



สถิติที่ใช้ในการวิจัย

การศึกษาวิธีการตรวจสอบค่าผิดปกติในการวิเคราะห์ความถดถอยเชิงเส้นในครั้งนี้เป็นการศึกษาเพื่อเปรียบเทียบความน่าจะเป็นของความผิดพลาดประเภทที่ 1, ความน่าจะเป็นซึ่งค่าผิดปกติที่ถูกตรวจพบเป็นค่าผิดปกติจริงทุกค่า (p_1), ความน่าจะเป็นซึ่งทำให้เกิดมาซคิงเอฟเฟ็ค (p_2) และความน่าจะเป็นซึ่งทำให้เกิดชวอมทิงเอฟเฟ็ค (p_3) ของตัวสถิติของ 4 วิธี การซึ่งได้แก่ วิธีของเมอร์วิน จี มาราสิง (Mervyn G. Marasinghe, 1985), วิธีของฮาดีและไซมันนอฟฟ์ (Hadi and Simonoff, 1993), วิธีเวียนเกิดโดยลำดับ (Sequential recursive method) และวิธีเวียนเกิดดัดแปร (Modified recursive method) ผู้วิจัยจะศึกษาในกรณีที่มีค่าผิดปกติ 3 กรณีคือ 1, 2 และ 3 ค่า ตามลำดับ เมื่อความคลาดเคลื่อนมีการแจกแจงเบ้ผู้วิจัยจะทำการแปลงข้อมูลให้เข้าสู่การแจกแจงปกติโดยใช้วิธีการแปลงภายใต้การแจกแจงข้อมูลของ ฟาเบียน แอร์นันเดซและริชาร์ด เอ จอห์นสัน (Fabian Hernandez and Richard A. Johnson, 1980) และจะทำการประมาณสัมประสิทธิ์การถดถอย b ด้วยวิธีกำลังสองน้อยสุด

สมการการถดถอยที่ศึกษามีรูปแบบดังนี้

$$y = X\beta + \varepsilon$$

เมื่อ n คือขนาดตัวอย่าง

y คือเวกเตอร์ของตัวแปรตามซึ่งมีขนาด $n \times 1$

X คือเมทริกซ์ของตัวแปรอิสระขนาด $n \times p$ และมี $\text{rank} = p$

p คือจำนวนพารามิเตอร์ที่ต้องการประมาณ

β คือเวกเตอร์ของสัมประสิทธิ์การถดถอยที่ไม่ทราบค่าขนาด $p \times 1$

และ ε คือเวกเตอร์ของความคลาดเคลื่อนซึ่งมีขนาด $n \times 1$

โดยมีข้อสมมติว่า

1. $E(\underline{\varepsilon}) = 0$
2. $\text{cov}(\underline{\varepsilon}) = E(\underline{\varepsilon}\underline{\varepsilon}') = \sigma^2 I_n$

การวิเคราะห์ความถดถอยเชิงเส้นเป็นการศึกษาและอธิบายแนวโน้มเวกเตอร์ค่าสังเกต y ด้วยเวกเตอร์การประมาณ \hat{y} ซึ่ง $\hat{y} = X\hat{b}$ เมื่อ \hat{b} เป็นค่าประมาณของสัมประสิทธิ์การถดถอย ด้วยวิธีกำลังสองน้อยสุด ดังนั้นเราจะได้ว่า

$$\hat{b} = (X'X)^{-1} X'y \quad (2.1)$$

∴ เราแทนค่า \hat{b} ของสมการ(2.1) ใน $\hat{y} = X\hat{b}$ จะได้ว่า

$$\hat{y} = X(X'X)^{-1} X'y = Hy \quad (2.2)$$

เมื่อ $H = X(X'X)^{-1} X'$ เป็นเมทริกซ์ฉายเชิงตั้งฉาก (projection matrix) หรือเรียกอีกอย่างหนึ่งว่าแฮตเมทริกซ์ (hat matrix) ซึ่ง H มีคุณสมบัติดังนี้

1. สมมาตร (symetry) กล่าวคือ $H = X(X'X)^{-1} X' = H'$
2. นิจพด (idempotent) กล่าวคือ $H^2 = H$
3. เมทริกซ์ H มี h_{ii} เป็นสมาชิกในเส้นทแยงมุม (diagonal) ซึ่งมีค่าระหว่าง 0 ถึง 1 และ $h_{ii} = \sum_{j=1}^n h_{ij}^2 = h_{ii}^2 + \sum_{j \neq i} h_{ij}^2$
4. $\text{rank}(H) = \text{rank}(X) = p$

