1, h+2, \dots, k$ อ่างน้อย 1 ค่า ที่ไม่เท่ากับ 0 จะต้องมีการพิจารณาต่อไปว่า β_j ตัวใดมีค่าไม่เท่ากับ 0 ซึ่งถ้าพบว่า β_j ทุกค่า มีค่าไม่เท่ากับ 0 หมด โมเดล (2) จะถูกนำไปใช้

ในทางปฏิบัติ จะสร้าง all possible models ซึ่งอาจรวมเทอมต่าง ๆ ของตัวแปรอิสระที่มีการเปลี่ยนข้อมูลแล้ว และเทอมของอิทธิพลร่วมระหว่างตัวแปรอิสระที่เป็นกลุ่มด้วย ขั้นมาเปรียบเทียบกัน และทำการวิเคราะห์ค่า deviance เพื่อเลือกโมเดลที่เหมาะสม พิจารณาได้จากตัวอย่างต่อไปนี้ คือ

จากการทดลองของ Hoblyn และ Palmer (1934) ซึ่งทดลองติด rak ของตันพิมพ์ผันธุ์ Common Mussel ในช่วงเวลาระหว่างเดือนตุลาคม 2477 และเดือนกุมภาพันธ์ 2478 ให้มีความยาว 2 ขนาด คือ ขนาด 12 และ 6 เซนติเมตร จำนวนขนาดละ 480 ล้วน และนำ

ครั้งหนึ่งของแต่ละขนาดไปปลูกกันที่ สำหรับอีกรั้งหนึ่งนำไปเพาะในทรายก้อนจนถึงคุณภาพในไม้ผลิแล้วจึงนำไปปลูก ทดลองจนถึงเดือนตุลาคม 2478 จังหวัดจำนวนต้นผลมีชีวิตครอบครอง ได้ข้อมูลดังตาราง 4.1

ตาราง 4.1 อัตราการอ่อนรู้สึกของต้นผลมจากการขยายพันธุ์โดยการตัดราก

ความยาวของรากที่ตัด	เวลาที่ปลูก	จำนวนที่อยู่รอด จากหั้งหมด 240 ส่วน	สัดส่วนที่อยู่รอด
6 เซนติเมตร	ปลูกกันที่ ปลูกในคุณภาพในไม้ผลิ	107	0.45
	ปลูกกันที่ ปลูกในคุณภาพในไม้ผลิ	31	0.13
12 เซนติเมตร	ปลูกกันที่ ปลูกในคุณภาพในไม้ผลิ	156	0.65
	ปลูกกันที่ ปลูกในคุณภาพในไม้ผลิ	84	0.35

Collett(1991) ได้นำข้อมูลนี้ไปวิเคราะห์การทดลอง ใช้โมเดลโลจิสติกเชิงเส้นจำนวน 5 โมเดล คือ

$$\text{โมเดล (1)} : \text{logit}(\hat{p}_{jk}) = b_0$$

$$\text{โมเดล (2)} : \text{logit}(\hat{p}_{jk}) = b_0 + b_1 X_j$$

$$\text{โมเดล (3)} : \text{logit}(\hat{p}_{jk}) = b_0 + b_2 X_k$$

$$\text{โมเดล (4)} : \text{logit}(\hat{p}_{jk}) = b_0 + b_1 X_j + b_2 X_k$$

$$\text{โมเดล (5)} : \text{logit}(\hat{p}_{jk}) = b_0 + b_1 X_j + b_2 X_k + b_3 X_j X_k$$

