

ข้อมูลที่เปลี่ยนแปลงค่าไปได้ง่าย จะต้องมื่อน้ำหนักหรือความน่าเชื่อถือน้อยกว่าข้อมูลที่ค่อนข้างคงที่

นอกจากนี้ความน่าเชื่อถือของข้อมูลยังเกี่ยวข้องกับหน่วยที่ใช้วัดอีกด้วย

2. ลักษณะความน่าเชื่อถือ

ข้อมูลแต่ละชนิดย่อมมีลักษณะ และขนาดของความน่าเชื่อถือแตกต่างกันไปตามธรรมชาติ ข้อมูลที่เกี่ยวกับเรื่องส่วนตัวที่น่าอับอาย ข้อมูลเกี่ยวกับทรัพย์สินหรือผลประโยชน์ และข้อมูลเกี่ยวกับความจำ และการคำนวณจะมีความน่าเชื่อถือต่ำ ส่วนข้อมูลที่เกี่ยวข้องกับความผูกพัน ความใกล้ชิดคุ้นเคยข้อมูลที่เกี่ยวข้องกับความเข้าใจ ตลอดจนข้อมูลที่ไม่เกี่ยวข้องกับผลประโยชน์หรือเรื่องที่ควรปกปิดจะมีความน่าเชื่อถือ หรือน้ำหนักสูงกว่า

ตัวอย่างเช่น ข้อมูลที่ได้รับจากคำถามเรื่องจำนวนที่ดินที่ใช้ในการเพาะปลูกจะมีความถูกต้องหรือน่าเชื่อใต้อน้อยกว่าข้อมูลจากคำถามเรื่องจำนวนที่ดินถือครอง (กรณีนี้คือกรณีของความจำหรือการคำนวณ) ข้อมูลจากคำถามเรื่องรายได้ของครอบครัวจะมีความถูกต้องหรือน่าเชื่อใต้อน้อยกว่าข้อมูลจากคำถามเรื่องเงินเดือนของผู้ตอบเอง (กรณีนี้คือกรณีที่เกี่ยวข้องกับผลประโยชน์)

3. ความน่าเชื่อถือของชุดข้อมูล (Reliability of a Specific Record)

กรณีนี้เป็นกรณีของข้อตกลงของ LSC กล่าวคือข้อมูลรายการใดที่ไม่สามารถตรวจสอบได้ด้วยวิธี Gross Check และวิธีตรวจสอบความแน่นอนได้หรือตรวจด้วยวิธีทั้งสองดังกล่าวแล้วไม่อาจหาข้อยุติได้ กรณีนี้ควรให้น้ำหนักรายการข้อมูลนั้นต่ำกว่ารายการอื่น

4. ชนิดของน้ำหนักใน LSC มี 2 ชนิดคือ w-weight เป็นน้ำหนักที่กำหนดให้แก่รายการข้อมูลและ u-weight เป็นน้ำหนักที่กำหนดให้แก่สมการแนบเนียน น้ำหนัก w ควรมากกว่าน้ำหนัก u

ค. การตัดสินใจใช้ข้อมูลทดแทน

เมื่อวิเคราะห์ได้เวกเตอร์ \bar{X} โดยวิธี LSC แล้ว ให้ใช้ข้อมูลทดแทน \bar{X} แทนข้อมูล Y เดิมทุกประการ โดยถือว่าข้อมูลทดแทน \bar{X} เป็นข้อมูลที่ถูกต้อง

13.2 Principal Component Method (PC)

ก. เหตุผลและความจำเป็น

ในการบรรณาธิกรณข้อมูลนั้น โดยปกติเราจำเป็นต้องกระทำกับข้อมูลทุกรายการ หลังจากการตรวจสอบแบบสอบถามครบทั้ง n ชุด เราสามารถสรุปได้ทันทีว่า รายการข้อมูลใดมีความน่าเชื่อถือรายการข้อมูลใดไม่น่าเชื่อถือและต้องหาทางแก้ไขรายการนั้นต่อไป

วิธีการที่จะใช้ในการแก้ไขก็คือ เมื่อตรวจสอบพบว่ารายการข้อมูลใดมีความบกพร่อง อยู่เช่น ขาดความแน่นอน ไม่สมเหตุสมผลหรือข้อมูลสูญหาย (โดยปกติจะมีมากกว่า 1 รายการเสมอ) ให้ผู้ตรวจแก้ดำเนินการตรวจสอบดูว่า รายการข้อมูลที่บกพร่องนั้นจะสามารถแก้ไข โดยอาศัยข้อสนเทศจากข้อมูลอื่น ๆ รายการใดบ้าง สมมุติว่าจากการตรวจสอบพบว่ารายการ ข้อมูลที่พบว่าบกพร่องนั้นสามารถได้โดยอาศัยข้อสนเทศจากข้อมูลรายอื่น ๆ k รายการ ปัญหา ที่จะต้องพิจารณากันต่อไปก็คือ จะนำข้อสนเทศเหล่านั้นมาใช้ประโยชน์ได้อย่างไร ทางออก สำหรับปัญหานี้จึงมีอยู่ 2 นัยคือ

(1) อาศัยวิธีการของสมการถดถอย (Regression Analysis) โดยถือว่าข้อมูลที่พบว่า บกพร่องเป็นค่าของตัวแปรตาม ส่วนข้อมูลจากรายการอื่น ๆ ที่เกี่ยวข้องกับรายการที่บกพร่อง ดังกล่าว เช่นช่วยอธิบายความผันผวนเปลี่ยนแปลงในค่าของรายการนั้นเป็นค่าของตัวแปร อิสระ

ให้ Y = รายการข้อมูลที่บกพร่อง

X_i 's ; $i = 1, 2, \dots, k$ คือรายการข้อมูลที่เกี่ยวข้องหรือเป็นรายการที่อธิบาย ความผันผวนใน Y ได้

ดังนั้นแบบจำลองที่พึงใช้ก็คือ

$$Y = \beta_0 + \sum_{i=1}^k \beta_i X_i + u \text{ เมื่อ } u \text{ คือ stochastic error}$$

โดยมีข้อตกลงว่า $E(u) = 0, V(u) = \sigma^2, E(u_i u_j) = \begin{cases} 0 & ; i \neq j \\ \sigma^2 & ; i = j \end{cases}$

และ X_i 's ต้องเป็นอิสระต่อกัน หรือเกี่ยวข้องกันในอัตราค่อนข้างต่ำกล่าวคือ X ต่างก็ทำหน้าที่อธิบายความผันผวนในค่าของ Y โดยที่ X 's ต้องไม่สัมพันธ์กันเองจะเกิดปัญหาที่สำคัญคือแยกไม่ออกจาก X ตัวใดแน่ที่มีอิทธิพลต่อ Y และที่สำคัญก็คือค่าประมาณของ β 's จะ