จากสัมประสิทธิ์การถดถอยที่ประมาณได้เราสามารถนำมาหาค่าตกค้าง (residual) e ซึ่งเป็นผลต่างระหว่างเวกเตอร์ค่าสังเกต y กับเวกเตอร์การประมาณ \hat{y} และเป็นค่าประมาณของความคลาดเคลื่อน ε ได้ดังนี้

$$\begin{aligned} \underline{e} &= \underline{y} - \hat{\underline{y}} \\ &= \underline{y} - X\hat{\underline{b}} \\ &= \underline{y} - X(X'X)^{-1} X'y \\ &= \underline{y} - H\underline{y} \\ &= (I - H)\underline{y} \end{aligned}$$

โดยที่ความแปรปรวนร่วมเกี่ยวของค่าตกค้าง (residual covariance) คือ

$$\text{cov}(\underline{e}) = (I - H) s^2$$

และมีความแปรปรวนค่าตกค้างของค่าสังเกตที่ i คือ

$$\text{Var}(e_i) = (1 - h_{ii}) s^2$$

เมื่อ s^2 เป็นตัวประมาณที่ไม่เอนเอียงของ σ^2 ซึ่ง $s^2 = \sum_{i=1}^n e_i^2 / (n - p)$; $i = 1, 2, \dots, n$

จากการวิเคราะห์ความถดถอยข้างต้น เราสามารถนำมาหาค่าตกค้าง(residual)ต่างๆได้ ดังนี้

ก) ค่าตกค้างมาตรฐานภายใน (Internally Studentized Residual)

$$R_i = e_i / (s \sqrt{1 - h_{ii}}) \quad ; i = 1, 2, \dots, n$$

ข) ค่าตกค้างมาตรฐานภายนอก (Externally Studentized Residual)

$$R_i^* = e_i / (s_{(-i)} \sqrt{1 - h_{ii}}) \quad ; i = 1, 2, \dots, n$$

เมื่อ $s_{(-i)}$ คือค่าตกค้างกำลังสองเฉลี่ย(residual mean square) เมื่อตัดค่าสังเกตที่ i ออกไป(วิธีการหาค่า $s_{(-i)}$ ดูที่ภาคผนวก ค.)

ค) ค่าตกค้างที่ปรับแล้ว (Adjusted Residual)

$$A_i = e_i / \sqrt{1 - h_{ii}} \quad ; i = 1, 2, \dots, n$$

ซึ่งเราจะนำค่าตกค้างเหล่านี้ไปใช้ในการตรวจสอบค่าผิดปกติในวิธีต่างๆ ดังนี้

สถาบันวิทยบริการ
จุฬาลงกรณ์มหาวิทยาลัย

1. วิธีของเมอร์วิน จี มาราสิง (Mervyn G. Marasinghe's method(MV))

ในปี ค.ศ. 1985 เมอร์วิน จี มาราสิง ได้เสนอวิธีการตรวจสอบค่าผิดปกติกรณีที่มีค่าผิดปกติหลายค่าในการวิเคราะห์ความถดถอยเชิงเส้น โดยมีหลักเกณฑ์คือ เราจะใช้ค่าตกค้างที่ปรับแล้ว (A_i) เป็นค่าบ่งชี้ค่าผิดปกติ โดยการนำค่าตกค้างที่ปรับแล้วซึ่งมีค่าสูงสุดในแต่ละรอบ (จำนวนค่าสูงสุดที่เป็นไปได้คือ k' รอบ) จำนวน k' ค่า โดยเปรียบเสมือนว่าค่าสังเกต k' ค่านี้เป็นค่าผิดปกติ จากนั้นจึงทำการคำนวณตัวสถิติทดสอบ

$$F_k' = (S - Q_k') / S$$

เมื่อ $S = (n - p) s^2$ เป็นผลบวกค่าตกค้างกำลังสอง (residual sum of square)

$$s^2 = \sum_{i=1}^n e_i^2 / (n - p) \text{ เป็นค่าตกค้างกำลังสองเฉลี่ย (residual mean square)}$$

และ $Q_k' = \sum_{i=1}^{k'} A_i^2$ เป็นผลบวกค่าตกค้างที่ปรับแล้วสูงสุดในแต่ละรอบ

โดยที่ $|A_i|$ คือค่าตกค้างที่ปรับแล้วสูงสุดรอบที่ i (maximum absolute adjusted residual) และ k' คือ จำนวนรอบสูงสุดที่เป็นไปได้ซึ่งเป็นจำนวนเต็มตั้งแต่ 2 ถึง 5 และ $k' \geq k$ เมื่อ k คือ จำนวนค่าผิดปกติ (ในการครั้งนี้เราจะศึกษากรณีที่ $k = 0, 1, 2$ และ 3 ตามลำดับ)