เมื่อ b_{jk} คือ ความน่าจะเป็นที่ต้นผลมีจากรากระดับชีวิตอยู่รอด เมื่อ rak ที่ปลูกมี
ขนาด j และปลูกในเวลา k ; $j = 1, 2$
 $k = 1, 2$

b_0 คือ ค่าคงที่

b_1 คือ สัมประสิทธิ์การลดถอยเนื่องจากความยาวราก

b_2 คือ สัมประสิทธิ์การลดถอย เนื่องจากเวลางอก

b_3 คือ สัมประสิทธิ์การลดถอย เนื่องจากอิทธิพลร่วมระหว่างความยาว
ราก และเวลาปลูก

X_j คือ ความยาวรากขนาดที่ j ; $X_1 = 0, X_2 = 1$

X_k คือ เวลาปลูกที่ k ; $X_1 = 0, X_2 = 1$

ผลการวิเคราะห์การลดถอยโลจิสติก ด้วยโปรแกรมสำเร็จรูป GLIM แสดงค่า deviance ของแต่ละโมเดล ในตาราง 4.2 และการวิเคราะห์ deviance ในตาราง 4.3

ตาราง 4.2 ค่า deviance ของแต่ละโมเดล

โมเดล	Deviance	d.f.
b_0	151.02	3
$b_0 + b_1 X_j$	105.18	2
$b_0 + b_2 X_k$	53.44	2
$b_0 + b_1 X_j + b_2 X_k$	2.29	1
$b_0 + b_1 X_j + b_2 X_k + b_3 X_j X_k$	0.00	0

ตาราง 4.3 การวิเคราะห์ deviance

Source of variation	Deviance	d.f.
ความยาวรากเมื่อไม่มีเวลาปลูก	$151.02 - 105.18 = 45.84$	1
เวลาปลูกเมื่อไม่มีความยาวราก	$151.02 - 53.44 = 97.58$	1
ความยาวรากเมื่อมีเวลาปลูก	$53.44 - 2.29 = 51.15$	1
เวลาปลูกเมื่อมีความยาวราก	$105.18 - 2.29 = 102.89$	1
มือทัชพลร่วม	2.29	1

เนื่องจากค่า deviance ที่ลดลง จากตาราง 4.3

พร้อมกับพิจารณาค่า

deviance จากตาราง 4.2 ควบคู่กันพบว่า ควรจะเลือกใช้ ไมเดล (4) ซึ่งมีเฉพาะอิทธิพลของปัจจัยหลัก คือความยาวรากกับเวลาที่ปลูก เพราะสามารถลดค่าของ deviance ลงได้มากที่สุด เท่ากับ 102.89 ซึ่งมีนัยสำคัญ และค่า deviance ของ ไมเดล (4) เท่ากับ 2.29 จะไม่มีนัยสำคัญที่ระดับ 10% (มีค่าน้อยกว่าค่า χ^2 จากตารางที่ d.f. = 1, ระดับนัยสำคัญ = 0.10) ซึ่งหมายความว่า ไมเดล (4) มีภาวะสารบดีกับค่าสังเกต

ค่าประมาณของพารามิเตอร์ จากโปรแกรมสำเร็จรูป GLIM เป็นดังนี้
 $\text{scale deviance} = 2.2938 \text{ at cycle 3}$
 $\text{d.f.} = 1$

	estimate	s.e.	parameter
1	-0.3039	0.1172	1
2	1.018	0.1455	LENG(2)
3	-1.428	0.1465	TIME(3)

แสดงว่า $b_0 = -0.3039$

$$b_1 = 1.018$$

$$b_2 = -1.428$$

จากโมเดล (4) : $\text{logit}(\hat{p}_{jk}) = b_0 + b_1 x_j + b_2 x_k$

จะได้ $\text{logit}(\hat{p}_{11}) = b_0 + b_1(0) + b_2(0) = -0.3039$

$$\text{logit}(\hat{p}_{12}) = b_0 + b_1(0) + b_2(1) = -0.3039 - 1.428 = -1.732$$

$$\text{logit}(\hat{p}_{21}) = b_0 + b_1(1) + b_2(0) = -0.3039 + 1.018 = 0.714$$

$$\begin{aligned} \text{logit}(\hat{p}_{22}) &= b_0 + b_1(1) + b_2(1) = -0.3039 + 1.018 - 1.428 \\ &= -0.714 \end{aligned}$$