เป็น indeterminate form ($\beta_i = \frac{0}{0}$) และมีความน่าเชื่อถือต่ำ ($V(\beta_i) \rightarrow \infty$)

(2) อาศัยวิธีการ PC

วิธี PC คือกรณีเฉพาะของวิธีวิเคราะห์ Factor Analysis กล่าวคือภายหลังจากรวบรวมพบว่ารายการข้อมูลที่บกพร่อง (Y) สามารถอ้างอิงได้โดยอาศัยข้อมูลสาเหตุจากรายการข้อมูลใด (X's) แล้ว ให้ตรวจสอบดูว่า X's มีความเกี่ยวข้องกันอยู่ในอัตราใด โดยอาศัยวิธีวัดค่าสหสัมพันธ์คือ r_{x_i, x_j} ; $i, j = 1, 2, \dots, k$ แล้วสร้างตัวแปรอิสระขึ้นมาใหม่โดยอาศัยสหสัมพันธ์เหล่านี้ในลักษณะที่ตัวแปรตัวใหม่เกิดขึ้นจากการอาศัยประโยชน์ของความสัมพันธ์ระหว่าง X's นั่นคือ $Z_i = \sum_j^k a_{ij} X_j$; $i = 1, 2, \dots, k$ และตัวแปรตัวใหม่เหล่านี้จะถูกสร้างขึ้นมาหลายตัวได้ตามความจำเป็นและความเหมาะสมซึ่ง Z's จะเป็นอิสระต่อกันเสมอ และ Z_1 จะมีความสำคัญสูงกว่า Z_2 Z_2 มีความสำคัญสูงกว่า Z_3 เป็นต้นนี้เรื่อยๆ ไป คำว่ามีความสำคัญสูงกว่า หมายความว่าเป็นตัวแปรที่ดึงดูดเอาอิทธิพลร่วมระหว่าง X's มาไว้ในตัวได้มากกว่า โดยนัยนี้ วิธี PC จึงมีข้อจำกัดน้อยกว่าวิธีวิเคราะห์ความถดถอย เพราะไม่ต้องสนใจว่า X's จะเป็นอิสระต่อกันหรือไม่ และด้วยเหตุที่โดยปกติแล้วตัวแปร X's มักมีส่วนเกี่ยวข้องเชื่อมโยงถึงกันอยู่เสมอ (Intercorrelate) การสร้างข้อมูลทดแทนด้วยวิธี PC จึงต้องตามเหตุผลและสถานะการณ์ได้มากกว่า

ข. ความหมายของ PC

ให้ $Y = f(X's)$ โดยที่ Y คือรายการที่ข้อมูลที่พบว่าบกพร่อง X's คือรายการข้อมูลที่ใช้เป็นแหล่งอ้างอิงของ y

ดังนั้น PC จึงหมายถึงเซตของตัวแปรตัวใหม่คือ Z_i ; $i = 1, 2, \dots, k$ ที่เกิดขึ้นจากการประกอบกัน (Linear Combination) ของตัวแปร X's โดยที่

$$Z_i = \sum_j^k a_{ij} X_j \quad ; \quad i = 1, 2, \dots, k$$

เรียก Z_i ว่า PC ที่ i และสัมประสิทธิ์ a_{ij} เรียกว่า loadings โดยที่ a_{ij} เป็นสัมประสิทธิ์ของ X_j ที่ถูกสร้างขึ้นหรือเลือกขึ้นมาในลักษณะที่ทำให้ PC Z_i สอดคล้องกับคุณสมบัติ 2 ประการคือ

(1) PC เป็นอิสระต่อกัน (Uncorrelated หรือ Orthogonal)

(2) PC ตัวแรกจะดึงดูดเอาอิทธิพลร่วมระหว่าง X's ใ้มากที่สุด PC ตัวลำดับถัดไปจะดึงดูดใ้ได้น้อยลงตามลำดับ

ค. วิธีดำเนินการสร้าง PC

การสร้าง PC ให้สอดคล้องกับคุณสมบัติของ loadings 2 ประการข้างต้นมีขั้นตอนดำเนินการกว้าง ๆ ขั้นตอนดังนี้

ขั้นที่ 1 กะประมาณค่าของ loadings a_{ij} ; $i, j = 1, 2, \dots, k$ ซึ่งค่า จะ

ช่วยแปลงรูป X's ให้เป็นตัวแปรอิสระ Z's (Orthogonal Artificial Variable) และตรวจสอบนัยสำคัญของ เพื่อที่ว่า Z แต่ละตัวควรเกิดขึ้นจากมาประกอบกันของ X's ก็ตัว

ขั้นที่ 2 ตัดสินใจคัดเลือก PC ว่าควรใช้ PC ตัวใด และควรจะคงจำนวน PC ใ้กี่ตัวจึงเพียงพอในการอธิบายความผันผวนของ Y โดยปกติจำนวน PC จะน้อยกว่าจำนวน X's กล่าวคือจำนวนสูงสุดของ PC จะไม่เกินจำนวน X's

ขั้นที่ 3 วิเคราะห์สมการถดถอย

$$Y = \gamma_0 + \sum_{i=1}^r r_i Z_i + u ; r \leq k ; Z_i = \sum_{j=1}^k \hat{a}_{ij} Z_j ; i = 1, 2, \dots, k$$

เพื่อกะประมาณค่าของพารามิเตอร์ γ_i ; $i = 0, 1, 2, \dots, r$ การกะประมาณใ้ใช้วิธี Ordinary Least Square (OLS) ทั้งนี้ถือข้อตกลงเกี่ยวกับ Stochastic error เช่นเดียวกับข้อตกลงของ u ในการวิเคราะห์ความถดถอย

$$\text{ใ้สมการ } \hat{Y} = \hat{\gamma}_0 + \sum_{i=1}^r \hat{\gamma}_i Z_i ; r \leq k$$

แล้วกะประมาณค่าของ Y โดยอาศัยค่า X's เดิม จากสมการถดถอยข้างต้น

ขั้นที่ 4 กะประมาณค่าของ Y ด้วยช่วงเชื่อมั่น $(1-\alpha)$ 100% ใ้ แล้วพิจารณาตัดสินใจเกี่ยวกับความถูกต้องของข้อมูล Y เดิมคือ

ก. ถ้า $\hat{Y}_L \leq Y \leq \hat{Y}_U$ ใ้ถือว่าข้อมูล Y เพิ่มถูกต้อง

ข. ถ้า $Y > \hat{Y}_U$ ใ้ใ้ \hat{Y}_U เป็นข้อมูลทดแทนของ Y

ค. ถ้า $Y < \hat{Y}_L$ ให้ใช้ \hat{Y}_L เป็นข้อมูลทดแทนของ Y

ค.1 การประมาณค่าของ loading a_{ij} ของ Z_i

การประมาณค่าของสัมพันธ์ประสิทธิ์ a_{ij} ให้ดำเนินการดังนี้

(1) หาสหสัมพันธ์ (Simple Correlation) ระหว่าง X's คือ $R_{X_i X_j}$;

$$i \neq j = 1, 2, \dots, k \quad (r_{X_i X_j} = 1 \text{ เมื่อ } i=j \text{ เมื่อ } r_{X_i X_i} = \frac{\sum_s^n X_{is} X_{js}}{\sqrt{\sum_s^n X_{is}^2 \sum_s^n X_{js}^2}} \text{ แล้วจัดเรียงไว้ใน$$