ขั้นตอนการทดสอบ

1. เราจะหาเซตของค่าสังเกตที่มี k' ค่า โดยพิจารณาจากค่าสัมบูรณ์ของค่าตกค้างที่ปรับแล้วสูงสุดในแต่ละรอบ ($|A_i|$) กล่าวคือ

1.1 คำนวณค่าตกค้างที่ปรับแล้วของแต่ละค่าสังเกตแล้วเลือกค่าที่ $\max |A_i|$ ตัดออกจากข้อมูลและนำค่าสังเกตดังกล่าวไปรวมไว้ในเซตย่อย K

1.2 คำนวณหา $\max |A_i|$ จากข้อมูลที่เหลือ $n - 1$ ค่าสังเกต แล้วตัดค่าสังเกตนั้นออกไปรวมไว้ในเซตย่อย K

1.3 ทำการคำนวณหา $\max |A_i|$ จากข้อมูลที่เหลือต่อไป จนกระทั่งได้เซตย่อย K ที่มีขนาด k'

2. คำนวณผลบวกค่าตกค้างกำลังสอง (residual sum of square) ในแต่ละรอบ (Q_k')

3. คำนวณค่าสถิติทดสอบ F_k'

4. ตรวจสอบค่าผิดปกติโดยมีสมมติฐานคือ

H : ไม่มีข้อมูลผิดปกติ

เทียบกับ K : มีข้อมูลผิดปกติอย่างมากที่สุด x^* ค่า

โดยนำค่า F_x^* ที่คำนวณได้เปรียบเทียบกับค่าขอบเขตวิกฤติในตารางของเมอร์วิน จี มาราซิง (1985) ถ้า F_x^* คำนวณ $<$ F_x^* ตาราง เราจะปฏิเสธสมมติฐานว่าง (H) ซึ่งแสดงว่าค่าสังเกตที่ i ที่มีค่าตรงกับค่า A_i เป็นค่าผิดปกติ

5. กรณีที่ปฏิเสธสมมติฐานว่าง เราจะตัดค่าสังเกตดังกล่าวในข้อ 4 ออกไปแล้ววิเคราะห์ข้อมูลขนาด $n - 1$ เช่นเดียวกับข้อ 1 ถึง 4 ทำซ้ำจนกว่าจะยอมรับสมมติฐานว่างจึงหยุดทำการทดสอบ

6. สรุปผลการทดสอบพร้อมสรุปค่าผิดปกติ



สถาบันวิทยบริการ
จุฬาลงกรณ์มหาวิทยาลัย

2. วิธีของฮาดีและไซมันอฟฟ์ (Hadi and Simonoff's method(HS))

วิธีนี้เป็นวิธีการตรวจสอบค่าผิดปกติหลายค่าในตัวแบบเชิงเส้น โดยการพยายามจัดข้อมูลออกเป็น 2 เซต คือ เซตของข้อมูลสะอาด (set of clean data point) และเซตของข้อมูลที่น่าจะเป็นค่าผิดปกติ (set of points that contain the potential outliers) และจะตรวจสอบข้อมูลที่น่าจะเป็นค่าผิดปกติโดยดูจากระยะที่มันเบี่ยงเบนไปจากเซตของข้อมูลสะอาด

ขั้นตอนการตรวจสอบ

1. เราจะหาเซตของข้อมูลสะอาด (แทนสัญลักษณ์ด้วยเซต M) ที่มีขนาด $b = [(n + p - 1) / 2]$ ดังนี้

1.1 คำนวณหาค่าสัมบูรณ์ของค่าตกค้างที่ปรับแล้ว ($|A_i|$) แล้วจัดเรียงค่าสังเกต ตาม $|A_i|$ จากน้อยไปมาก

1.2 แบ่งข้อมูลออกเป็น 2 เซตย่อยคือ

ก. **เซตย่อยพื้นฐาน** (basic subset : B) ซึ่งมีขนาด $s = p + 1$ ค่าสังเกตแรกที่ได้จากข้อ 1.1 มี X_B เป็นเมทริกซ์ไม่เอกฐาน แต่ถ้า X_B เป็นเมทริกซ์เอกฐาน เราจะทำให้เป็นเมทริกซ์ไม่เอกฐานโดยการรวมค่าสังเกตที่อยู่ถัดจาก $p + 1$ ค่าสังเกตแรกเข้าไปในเซตย่อยพื้นฐานจนกว่า X_B จะเป็นเมทริกซ์ไม่เอกฐาน และกำหนดให้ \tilde{b}_B เป็นตัวประมาณสัมประสิทธิ์การถดถอยของเซตย่อยพื้นฐาน