และ $\hat{p}_{11} = [1 + \exp(-0.3039)]^{-1} = 0.425$

$$\hat{p}_{12} = [1 + \exp(-1.732)]^{-1} = 0.150$$

$$\hat{p}_{21} = [1 + \exp(-0.714)]^{-1} = 0.671$$

$$\hat{p}_{22} = [1 + \exp(-0.714)]^{-1} = 0.329$$

ซึ่งเปรียบเทียบกับค่าในตาราง 4.1 ได้ตาราง 4.4 เป็น

ตาราง 4.4 เปรียบเทียบความน่าจะเป็นระหว่างค่าที่ลังเกตได้ กับค่าประมาณ

ความยาวของ rakที่ตัด	เวลาที่ปลูก	ความน่าจะเป็นที่จะอยู่รอด ค่าลังเกต	ค่าประมาณ
6 เซนติเมตร	ปลูกทันที	0.446	0.425
	ปลูกในถุงใบไม้ผลิ	0.129	0.150
12 เซนติเมตร	ปลูกทันที	0.650	0.671
	ปลูกในถุงใบไม้ผลิ	0.350	0.329

การทดสอบสมมติฐานเกี่ยวกับ β_j อาจทำการทดสอบได้อีกวิธีหนึ่ง โดยการพิจารณา การแจกแจงของตัวประมาณ $\hat{\beta}_j = b_j$ ซึ่งมีการแจกแจงใกล้เคียงการแจกแจงแบบปกติ² จะได้ตัวทดสอบสถิติ เป็น

$$Z = \frac{b_j}{\sqrt{\text{Var}(b_j)}}$$

ซึ่งมีการแจกแจงใกล้เคียงการแจกแจงแบบปกติมาตรฐาน

หรือ

$$Z^2 = \frac{b_j^2}{\text{Var}(b_j)}$$

ซึ่งมีการแจกแจงใกล้เคียงการแจกแจงแบบไคสแควร์ ด้วย d.f. = 1

² McCullagh and Nelder. Generalized liner models (1983) pp.83

4.6 การคำนวณค่าความน่าจะเป็นจากโมเดล

วัตถุประสงค์หลักของการวิเคราะห์การผลถอยก็คือ การคำนวณค่า (prediction) หรือการประมาณค่า (estimation) ของตัวแปรตามจากโมเดลประมาณที่วิเคราะห์ได้ ในการวิเคราะห์การผลถอยโลจิสติก ค่าคำนวณที่ต้องการคือความน่าจะเป็นที่จะได้ข้อมูลที่มีค่าเท่ากับ 1 หรือ \hat{p}_1 .

$$\text{โดย } \hat{p}_1 = [1 + \exp(-\sum_{j=0}^{q-1} b_j X_{j+1})]^{-1}$$

เนื่องจากการวิเคราะห์การผลถอยโลจิสติก จะทำการแปลงค่า \hat{p}_1 ให้อยู่ในรูปของ $\text{logit}(\hat{p}_1)$ ซึ่งจะได้ว่า $\text{logit}(\hat{p}_1)$ มีความสัมพันธ์เชิงเส้นกับตัวแปรอิสระ X_{j+1} มีสมการ

ผลถอยโลจิสติกเชิงเส้น เป็น

$$\text{logit}(\hat{p}_1) = u_1 = \sum_{j=0}^{q-1} b_j X_{j+1} ; X_{01} = 1 \quad (4.6.1)$$

การหาค่าประมาณของ P_0 ณ ค่า X_{j_0} ที่กำหนด กระทำได้โดยการประมาณค่า $\text{logit}(\hat{p}_1)$ ณ ค่า X_{j_0} ที่กำหนด จากสมการ (4.6.1) จะได้ค่า $\text{logit}(\hat{p}_0)$ หรือ u_0 เป็นค่าประมาณของ $\text{logit}(P_0)$ หรือ U_0 โดยความแปรปรวนของ U_0 คือ $V(u_0)$ เป็น

$$V(u_0) = \sum_{j=0}^{q-1} X_{j_0}^2 V(b_j) + 2 \sum_{j,h=0}^{q-1} X_{j_0} X_{h_0} \text{Cov}(b_j, b_h) ; X_{00} = 1$$