ตารางเรียกว่าตารางสหพันธ์ (Correlation Table) หรือเมตริกซ์สหพันธ์ (Correlation Matrix)

เมตริกซ์ดังกล่าวจะเป็น Symmetric Matrix เพราะ $r_{X_i X_j} = r_{X_j X_i}$ สำหรับทุกค่าของ $i \neq j =$

$1, 2, \dots, k$ และสมาชิกในแนวทแยง (Main Diagonal) จะเท่ากับ 1 เสมอเพราะ $r_{X_i X_i} =$

$r_{X_i X_i} = 1$ สำหรับทุกค่าของ i, j ที่ $i = j = 1, 2, \dots, k$ ดังนี้

	X_1	X_2	X_3	X_4	.	.	.	X_k
X_1	1	$r_{X_1 X_2}$	$r_{X_1 X_3}$	$r_{X_1 X_4}$.	.	.	$r_{X_1 X_k}$
X_2	$r_{X_2 X_1}$	1	$r_{X_2 X_3}$	$r_{X_2 X_4}$.	.	.	$r_{X_2 X_k}$
X_3	$r_{X_3 X_1}$	$r_{X_3 X_2}$	1	$r_{X_3 X_4}$.	.	.	$r_{X_3 X_k}$
X_4	$r_{X_4 X_1}$	$r_{X_4 X_2}$	$r_{X_4 X_3}$	1	.	.	.	$r_{X_4 X_k}$
.
.
.
X_k	$r_{X_k X_1}$	$r_{X_k X_2}$	$r_{X_k X_3}$	$r_{X_k X_4}$.	.	.	1
c_j	$\sum_i^k r_{X_i X_1}$	$\sum_i^k r_{X_i X_2}$	$\sum_i^k r_{X_i X_3}$	$\sum_i^k r_{X_i X_4}$.	.	.	$\sum_i^k r_{X_i X_k}$
loadings	l_1	l_2	l_3	l_4	.	.	.	l_k

$$1. x_{is} = X_{is} - \bar{X}_i, x_{js} = X_{js} - \bar{X}_j$$

(2) รวมค่าของสหสัมพันธ์ในทุก ๆ สดมภ์ (หรือจะใช้แถวแทนสดมภ์

ก็ได้) ดังนั้น ผลรวมของสหสัมพันธ์ในสดมภ์ที่ j คือ $c_j = \sum_i r_{x_i x_j} ; j = 1, 2, \dots, k$

(3) หายอดรวมทั้งหมดของสหสัมพันธ์ (Grand Total) แล้วถอดรากที่

สองนั่นคือ

$$\sqrt{GT} = \left(\sum_j \sum_i r_{x_i x_j} \right)^{1/2}$$

(4) คำนวณหาค่าประมาณของ loading a_{ij} ของ Z_i (PC ตัวที่ 1)

ด้วยการนำ \sqrt{GT} ไปหารผลรวม c_j

$$\text{นั่นคือ} \quad \hat{l}_i = \hat{a}_{ij} = \frac{c_j}{\sqrt{GT}} ; j = 1, 2, \dots, k$$

$$\text{หรือ} \quad \hat{l}_i = \hat{a}_{ij} = \left(\sum_j r_{x_i x_j} \right) / \left(\sum_j \sum_i r_{x_i x_j} \right)^{1/2} ; j = 1, 2, \dots, k$$

$$\begin{aligned} \text{ดังนั้น} \quad Z_i &= \hat{l}_1 X_1 + \hat{l}_2 X_2 + \dots + \hat{l}_k X_k && \text{หรือ} \\ &= \hat{a}_{i1} X_1 + \hat{a}_{i2} X_2 + \dots + \hat{a}_{ik} X_k \end{aligned}$$

ข้อสังเกต 1. Subscript 1 ของ \hat{a}_{ij} หมายถึง PC ตัวที่ 1

2. ขอให้สังเกตว่า \hat{l}_i ก็คือสัมประสิทธิ์สหสัมพันธ์รูปหนึ่ง

(5) หา latent root (หรือ eigen value หรือ Charaxctertic root) ของ PC ตัวที่ 1 โดยที่ latent root λ คือผลรวมกำลังสอง (Sum of Square, SS) ของ loadings กล่าวคือ latent root ของ PC ตัวที่ m คือ

$$\lambda_m = \sum_j \hat{a}_{mj}^2 ; m \leq k$$

ดังนั้น latent root ของ Z_i คือ

$$\lambda_i = \sum_j \hat{a}_{ij}^2$$

ด้วยเหตุที่ผลรวมของ latent root ของ PC ทุกตัวมีค่าเท่ากับจำนวนสูงสุดของ

FC (หรือจำนวน X's กล่าวคือ $\sum_{m=1}^k \lambda_m = k$)

ดังนั้น latent root ของ PC ตัวใดจึงเป็นค่าแสดงความสำคัญของ PC ตัวนั้น เพราะ latent root จะแสดงอัตราความผันแปรรวม (Total Variation) ที่ PC ตัวนั้นดึงดูดออกมาจากกลุ่มของ X's ปริมาณดังกล่าวสามารถวัดเป็นปริมาณของอัตราความผันแปรรวม (Percentage of Total Variation) ได้ดังนี้

$$\text{อัตราของความผันแปรที่ } Z_m \text{ ดึงดูดเอาไว้ได้} = \frac{\lambda_m}{k} \times 100\% ; m = 1, 2, \dots, k$$

โดยนัยดังกล่าว Z_1 จึงมี latent root สูงกว่า Z_2 แสดงว่า Z_1 มีความสำคัญหรือดึงดูดเอาความผันผวนหรือมีอิทธิพลรวมของ X's ไว้ได้มากที่สุด Z_2 มีความสำคัญรองลงไป เป็นเช่นนี้เรื่อย ๆ ไป หรือนัยหนึ่ง $\lambda_1 > \lambda_2 > \lambda_3 > \dots > \lambda_k$

(6) จากขั้นตอนการทั้ง 5 ขั้น ทำให้ได้ Z_1 พร้อมทั้งปริมาณอิทธิพลรวม X's ดึงดูดเอาไว้ได้คือ λ_1 ดังนี้

$$Z_1 = \hat{a}_{11} X_1 + \hat{a}_{12} X_2 + \dots + \hat{a}_{1k} X_k ; \lambda_1 ; \text{ความสำคัญ} = \frac{\lambda_1}{k} \times 100\% \text{ ขั้น}$$