ข. **เซตย่อยไม่พื้นฐาน** (non-basic subset) มีขนาด $n - p - 1$ ค่าสังเกตหลังของข้อ 1.1 หรือมีขนาดเท่ากับจำนวนค่าสังเกตที่เหลือจากการจัดเซตย่อยพื้นฐานแล้ว

1.3 ทำการวิเคราะห์ความถดถอยเชิงเส้นของเซตย่อยพื้นฐานจะได้

$$e_{B_i} = y_i - x_i' \tilde{b}_B$$

$$\text{ซึ่งมี } \text{Var}(e_{B_i}) = \begin{cases} 1 - x_i' (X_B' X_B)^{-1} x_i & ; i \in B \\ 1 + x_i' (X_B' X_B)^{-1} x_i & ; i \notin B \end{cases}$$

เมื่อ e_{B_i} คือค่าตกค้างที่เกิดจาเซตย่อยพื้นฐาน

X_B คือเมทริกซ์ของตัวแปรอิสระในเซตย่อยพื้นฐาน

และ x_i' คือสมาชิกในแถวที่ i ของ X

* เซตของข้อมูลสะอาด หมายถึง เซตที่ไม่ใช่ข้อมูลผิดปกติรวมอยู่

1.4 คำนวณค่าสัมบูรณ์ค่าตกค้างมาตรฐาน ดังนี้

$$|B_i| = \begin{cases} \frac{|y_i - x_i' \hat{b}_B|}{\sqrt{1 - x_i' (X_B' X_B)^{-1} x_i}} & ; i \in B \\ \frac{|y_i - x_i' \hat{b}_B|}{\sqrt{1 + x_i' (X_B' X_B)^{-1} x_i}} & ; i \notin B \end{cases}$$

1.5 จัดเรียงค่าสังเกตตามค่า $|B_i|$ จากน้อยไปมาก แล้วจัดเป็นเซตย่อยพื้นฐานที่มีขนาด $s + 1$ ค่าสังเกต

1.6 คำนวณหาเซตย่อยพื้นฐานต่อไปเหมือนข้อ 1.3 ถึง 1.5 จนกระทั่งเซตย่อยพื้นฐานมีขนาดเท่ากับเซต M คือมีขนาด $h = (n + p - 1) / 2$

2. หาคำแบบที่เหมาะสมของเซต M แล้วคำนวณหาค่าตกค้างมาตรฐาน (M_i) ของเซต M โดยที่

$$|M_i| = \begin{cases} \frac{|y_i - x_i' \hat{b}_M|}{s_M \sqrt{1 - x_i' (X_M' X_M)^{-1} x_i}} & ; i \in M \\ \frac{|y_i - x_i' \hat{b}_M|}{s_M \sqrt{1 + x_i' (X_M' X_M)^{-1} x_i}} & ; i \notin M \end{cases}$$

3. จัดเรียงค่าสังเกตตาม $(|M_i|)$ จากน้อยไปมาก และกำหนด $M_{(n+1)}$ เป็นตัวสถิติอันดับที่ $h + 1$ ของ $|M_i|$ โดยที่ h คือขนาดของเซต M ปัจจุบัน เราจะทำการทดสอบสมมติฐานว่า

H : ค่าสังเกตที่ i เป็นค่าผิดปกติ ; $i = h + 1, h + 2, \dots, n$
 เทียบกับ K : ค่าสังเกตที่ i ไม่เป็นค่าผิดปกติ

ก) ถ้า $|M_{(n+1)}| \geq t_{(\alpha/2, (n+1), h-p)}$ เราจะสรุปว่าค่าสังเกตทั้งหมดที่ $|M_{(i)}| \geq t_{(\alpha/2, (n+1), h-p)}$ เป็นค่าผิดปกติและหยุดทำการทดสอบ

ข) กรณีอื่นให้จัดเซตย่อย M ใหม่ โดยให้มีขนาด $h+1$ ค่าสังเกต ถ้า $n = h+1$ สรุปว่าไม่มีค่าผิดปกติเกิดขึ้น จึงหยุดทำการทดสอบ แต่ถ้า $n \neq h + 1$ กลับไปเริ่มทำจากข้อ 2 ใหม่

3. วิธีของไคนิฟาร์ดและซวอลโด (Kiamifard and Swallow's method)

ในปี ค.ศ. 1990 ไคนิฟาร์ด(Kiamifard) และซวอลโด(Swallow) ได้เสนอวิธีการตรวจสอบค่าผิดปกติโดยใช้เกณฑ์ค่าตกค้างเวียนเกิด (recursive residual) ซึ่งค่าตกค้างเวียนเกิดคือ

$$|W_i| = \frac{|y_i - x_i' b_{i-1}|}{\sqrt{1 + x_i' (X_{i-1}' X_{i-1})^{-1} x_i}} \quad ; i = p+1, p+2, \dots, n$$