และได้ช่วงความเชื่อมั่น $(1 - \alpha)\%$ สำหรับค่าประมาณของ U_0 มีค่าเป็น

$$u_0 - Z_{\alpha/2} \sqrt{V(u_0)} < U_0 < u_0 + Z_{\alpha/2} \sqrt{V(u_0)}$$

ถ้าให้ $u_{OL} = u_o - Z_{\alpha/2} \sqrt{V(u_o)} =$ ชีดจำากัดล่างของค่าประมาณของ U_o

$$u_{OU} = u_o + Z_{\alpha/2} \sqrt{V(u_o)} =$$
 ชีดจำากัดบนของค่าประมาณของ U_o

จะได้ $p_{OL} = [1 + \exp(-u_{OL})]^{-1} =$ ชีดจำากัดล่างของค่าประมาณของ P_o

$$p_{OU} = [1 + \exp(-u_{OU})]^{-1} =$$
 ชีดจำากัดบนของค่าประมาณของ P_o

4.7 การตรวจสอบโมเดล

เมื่อประมาณโมเดลการถดถอยโลจิสติกเชิงเส้น และวิเคราะห์ค่า deviance เพื่อทดสอบพารามิเตอร์ B_j จะได้โมเดลที่ต้องการแล้ว อาจพบว่าค่า deviance ของการทดสอบภาวะสารภูมิ ยังคงมีนัยสำคัญอยู่ ซึ่งหมายความว่าโมเดลที่ประมาณได้ยังไม่พอดีกับค่าลังเกต ควรทำการตรวจสอบโมเดล เพื่อหาสาเหตุของความไม่พอดี หรือความไม่เหมาะสมของโมเดล ซึ่งจะกระทำได้ในทำนองเดียวกับการพิจารณาสาเหตุของความไม่เหมาะสมของโมเดล ในการวิเคราะห์การถดถอยโลจิสติกเชิงเส้น ที่กล่าวในบทที่ 2

สำหรับความคลาดเคลื่อนที่ใช้ในการตรวจสอบความเหมาะสมของโมเดล ในการวิเคราะห์การถดถอยโลจิสติกเชิงเส้น ได้แก่

1. Pearson residual

$$X_i = \frac{y_i - n_i \hat{p}_i}{\sqrt{n_i \hat{p}_i (1 - \hat{p}_i)}}$$

ซึ่ง $\sum_{i=1}^n X_i^2$ ก็คือค่า Pearson's X^2 -statistic

และนี่ standardized Pearson residual เป็น

$$r_{Pi} = \frac{y_i - n_i \hat{p}_i}{\sqrt{v_i (1 - h_i)}}$$

$$\text{เมื่อ } v_i = n_i \hat{p}_i (1 - \hat{p}_i)$$

h_i = คือค่าในแนวทั่วไปของแมทริกซ์ H ขนาด $n \times n$

$$H = W^{1/2} X(X'WX)^{-1} X'W^{1/2}$$

W = weighted matrix ขนาด $n \times n$ ที่ใช้ในการประมาณค่าพารามิเตอร์ β_j

2. Deviance residual

$$d_i = \text{sgn}(y_i - \hat{y}_i) [2 \sum_{j=1}^n y_{ij} \ln \frac{y_{ij}}{\hat{y}_{ij}} + (n_i - y_i) \ln \frac{n_i - y_i}{n_i - \hat{y}_i}]^{1/2}$$

$\text{sgn}(y_i - \hat{y}_i)$ คือฟังก์ชันที่ทำให้ค่า d_i มีค่าเป็นบวก เมื่อ $y_i \geq \hat{y}_i$

และมีค่าเป็นลบ เมื่อ $y_i < \hat{y}_i$ และ $\sum_{i=1}^n d_i^2$ ก็คือค่า deviance

และมี standardized deviance residual เป็น

$$r_{di} = \frac{d_i}{\sqrt{1 - h_i}}$$