ต่อไปดำเนินการหา Z_2

การหา PC ตัวที่ 2 หรือ Z_2 ให้ดำเนินการเป็นขั้น ๆ ดังนี้

(1) สร้างตาราง Residual Correlation หรือ Residual Correlation Matrix โดยหักลบอิทธิพลของ Z_1 ออกจากตารางสหพันธ์ที่ใช้หา Z_1 ค่าสหพันธ์ผลลัพธ์เรียกว่า Residual Correlation $r_{r_i X_j}^*$ โดยที่ $r_{r_i X_j}^* = r_{X_i X_j} - \hat{a}_{1i} \hat{a}_{1j}$; $i, j = 1, 2, \dots, k$

ทั้งนี้เพราะ \hat{a}_{ij} คือสัมประสิทธิ์สหสัมพันธ์รูปหนึ่ง

ตัวอย่างเช่น เดิม $r_{X_1 X_2} = r_{X_2 X_1} = .40$, $\hat{a}_{11} = .30$, $\hat{a}_{12} = .50$ ดังนั้น Residual Correlation

$$r_{X_1 X_2}^* - r_{X_1 X_2} = \hat{a}_{11} \hat{a}_{12} = 40. - (30 - 50) = 25$$

ดังนั้น Residual Correlation Table สำหรับ Z_2 จึงปรากฏดังนี้

	X_1	X_2	X_k
X_1	$r_{X_1X_1}^* = 1 - \hat{a}_{11}^2$	$r_{X_1X_2}^* = r_{X_1X_2} - \hat{a}_{11} \hat{a}_{12} \dots$		$r_{X_1X_k}^* = r_{X_1X_k} - \hat{a}_{11} \hat{a}_{1k}$
X_2	$r_{X_2X_1}^* = r_{X_2X_1} - \hat{a}_{12} \hat{a}_{11}$	$r_{X_2X_2}^* = 1 - \hat{a}_{12}^2 \dots$		$r_{X_2X_k}^* = r_{X_2X_k} - \hat{a}_{12} \hat{a}_{1k}$
X_3	$r_{X_3X_1}^* = r_{X_3X_1} - \hat{a}_{13} \hat{a}_{11}$	$r_{X_3X_2}^* = r_{X_3X_2} - \hat{a}_{13} \hat{a}_{12} \dots$		$r_{X_3X_k}^* = r_{X_3X_k} - \hat{a}_{13} \hat{a}_{1k}$
\vdots	\vdots	\vdots	\vdots	\vdots
\vdots	\vdots	\vdots	\vdots	\vdots
\vdots	\vdots	\vdots	\vdots	\vdots
X_k	$r_{X_kX_1}^* = r_{X_kX_1} - \hat{a}_{1k} \hat{a}_{11}$	$r_{X_kX_2}^* = r_{X_kX_2} - \hat{a}_{1k} \hat{a}_{12} \dots$		$r_{X_kX_k}^* = 1 - \hat{a}_{1k}^2$
รวม	$\sum_i^k r_{X_iX_1}^*$	$\sum_i^k r_{X_iX_2}^*$...	$\sum_i r_{X_iX_i}$
loadings	\hat{a}_{21}	\hat{a}_{22}	...	\hat{a}_{2k}

ขอให้สังเกตว่า $r_{X_iX_i}^* \neq 1$ เมื่อ $i=j$

(2) จาก Residual Correlation Matrix Loading ได้

$$\hat{a}_{2j} = \left(\sum_i^k r_{X_iX_j}^* \right) / \left(\sum_i^k \sum_i^k r_{X_iX_i}^* \right)^{1/2} ; j = 1, 2, \dots, k$$

(3) คำนวณหา latent root ของ Z_2 ได้

$$\lambda_2 = \sum_j^k \hat{a}_{2j}^2$$

พร้อมทั้งปริมาณความผันแปรในอิทธิพลของ X 's ที่ Z_2 ดึงดูดไว้ภายหลังจากที่ Z_1 ได้ดึงดูดไปแล้วนั้นคือ

$$\text{อัตราความผันแปรที่ } Z_2 \text{ ดึงดูดไว้ได้} = \frac{\lambda_2}{k} \times 100\%$$

(4) คำนวณหา Z_2 จะพบว่า

$$Z_2 = \hat{a}_{21} X_1 + \hat{a}_{22} X_2 + \dots + \hat{a}_{2k} X_k ; \text{ความสำคัญ} = \frac{\lambda_2}{k} \times 100\%$$

การหา Z_3 ให้หาโดยการสร้าง Residual Correlation Matrix ขึ้นใหม่จากตาราง Residual Correlation ของ Z_2 แล้วดำเนินการแบบเดิมจะได้

$$Z_3 = \hat{a}_{31} X_1 + \hat{a}_{32} X_2 + \dots + \hat{a}_{3k} X_k ; \text{ความสำคัญ} = \frac{\lambda_3}{k} \times 100\%$$

โดยที่

$$\hat{a}_{3j} = \left(\sum_i^k r_{X_i X_i}^{**} \right) / \left(\sum_i^k \sum_j^k r_{Z_i Z_j}^{**} \right)^{1/2} ; j=1, 2, \dots, k$$

$$r_{X_i X_j}^{**} = r_{X_i X_j}^* - \hat{a}_{2i} \hat{a}_{2j} ; i, j = 1, 2, \dots, k$$

$$\lambda_3 = \sum_j^k \hat{a}_{3j}^2$$

สำหรับ Z_4, Z_5, \dots, Z_k ก็ดำเนินการแบบเดียวกันนี้ กล่าวคือสร้าง Residual Correlation Matrix ขึ้นมาโดยอาศัย Residual Correlation Matrix ของ PC ตัวก่อนหรือลำดับก่อน

ค.2 การตัดสินใจคัดเลือกองค์ประกอบของ Z_i

การพิจารณาคัดเลือกองค์ประกอบของ Z_i หมายถึงการตัดสินใจว่า Z_i ควรประกอบไปด้วยตัวแปร X 's ที่ตัวจริงจะถือว่าเหมาะสมหรือเพียงพอ เกณฑ์การตัดสินใจก็คือ นัยสำคัญของ a_{ij} ถ้า a_{ij} มีนัยสำคัญ แสดงว่า X_j ควรคงอยู่หรือควรถือว่า X_j เป็นองค์ประกอบของ Z_i หรือ X_j ควรเป็นปัจจัยที่มีอิทธิพลร่วมให้เกิดตัวแปร Z_i

ดังนั้น จาก PC ใด ๆ คือ

$$Z_i = \sum_j^k \hat{a}_{ij} X_j ; i=1, 2, \dots, k$$

$$\text{หรือ } Z_i = \hat{a}_{i1} X_1 + \hat{a}_{i2} X_2 + \dots + \hat{a}_{ik} X_k$$