โดยที่ $b_{i-1} = (X_{i-1}' X_{i-1})^{-1} X_{i-1}' y_{i-1}$

X_{i-1} คือเมทริกซ์ขนาด $(i-1) \times p$ ที่ประกอบด้วยสมาชิก $i-1$ แถวแรกของ X

และ y_{i-1} คือเวกเตอร์ย่อยที่ประกอบด้วยสมาชิก $i-1$ แถวแรกของ y

ในปี ค.ศ. 1950 แพลคเกตต์ (Plackett), และในปี ค.ศ. 1975 บราวน์และคณะ (Brown et al.) ได้เสนอวิธีการหาค่า W_i โดยการหา b_i จาก b_{i-1} ดังนี้

$$b_i = b_{i-1} + \frac{(X_{i-1}' X_{i-1})^{-1} x_i (y_i - x_i' b_{i-1})}{1 + x_i' (X_{i-1}' X_{i-1})^{-1} x_i} \quad ; i = p+1, p+2, \dots, n$$

$$\text{เมื่อ } (X_i' X_i)^{-1} = (X_{i-1}' X_{i-1})^{-1} - \frac{(X_{i-1}' X_{i-1})^{-1} x_i x_i' (X_{i-1}' X_{i-1})^{-1}}{1 + x_i' (X_{i-1}' X_{i-1})^{-1} x_i}$$

และ $S_i = S_{i-1} + W_i^2 = (y - X_i b_i)' (y - X_i b_i)$; S_i คือผลบวกค่าตกค้างกำลังสองของค่าสังเกต i ค่า

วิธีการตรวจสอบค่าผิดปกติโดยใช้เกณฑ์ค่าตกค้างเวียนเกิดที่ ไคนิฟาร์ดและซวอลโดเสนอในปี 1990 มี 2 วิธี ได้แก่ วิธีเวียนเกิดโดยลำดับ (Sequential recursive method) และวิธีเวียนเกิดดัดแปร (modified recursive method) ซึ่งแต่ละวิธีมีขั้นตอนดังนี้

3.1 วิธีเวียนเกิดโดยลำดับ (sequential recursive method (SRM))

1. คำนวณหาค่าสัมบูรณ์ของค่าคงที่ที่ปรับแล้ว ($|A_i|$) แล้วจัดเรียงค่าดังกล่าวตาม $|A_i|$ จากน้อยไปมาก
2. ใช้ p ค่าดังกล่าวในข้อ 1 มาจัดเป็นเซตพื้นฐานเพื่อใช้คำนวณค่าคงที่เวียนเกิด (W_i)

3. คำนวณค่า W_i และคำนวณหาค่า $W_i / s_{(i)}$; $i = p+1, p+2, \dots, n$
4. คำนวณหาค่า $\max |W_i / s_{(i)}|$
5. ตรวจสอบค่าผิดปกติโดยใช้สมมติฐาน

H : ไม่มีข้อมูลผิดปกติ

เทียบกับ K : มีข้อมูลผิดปกติอย่างน้อย 1 ค่า

โดยเราจะทำการเปรียบเทียบค่าสถิติที่คำนวณได้กับขอบเขตวิกฤติจากตารางการแจกแจงที่ ถ้า $\max |W_i / s_{(i)}| > t_{\alpha/2, n-p-1}$ เราจะปฏิเสธสมมติฐานว่างและสรุปว่าค่าดังกล่าวที่มี $\max |W_i / s_{(i)}|$ เป็นค่าผิดปกติ

6. กรณีปฏิเสธสมมติฐานว่าง เราจะตัดค่าดังกล่าวในข้อ 5 ออกแล้ววิเคราะห์ข้อมูลขนาด $n - 1$ จากข้อ 1 ถึง 5 อีกครั้ง และทำไปจนกว่าจะยอมรับสมมติฐานว่างจึงหยุดทำการทดสอบ

7. สรุปผลการทดสอบพร้อมสรุปค่าผิดปกติ

3.2 วิธีเวียนเกิดดัดแปร (modified recursive method (MRM))

1. ในขั้นแรกเราจะทำเช่นเดียวกับข้อ 1 ถึง 5 ในวิธี SRM
2. ถ้าปฏิเสธสมมติฐานว่าง เราจะเปรียบเทียบ $|W_i / s_{(i)}|$ ที่เหลือกับ $t_{\alpha/2, n-p-1}$ ถ้าค่าดังกล่าวที่มี $|W_i / s_{(i)}| > t_{\alpha/2, n-p-1}$ แสดงว่าค่าดังกล่าวเป็นค่าผิดปกติ
3. สรุปผลการทดสอบพร้อมกับสรุปค่าผิดปกติ