สิ่งที่ต้องการคือ $a_{ij} ; j = 1, 2, \dots, k$ มีนัยสำคัญ ระดับนัยสำคัญ α หรือไม่ นั่นคือเราต้องทำการทดสอบสมมติฐานต่อไปนี้

$H_0 : a_{ij} = 0$ vs $H_1 : a_{ij} \neq 0 ; j = 1, 2, \dots, k$ สำหรับ i ใด ๆ โดยที่

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

เช่น ต้องการทดสอบว่า Z_i ควรอาศัยอิทธิพลร่วมของ X ตัวใดบ้าง ดังนั้นจากโครงสร้าง

$$Z_i = \sum_j^k a_{ij} X_j \quad \text{หรือ}$$

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

เราต้องทำการทดสอบสมมติฐานต่อไปนี้

$$H_0 : a_{ij} = 0 \text{ vs } H_1 : a_{ij} \neq 0 : j = 1, 2, \dots, k$$

การตรวจสอบนัยสำคัญของ $a_{ij} ; j = 1, 2, \dots, k$ สำหรับ fixed $i \quad i = 1, 2, \dots, k$ นั้นมีวิธีการต่าง ๆ หลายวิธีดังนี้

(1) Empirical Test

การทดสอบวิธีนี้เป็นวิธีที่ค่อนข้างหยาบแต่ใช้ได้ผลดีสำหรับการเก็บขนาดตัวอย่าง n มีค่าน้อยเท่ากับ 50 หน่วย ($n \leq 50$) กล่าวคือ loadings a_{ij} ใดจะมีนัยสำคัญก็ต่อเมื่อ $|a_{ij}| > .30$ หรือ $-.30 > a_{ij} > .30$ หรือนัยหนึ่งปฏิเสธสมมติฐานหลัก $H_0 : a_{ij} = 0$ และยอมรับสมมติฐานรอง $H_1 : a_{ij} \neq 0$ เมื่อ $|a_{ij}| > .30$ หรือ a_{ij} มากกว่า $\pm .30$

(2) Standard error Test

$$\text{เนื่องจาก } \hat{a}_{ij} = \left(\sum_j^k r_{x_i x_j} \right) / \left(\sum_j^k \sum_j^k r_{x_i x_j} \right)^{1/2} ; i = 1, 2, \dots, k$$

คือรูปหนึ่งของสัมประสิทธิ์สหสัมพันธ์ การทดสอบ $H_0 : a_{ij} = 0$ vs $H_1 : a_{ij} \neq$

$0 : j = 1, 2, 3, \dots, k$ สำหรับ fixed $i \quad i = 1, 2, 3, \dots, k$ จึงทดสอบเช่นเดียวกับการทดสอบ

สหสัมพันธ์ r (Pearson's Product Moment) กล่าวคือใช้ตารางค่าวิกฤติของ r ตรวจสอบนัยสำคัญของ a_{ij} ค่าวิกฤติของ r ก็คือส่วนเบี่ยงเบนมาตรฐานของ r ซึ่งจะใช้เป็นส่วนเบี่ยงเบน

มาตรฐานของ a_{ij} ต่อไปนี้เอง

การตรวจสอบนัยสำคัญของ a_{ij} ให้กระทำดังนี้

จาก $H_0 : a_{ij} = 0$ vs $H_1 : a_{ij} \neq 0 ; j = 1, 2, \dots, k$ สำหรับ fixed $i \quad i = 1, 2, \dots, k$

เราจะปฏิเสธสมมติฐานหลักและยอมรับสมมติฐานรองเมื่อ $|a_{ij}| > r_{n,\alpha}$ เมื่อ $r_{n,\alpha}$ คือค่าวิกฤติ ณ. ระดับนัยสำคัญ α องศาความเป็นอิสระ n ซึ่งเป็นค่าในตารางวิกฤติของ

ขนาดตัวอย่าง	ค่าวิกฤติสหสัมพันธ์ ณ. ระดับนัยสำคัญ	
	5 %	1 %
5	.755	.875
10	.676	.714
15	.603	.605
20	.545	.538
25	.498	.488
30	.463	.440
35	.435	.417
40	.412	.394
45	.392	.370
50	.375	.346
60	.352	.328
70	.335	.308
80	.322	.290
90	.310	.272
100	.300	.255
150	.272	.209
200	.252	.182
250	.238	.163
500	.206	.115

(2) Burt-Bank Test (B-B)

วิธีทดสอบที่ (1) และ (2) เป็นวิธีที่ค่อนข้างหายากทั้งนี้เพราะมิได้คำนึงถึงจำนวนตัวแปร X 's ที่ร่วมเป็นองค์ประกอบร่วมของ Z และลำดับความสำคัญของ Z เอง ทำให้การตรวจจับนัยสำคัญของวิธีทดสอบดังกล่าวมีช่องทางให้คัดค้านได้ วิธีทดสอบของ B-B จะแก้จุดอ่อนข้อนี้ได้โดยนำทั้งจำนวนตัวแปร X และลำดับความสำคัญของ Z มาร่วมพัฒนาตัวทดสอบเพื่อตรวจสอบนัยสำคัญของ loadings โดยนำค่าทั้งสองดังกล่าวมาปรับปรุงค่าส่วนเบี่ยงเบนมาตรฐานของ r ให้รัดกุมเสียก่อน ก่อนที่จะนำค่าดังกล่าวไปใช้เป็นค่าวิกฤติ

ให้ $k =$ จำนวนตัวแปร X

$m =$ Subscript ของ Z ซึ่งแสดงลำดับความสำคัญของ Z (อย่าลืมว่า Z_1 สำคัญกว่า Z_2 Z_2 สำคัญกว่า Z_3 , ..., Z_{k-1} สำคัญกว่า Z_k)

$s(r_{x_i x_j}) =$ ค่าวิกฤติ (ส่วนเบี่ยงเบนมาตรฐาน) ของ r ณ. ระดับสำคัญ α อิสระ n (ขนาดตัวอย่าง) ซึ่งในที่นี้ $s(r_{x_i x_j})$ ก็คือค่าวิกฤติของ a_{ij} หรือ $S(r_{x_i x_j}) = r_{n, \alpha}$ ดังนั้นค่าวิกฤติของ a_{ij} ที่ปรับปรุงแล้วคือ $s(a_{mj})$ โดยที่

$$\begin{aligned} S(a_{mj}) &= s(r_{x_i x_j}) \sqrt{\frac{k}{k+1-m}} \\ &= r_{n, \alpha} \sqrt{\frac{k}{k+1-m}} \end{aligned}$$

นั่นคือวิกฤติของ loading ตัวที่ $j; j = 1, 2, \dots, k$ ของ PC ตัวที่ m คือ

$$s(a_{mj}) = r_{n, \alpha} \sqrt{\frac{k}{k+1-m}} ; j=1, 2, \dots, k$$

ซึ่งเราจะปฏิเสธสมมติฐานหลัก $H_0: a_{ij} = 0 ; j = 1, 2, \dots, k ; i = m$ และยอมรับสมมติฐาน

รอง $H_1: a_{ij} \neq 0 ; j = 1, 2, \dots, k ; i=m$ เมื่อ

$$|a_{mj}| > s(a_{mj}) : i = 1, 2, \dots, k \text{ สำหรับ fixed } i ; i = m$$

ค.3. การพิจารณาคัดเลือกจำนวน PC ที่จะเหลือไว้ในการวิเคราะห์

โดยปกติจำนวน PC จะอยู่ไม่เกินจำนวนตัวแปร X และเหตุที่สมการ $Y = f(X's)$ นั้นควรจะคัดเลือกตัวแปร Z ไว้เฉพาะตัวที่มีความสำคัญจริง ๆ หรือตัวที่ดึงดูดเอาอิทธิพลร่วมระหว่าง X ไว้ได้มากจริง ๆ เท่านั้น วิธีการตัดสินใจคัดเลือกมีอยู่หลายวิธี จะกล่าวถึงเพียง 2 วิธี ดังนี้

(1) Kaiser's Criterion

วิธีการคัดเลือกของไกเซอร์พิจารณาเลือก PC เฉพาะตัวที่ให้ latent root มากกว่า 1 เท่านั้นตัวที่ให้ latent root น้อยกว่า 1 จะถูกตัดทิ้ง

นั่นคือคัดเฉพาะ Z_m ที่ $\lambda_m > 1$ ไว้ในการวิเคราะห์ แต่โดยเหตุที่ $\lambda_1 > \lambda_2 > \dots > \lambda_k$ ดังนั้นเราจะหยุดวิเคราะห์ PC ตัว $(r+1)$ ถ้าพบว่า $\lambda_r < 1$ ซึ่งจะช่วยประหยัดเวลาและแรงงานลงไปได้มากทั้งนี้เพราะ $\lambda_{r+1}, \lambda_{r+2}, \dots, \lambda_k$ จะมีค่าน้อยกว่า 1 อย่างแน่นอน

กฎเกณฑ์ของไกเซอร์เหมาะที่จะใช้เป็นหลักในการคัดเลือก PC หรือมีความน่าเชื่อถือได้เฉพาะเมื่อ $20 < k < 50$ หรือมีจำนวนตัวแปร X อยู่ระหว่าง 20-50 ตัว ถ้า $k > 50$ จะมีแนวโน้มให้เหลือ PC ไว้มากเกินไป ขณะที่ $k < 20$ จะมีแนวโน้มให้เหลือ PC ไว้น้อยเกินไป

(2) Bartlett's Criterion

กฎเกณฑ์ของบาร์ตเลตต์มีลักษณะการคัดเลือกคล้ายวิธีของไกเซอร์ในแง่ที่พิจารณาความสำคัญของ PC โดยพิจารณาจาก latent root โดยมีหลักเกณฑ์การคงจำนวน PC ไว้เช่นเดียวกัน ต่างกับวิธีของไกเซอร์ที่บาร์ตเลตต์มีวิธีการคัดเลือก PC ตัวต่อไปได้รับกุ่มกว่า ดังนี้

สมมติว่า $\lambda_1, \lambda_2, \dots, \lambda_r ; r < k$ มีค่าสูงและแตกต่างกัน ซึ่งทำให้เชื่อได้ว่า Z_1, Z_2, \dots, Z_r ควรจะคงไว้ในสมการ $Y = f(Z's)$ ปัญหาก็คือ PC ที่เหลืออยู่อีก $k-r$ ดังนั้นควรจะคงไว้ด้วยหรือไม่ ควรจะนำไปเพิ่มในสมการ $Y = f(X's)$ ได้อีกหรือไม่ ทั้งนี้ให้ถือว่า Z_r มีความสำคัญน้อยที่สุดและไม่ควรนำเข้าไปในสมการโครงสร้าง

ด้วยปัญหาดังกล่าว บาร์ตเลตต์จึงเสนอให้ทำการทดสอบสมมุติฐานดังนี้

H_0 : PC ที่เหลืออยู่ $(k-r)$ ตัวไม่มีความสำคัญ หรือเสนอเป็นสัญลักษณ์ได้ดังนี้

$$H_0 : \lambda_{r+1} = \lambda_{r+2} = \dots = \lambda_k$$

ขณะที่สมมติฐานรองคือ

H_1 : ควรเพิ่มจำนวน PC เข้าไปในสมการโครงสร้างได้อีก หรือเสนอเป็นสัญลักษณ์ได้ ดังนี้

H_1 : λ_i not all equal ; $i = r+1, r+2, \dots, k$

ตัวสถิติที่ใช้ทดสอบคือ

$$X^{*2} = n \ln \left\{ (\lambda_{r+1} \cdot \lambda_{r+2} \cdot \dots \cdot \lambda_k)^{-1} \left(\frac{\lambda_{r+1} + \lambda_{r+2} + \dots + \lambda_k}{k-r} \right)^{k-r} \right\}$$

โดยที่ $X^{*2} \sim X^2_{\gamma, \alpha}$ เมื่อ $\gamma = \frac{1}{2} (k-r-1) (k-r+2)$

และปฏิเสธสมมติฐานหลักเมื่อ $X^{*2} > X^2_{\gamma, \alpha}$ ซึ่งหมายความว่าควรเพิ่ม PC เข้า

ไปในสมการโครงสร้างได้ การเพิ่มจะต้องเพิ่มควรวละตัวแล้วทดสอบสมมติฐานต่อไปจนกว่าจะยอมรับสมมติฐานหลักซึ่งหมายความว่า PC ที่เหลืออยู่มีความสำคัญมากไม่ควรนำเพิ่มเข้าไปในสมการ

ค. วิเคราะห์สมการถดถอย $Y = \gamma_0 + \sum_i^r \gamma_i Z_i + u$

การวิเคราะห์สมการถดถอย $Y = \gamma_0 + \sum_i^r \gamma_i Z_i + u$ ยึดถือวิธี OLS ในการ