การแปลงข้อมูลภายใต้การแจกแจงความน่าจะเป็น

ในกรณีที่การแจกแจงของความคลาดเคลื่อนไม่มีการแจกแจงปกติหรือกรณีทั่วไปที่ สเกลการวัดของตัวแปรไม่เหมาะสมกับสถิติวิเคราะห์ โดยทั่วไปนักสถิติจะแก้ปัญหาโดยการแปลง ข้อมูลเพื่อให้มีคุณสมบัติที่เหมาะสม การแปลงหลักๆ ที่ใช้กันอยู่คือ การแปลงแบบยกกำลัง (Power transformation) ของ Box และ Cox(1964) ซึ่งมีรูปแบบการแปลง คือ

$$y^{(\lambda)} = \begin{cases} y^\lambda - 1 & ; \lambda \neq 0 \\ \log y & ; \lambda = 0 \end{cases} ; y > 0$$

ซึ่งจะทำให้ $y^{(\lambda)}$ เข้าสู่การแจกแจงปกติ สำหรับวิธีการเลือกค่า λ ที่จะใช้ในการแปลงข้อมูลของ Box และ Cox จะเลือกจากค่า λ ที่ทำให้ $L(\lambda)$ มีค่าสูงสุด โดยที่

$$L(\lambda) = \begin{cases} \frac{n}{2} \log(\lambda^2) - \frac{n}{2} \log(RSS_\lambda) + (\lambda - 1) \sum_{i=1}^n \log(y_i) & ; \lambda \neq 0 \\ -\frac{n}{2} \log(RSS_\lambda) - \sum_{i=1}^n \log(y_i) & ; \lambda = 0 \end{cases}$$

เมื่อ RSS_λ คือผลบวกตกค้างกำลังสอง ซึ่งเราสามารถเลือกค่า λ โดยเลือกจากค่า λ ที่ทำให้ RSS_λ ต่ำสุดได้เช่นกัน

วิธีการแปลงของ Box และ Cox เป็นการแปลงข้อมูลภายใต้ตัวอย่างสุ่มของตัวแปรสุ่มเชิงบวกขนาด n โดยเราทำการเลือกพารามิเตอร์สำหรับการแปลงที่เหมาะสมกับตัวอย่างสุ่มชุดนั้น ดังนั้นตัวอย่างสุ่มแต่ละชุดในการแจกแจงเดียวกันก็จะทำให้ได้พารามิเตอร์ในการแปลงที่แตกต่างกัน ในปี ค.ศ. 1980 ฟาเบียน เฮอร์นันเดซและริชาร์ด เอ จอห์นสัน (Fabian Hernandez and Richard A. Johnson) ได้เสนอวิธีการแปลงข้อมูลเพื่อให้มีการแจกแจงปกติเมื่อทราบการแจกแจงของข้อมูล โดยในการแปลงจะเป็นการแปลงตัวอย่างสุ่มภายใต้การแจกแจงความน่าจะเป็น และใช้จำนวนสารนิเทศกูลแบค-เลเบอร์ (Kullback - Leibler, 1968) เป็นหลักในการเลือกพารามิเตอร์ λ ที่จะใช้ในการแปลง ซึ่งมีรายละเอียดดังนี้

กำหนดให้ h_1 และ h_2 เป็นฟังก์ชันความหนาแน่นของความน่าจะเป็นและต่อเนื่องของตัวแปร Z มีค่าความแตกต่างระหว่างการแจกแจง คือ

$$I[h_1, h_2] = \int h_1(t) \log \left\{ \frac{h_1(t)}{h_2(t)} \right\} dt$$

เราเรียกความแตกต่างระหว่างการแจกแจงดังกล่าวว่า จำนวนสารสนเทศจุดแบบค-เอนเบอร์

ถ้า y เป็นตัวแปรสุ่มเชิงบวก (positive random variable) ที่มีฟังก์ชันความหนาแน่น (probability density function (pdf)) คือ $g(\cdot)$ ในการเลือก λ ที่จะใช้ในการแปลง y ให้มีการแจกแจงปรกติ เราจะเลือก λ ที่ทำให้จำนวนสารสนเทศ $I[f_\lambda; \phi_{\mu\sigma}]$ ซึ่งแสดงค่าความแตกต่างระหว่างการแจกแจงของข้อมูลที่แปลงด้วย λ กับข้อมูลที่มีการแจกแจงปรกติด้วยค่าเฉลี่ย μ และส่วนเบี่ยงเบนมาตรฐาน σ มีค่าต่ำสุดโดยที่

$$I[f_\lambda; \phi_{\mu\sigma}] = \int f_\lambda(a) \log \left\{ \frac{f_\lambda(a)}{\phi_{\mu\sigma}(a)} \right\} da$$

เมื่อ $f_\lambda(a)$ เป็นฟังก์ชันความหนาแน่นของความน่าจะเป็นหนาแน่นของตัวแปร $A = y^{(\lambda)}$ และ $\phi_{\mu\sigma}$ เป็นฟังก์ชันความหนาแน่นของความน่าจะเป็นหนาแน่นของการแจกแจงปรกติด้วยค่าเฉลี่ย μ และส่วนเบี่ยงเบนมาตรฐาน σ

ในขั้นแรกเราจะเลือกค่า μ และ σ ด้วย f_λ ซึ่งจะทำให้ A มีการแจกแจงปรกติโดยประมาณก่อนแล้วจึงจะหาค่า λ ที่ทำให้จำนวนสารสนเทศมีค่าต่ำสุด หรือมีระยะทางน้อยที่สุด กล่าวคือ

กำหนด λ ดังที่ และ ให้

$$y^{(\lambda)} = \begin{cases} y^\lambda - 1 & ; \lambda \neq 0 \\ \log y & ; \lambda = 0 \end{cases} ; y > 0$$

สมมติว่าค่าคาดหวัง $E_p(y^{2\lambda})$ และ $E_p[(\log y)^2]$ มีค่าจำกัด ดังนั้นค่า μ และ σ ที่ทำให้จำนวนสารสนเทศ $I[f_\lambda; \phi_{\mu\sigma}]$ มีค่าต่ำสุดคือ

$$\mu(\lambda) = E_p(y^{(\lambda)})$$

$$\text{และ} \quad \sigma^2(\lambda) = E_p[y^{(2\lambda)} - E_p(y^{(\lambda)})]^2 = V_p(y^{(\lambda)})$$

สำหรับ $\mu(\lambda)$ และ $\sigma^2(\lambda)$ ที่เลือกแล้วเราจะได้

$$\begin{aligned} G(\lambda) &= \min_{\mu, \sigma} I[f_\lambda; \phi_{\mu\sigma}] & (2.3) \\ &= \frac{1}{2} [\log(2\pi) + 1] + E_g \{ \log[g(y)] \} \\ &\quad + (1-\lambda) E_g \{ \log(y) \} + \frac{1}{2} \log[\text{Var}_g(y^{(\lambda)})] \end{aligned}$$

ซึ่ง λ ที่ทำให้ $G(\cdot)$ มีค่าต่ำสุดเป็นพารามิเตอร์ที่เราจะใช้ในการแปลง y ให้มีการแจกแจงปรกติ

ตัวอย่างที่ 1 ให้อ่านงานสารนิเทศ $I[f_\lambda; \phi_{\mu\sigma}]$ เพื่อใช้หาค่า λ . สำหรับการแจกแจงแกมมา งานสารนิเทศๆ สำหรับการแจกแจงแกมมา หาได้ดังนี้ การแจกแจงแกมมามีสูตรของความหนาแน่นดังนี้

$$g(y) = \frac{\beta^\alpha}{\Gamma(\alpha)} y^{\alpha-1} e^{-\beta y} \quad ; y > 0 ; \alpha, \beta > 0$$

ถ้าเรากำหนดให้ $\beta = 1$ จะได้ฟังก์ชันความหนาแน่นของการแจกแจงแกมมา

$$\text{อยู่ในรูปของ } g(y) = \frac{1}{\Gamma(\alpha)} y^{\alpha-1} e^{-y} \quad ; y > 0 ; \alpha > 0$$

จากสมการที่ (2.3) เราจะหาค่าของแต่ละพจน์ดังนี้

$$\begin{aligned} E_g \{ \log[g(y)] \} &= E_g \{ (\alpha - 1) \log(y) - y - \log[\Gamma(\alpha)] \} \\ &= E_g \{ (\alpha - 1) \log(y) \} - E_g(y) - E_g \{ \log[\Gamma(\alpha)] \} \\ &= (\alpha - 1) \int_0^\infty \frac{\log(y) y^{\alpha-1} e^{-y}}{\Gamma(\alpha)} dy - \int_0^\infty \frac{y y^{\alpha-1}}{\Gamma(\alpha)} dy - \log[\Gamma(\alpha)] \\ &= (\alpha - 1) \left[\frac{\Gamma'(\alpha)^*}{\Gamma(\alpha)} \right] - \frac{\Gamma(\alpha + 1)}{\Gamma(\alpha)} - \log[\Gamma(\alpha)] \\ &= (\alpha - 1) \psi(\alpha) - \alpha \frac{\Gamma(\alpha)}{\Gamma(\alpha)} - \log[\Gamma(\alpha)] \\ &= \alpha \psi(\alpha) - \psi(\alpha) - \alpha - \log[\Gamma(\alpha)] \end{aligned}$$