ประมาณค่าพารามิเตอร์ γ_i ; $i = 0, 1, 2, \dots, r$

$$\text{โดยที่ } \hat{\gamma} = (\hat{\gamma}_0, \hat{\gamma}_1, \dots, \hat{\gamma}_r)^T = (Z^T Z)^{-1} Z^T Y$$

$$\hat{\sigma}^2 = \frac{1}{n-r} (Y^T Y - \hat{\gamma}^T Z^T Y)$$

$$\text{และ } V(\hat{\gamma}) = \hat{\sigma}^2 (Z^T Z)^{-1}$$

ทั้งนี้อาศัยการวิเคราะห์แบบ Linear Hypothesis Model of Full Rank เนื่องจากตัวแปร Z's เป็นอิสระต่อกัน (Uncorrelate หรือ Orthogonal) โดยวิเคราะห์จากแบบจำลอง

$$Y_j = \gamma_0 + \sum_i^r \gamma_i Z_{ij} + u_j \quad ; j = 1, 2, 3, \dots, n$$

$$\text{เมื่อ } Z_i = \sum_j^k \hat{\alpha}_{ij} X_j ; i = 1, 2, 3, \dots, r$$

หรือจัดในรูปสมการเมตริกซ์

$$Y_{n \times 1} = Z_{n \times (r+1)} \gamma_{(r+1) \times 1} + u_{n \times 1}$$

และด้วยเหตุที่งานวิเคราะห์นี้มีจุดมุ่งหมายเพื่อการพยากรณ์ คือ คาดหมายข้อมูลทดแทน เครื่องมือที่ใช้ในการตัดสินใจเลือก Best fitted model จึงใช้สัมประสิทธิ์แห่งการตัดสินใจ \bar{R}^2 (coefficient of Determination) โดยที่

$$\bar{R}^2 = 1 - \{ (Y^T Y - \hat{\gamma}^T Z^T Y) / (n-r) \} / \{ (Y^T Y - n\bar{Y}^2) / (n-1) \}$$

นั่นคือ Estimated model ใดให้ \bar{R}^2 สูงสุดก็พึงเลือกใช้ model นั้น ทั้งนี้โดยมีต้องคำนึงว่า $\gamma_i ; i = 1, 2, \dots, r$ จะมีนัยสำคัญ หรือไม่มากนัก

สมการพยากรณ์ที่ต้องการคือ

$$\hat{Y} = \hat{\gamma}_0 + \hat{\gamma}_1 Z_1 + \hat{\gamma}_2 Z_2 + \dots + \hat{\gamma}_r Z_r$$

$$\begin{matrix} (s_{\hat{\gamma}_0}) & (s_{\hat{\gamma}_1}) & (s_{\hat{\gamma}_2}) & (s_{\hat{\gamma}_3}) & \dots & (s_{\hat{\gamma}_r}) \\ & (t_{\hat{\gamma}_1}) & (t_{\hat{\gamma}_2}) & (t_{\hat{\gamma}_3}) & \dots & (t_{\hat{\gamma}_r}) \end{matrix}$$

$$\text{โดยที่ } ; i = 1, 2, \dots, r-1, s_{\hat{\gamma}_i} = \sqrt{c_{ii}}$$

เมื่อ c_{ii} คือสมาชิกในตำแหน่งที่ $(i+1, i+1)$ ของเมตริกซ์

$$\sigma^2 (Z^T Z)^{-1} \text{ และ } t_{\hat{\gamma}_i} = \frac{\hat{\gamma}_i - \gamma_i}{\sqrt{c_{ii}}} \sim t_{n-1, 1-\alpha/2} \text{ ซึ่งใช้ตรวจสอบนัยสำคัญของ } \gamma_i ; i = 1, 2,$$

\dots, r หรือสมมติฐาน $H_0 : \gamma_i = 0 ; \text{ vs } H_1 : \gamma_i \neq 0 ; i = 1, 2, \dots, r$ และปฏิเสธสมมติฐาน

หลัก $H_0 : \gamma_i = 0 \quad i = 1, 2, \dots, r$ เมื่อ $|t_{\hat{\gamma}_i}| > t_{n-1, 1-\alpha/2}$

สำหรับข้อตกลงเบื้องต้นเกี่ยวกับ Stochastic error ให้ยึดถือข้อตกลงของ u ใน OLS ทุกประการคือ

$$E(u) = 0, V(u) = \sigma^2, E(u_i u_j) = \begin{cases} 0 & ; i \neq j \\ \sigma^2 & ; i = j \end{cases}, u \sim N(0, \sigma^2)$$

ค.4 การพยากรณ์ค่าของ Y ด้วยช่วงเชื่อมั่น

ให้ Z_{iF} = ค่าพยากรณ์ของ Z_i ใน forecasting period ในที่นี้จะถือว่า $Z_{iF} = Z_i$

n = ขนาดตัวอย่าง ในที่นี้ n คือจำนวนชุดของแบบสอบถามหรือจำนวน record

$\hat{\sigma}^2$ = ค่าประมาณของ σ^2

ดังนั้นค่า ค่าพยากรณ์ของค่าเฉลี่ยของ Y_F คือ

$$\hat{Y}_{F_i} = \hat{\gamma}_0 + \hat{\gamma}_1 Z_{1F_i} + \hat{\gamma}_2 Z_{2F_i} + \dots + \hat{\gamma}_r Z_{rF_i}; i = 1, 2, \dots, n$$

ให้ $Z_F = (1, Z_{1F}, Z_{2F}, \dots, Z_{rF})$ เมื่อ Z_F เป็นเมตริกซ์ขนาด $n \times (r+1)$

ซึ่งสมาชิกของ Z_F คือ Column Vector ขนาด $n \times 1$ กล่าวคือ

$$1 = (1, 1, 1, \dots, 1)^T, Z_{1F} = (Z_{21}, Z_{22}, Z_{23}, \dots, Z_{2n})^T$$

$$Z_{2F} = (Z_{31}, Z_{32}, Z_{33}, \dots, Z_{3n})^T, \dots, Z_{rF} = (Z_{(r+1) \times 1}, Z_{(r+1) \times 2}, \dots, Z_{(r+1) \times n})^T$$

$$\text{ให้ } \hat{\gamma} = (\hat{\gamma}_0, \hat{\gamma}_1, \dots, \hat{\gamma}_r)^T \quad Y = (Y_1, Y_2, Y_3, \dots, Y_n)^T$$

ดังนั้น

$$Y = Z \hat{\gamma} = (1, Z_{1F}, Z_{2F}, \dots, Z_{rF}) \begin{pmatrix} \hat{\gamma}_0 \\ \hat{\gamma}_1 \\ \hat{\gamma}_2 \\ \vdots \\ \hat{\gamma}_r \end{pmatrix}$$

หรือ

$$\begin{pmatrix} \hat{Y}_1 \\ \hat{Y}_2 \\ \hat{Y} \\ \vdots \\ \hat{Y}_n \end{pmatrix} = \begin{pmatrix} 1 & Z_{21} & Z_{31} & \dots & Z_{r1} \\ 1 & Z_{22} & Z_{32} & \dots & Z_{r2} \\ 1 & Z_{23} & Z_{33} & \dots & Z_{r3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & Z_{2n} & Z_{3n} & \dots & Z_{rn} \end{pmatrix} \begin{pmatrix} \hat{\gamma}_0 \\ \hat{\gamma}_1 \\ \hat{\gamma}_2 \\ \vdots \\ \hat{\gamma}_r \end{pmatrix}$$