$$\begin{aligned} E_g \{ \log(y) \} &= \int_0^\infty \frac{\log(y) y^{\alpha-1} e^{-y}}{\Gamma(\alpha)} dy \\ &= \psi(\alpha) \end{aligned}$$

$$* \int_0^\infty \log(y) y^{\alpha-1} e^{-y} dy = \Gamma'(\alpha)$$

$$** \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} = \psi(\alpha)$$

$$\begin{aligned} \text{Var}_g(y^{(\lambda)}) &= V_g\left(\frac{y^{2\lambda}-1}{\lambda}\right) \\ &= \frac{1}{\lambda^2} V(y^{2\lambda}) \end{aligned}$$

จาก $\text{Var}(y) = E(y^2) - \{E(y)\}^2$ เราจะได้ว่า

$$\begin{aligned} \text{Var}_g(y^{2\lambda}) &= E_g(y^{2\lambda}) - \{E_g(y^{2\lambda})\}^2 \\ &= \int_0^{\infty} \frac{y^{2\lambda} y^{\alpha-1} e^{-y}}{\Gamma(\alpha)} dy - \left\{ \int_0^{\infty} \frac{y^{\lambda} y^{\alpha-1} e^{-y}}{\Gamma(\alpha)} dy \right\}^2 \\ &= \frac{\Gamma(2\lambda + \alpha)}{\Gamma(\alpha)} - \left\{ \frac{\Gamma(\lambda + \alpha)}{\Gamma(\alpha)} \right\}^2 \\ &= \frac{\Gamma(\alpha)\Gamma(2\lambda + \alpha) - \{\Gamma(\lambda + \alpha)\}^2}{\{\Gamma(\alpha)\}^2} \end{aligned}$$

ดังนั้นเราจะได้
$$\text{Var}_g(y^{(\lambda)}) = \frac{\Gamma(\alpha)\Gamma(2\lambda + \alpha) - \{\Gamma(\lambda + \alpha)\}^2}{\lambda^2 \{\Gamma(\alpha)\}^2}$$

เมื่อเราแทนค่าผลลัพธ์ของพจน์ต่างๆในสมการ(2.3) จะได้ว่า

$$\begin{aligned} G(\lambda) &= \frac{1}{2} [\log(2\pi) + 1] + \{\alpha\psi(\alpha) - \psi(\alpha) - \alpha - \log\Gamma(\alpha)\} \\ &\quad + \{(1-\lambda)\psi(\alpha)\} + \left\{ \frac{1}{2} \log \left[\frac{\Gamma(\alpha)\Gamma(2\lambda + \alpha) - [\Gamma(\lambda + \alpha)]}{\lambda^2 [\Gamma(\alpha)]^2} \right] \right\} \\ &= \frac{1}{2} [\log(2\pi) + 1] + \alpha\psi(\alpha) - \psi(\alpha) - \alpha - \log\Gamma(\alpha) \\ &\quad + \psi(\alpha) - \lambda\psi(\alpha) + \frac{1}{2} \log \left[\frac{\Gamma(\alpha)\Gamma(2\lambda + \alpha) - [\Gamma(\lambda + \alpha)]}{\lambda^2} \right] \\ &\quad - \frac{1}{2} \log\{[\Gamma(\alpha)]^2\} \\ &= \frac{1}{2} \log(2\pi) + 1 - 2\log\Gamma(\alpha) + \alpha[\psi(\alpha) - 1] \\ &\quad - \lambda\psi(\alpha) + \frac{1}{2} \log \left\{ \frac{\Gamma(\alpha)\Gamma(2\lambda + \alpha) - [\Gamma(\lambda + \alpha)]^2}{\lambda^2} \right\} \end{aligned}$$

หมายเหตุ ψ คือฟังก์ชันไดแกมมา*(digamma function) ซึ่งมีรูปแบบฟังก์ชัน คือ

$$\begin{aligned}\psi(\alpha) &= \frac{d}{d\alpha} \log[\Gamma(\alpha)] \\ &= \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} \\ &= \ln \alpha - \frac{1}{2\alpha} - \frac{1}{12\alpha^2} + \frac{1}{120\alpha^4} + \frac{1}{252\alpha^6} + \dots\end{aligned}$$

ค่าประมาณของฟังก์ชันไดแกมมาเราจะดูจากตารางที่ ง.1 ในภาคผนวก



สถาบันวิทยบริการ
จุฬาลงกรณ์มหาวิทยาลัย

* William H. Beyer, CRC Standard Mathematical Tables, 26th ed. (Florida: CRC Press, 1981) , p. 399