และเนื่องจาก $\hat{Y} = (1, Z_{1F}, Z_{2F}, \dots, Z_{rF}) (\hat{\gamma}_0, \hat{\gamma}_1, \hat{\gamma}_2, \dots, \hat{\gamma}_r)^T$ โดยกำหนดให้ Z_{iF}

เป็นแวกเตอร์ของค่าคงที่ของ X ซึ่งในที่นี้ถือว่า $Z_{iF} = Z_i$ หมายความว่าถือว่าค่าเดิมของ Z_i เป็นค่าเดียวกับค่าพยากรณ์ของ Z_i ดังนั้น \hat{Y} จึงเป็น linear Combination ของตัวแปรสุ่ม $\hat{\gamma}_0, \hat{\gamma}_1, \hat{\gamma}_2, \dots, \hat{\gamma}_r$

$$\text{ดังนั้น } V(\hat{Y}_{Fi}) = V(\hat{\gamma}_0 + \hat{\gamma}_1 Z_{1i} + \hat{\gamma}_2 Z_{2i} + \dots + \hat{\gamma}_r Z_{ri}) ; i = 1, 2, \dots, n$$

$$\text{หรือ} = \sum_{i=0}^{r+1} \sum_{j=0}^{r+1} V_{ij} Z_i Z_j ; V_{ij} = \begin{cases} V(\gamma_i) ; i=j \\ \text{Cov}(\gamma_i, \gamma_j) ; i \neq j \end{cases}$$

จะเห็นได้ว่า $V(\hat{Y}_{Fi})$ เป็น Quadratic form

ดังนั้น

$$V(\hat{Y}_{Fi}) = (1, Z_{2i}, Z_{3i}, \dots, Z_{ri}) \begin{pmatrix} V(\lambda_0) & \text{Cov}(\gamma_0, \gamma_1) & \text{Cov}(\gamma_0, \gamma_2) & \dots & \text{Cov}(\gamma_0, \gamma_r) \\ \text{Cov}(\gamma_1, \gamma_0) & V(\gamma_1) & \text{Cov}(\gamma_1, \gamma_2) & \dots & \text{Cov}(\gamma_1, \gamma_r) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(\gamma_r, \gamma_0) & \text{Cov}(\gamma_r, \gamma_1) & \text{Cov}(\gamma_r, \gamma_2) & \dots & V(\gamma_r) \end{pmatrix} \begin{pmatrix} 1 \\ Z_{2i} \\ \vdots \\ \vdots \\ Z_{ri} \end{pmatrix}$$

$i = 1, 2, \dots, n$

$$= (1, Z_{2i}, Z_{3i}, \dots, Z_{ri}) \cdot (Z^T Z)^{-1} \sigma^2 \cdot (1, Z_{2i}, Z_{3i}, \dots, Z_{ri})^T ; i = 1, 2, \dots, n$$

$$= \sigma^2 (1, Z_{2i}, Z_{3i}, \dots, Z_{ri}) (Z^T Z)^{-1} (1, Z_{2i}, Z_{3i}, \dots, Z_{ri})^T ; i = 1, 2, \dots, n$$

$$\text{โดยที่ } \sigma^2 = \frac{1}{n-r} (Y^T Y - \hat{Y}^T Z^T Y)$$

1 Wilk, SS., "Mathematical Statistics" (John Wiley & Sons Inc, Tokyo 1962) p. 82

ดังนั้นช่วงพยากรณ์ $(1 - \alpha) 100\%$ ของ Y_{Fi} ; $i = 1, 2, \dots, n$ คือ

$$\hat{Y}_{Fi} + t_{n-r, \alpha/2} \sqrt{V(\hat{Y}_{Fi})} < Y_{Fi} < \hat{Y}_{Fi} + t_{n-r, 1-\alpha/2} \sqrt{V(\hat{Y}_{Fi})}$$

หรือ $Y_L < Y_i < Y_U$

โดยที่ $\hat{Y}_{Fi} = \hat{\gamma}_0 + \hat{\gamma}_1 Z_{1i} + \hat{\gamma}_2 Z_{2i} + \dots + \hat{\gamma}_r Z_{ri}$

$$V(\hat{Y}_{Fi}) = \hat{\sigma}^2 (1, Z_{2i}, Z_{3i}, \dots, Z_{ri}) (Z^T Z)^{-1} (1, Z_{2i}, Z_{3i}, \dots, Z_{ri})$$

$$\hat{\sigma}^2 = \frac{1}{n-r} (Y^T Y - \hat{\gamma}^T Z^T Y)$$

ค่าช่วงพยากรณ์ $(1-\alpha) 100\%$ นี้เรียกว่าช่วงของข้อมูลทดแทน (Y_L, Y_U)

ค. 5 การตัดสินใจเกี่ยวกับความถูกต้องน่าเชื่อถือของข้อมูลเดิม
ให้ Y_i คือข้อมูลเดิม และ Y_{Fi} คือข้อมูลทดแทนของ Y_i ที่ได้จากวิธีการ PC ดังนั้น ถ้า

(1) $Y_L < Y_i < Y_U$ 0 0 ให้ถือว่าข้อมูลเดิมถูกต้อง

(2) ถ้า $Y_i > Y_U$ ให้ใช้ $Y_U = \hat{Y}_{Fi} + t_{n-r, 1-\alpha/2} \sqrt{V(\hat{Y}_{Fi})}$ เป็นข้อมูลทดแทน

(3) ถ้า $Y_i < Y_L$ ให้ใช้ $Y_L = \hat{Y}_{Fi} + t_{n-r, \alpha/2} \sqrt{V(\hat{Y}_{Fi})}$ เป็นข้อมูลทดแทน

1.3.3 วิธีบรรณาการเชิงตรรกวิทยา (Logical Editing) ¹

ก. เหตุผลและความจำเป็น

ในกรณีที่ข้อมูลต้องการบรรณาการและแก้ไขเป็นข้อมูลเชิงคุณภาพ (Qualitative Data) เช่นเพศ สีผิว เชื้อชาติ ศาสนา สถานภาพสมรส ความสัมพันธ์กับหัวหน้าครอบครัว ภูมิภาค ฯลฯ การบรรณาการและการสร้างข้อมูลทดแทนย่อมมีวิธีการแตกต่างไปจากวิธีทั้งสองข้างต้น เพราะวิธี LSC และ PC เหมาะสำหรับใช้กับกรณีของข้อมูลเชิงปริมาณ (Quantitative Data or Metric) ทั้งนี้ด้วยเหตุผลประการสำคัญคือข้อมูลทดแทนที่ได้รับจาก

1. P. Fellegi and D. Holt, A Systematic Approach to Automatic Edit and Imputation (Jour. Am. Stat. Asso.) P.17

วิธีดังกล่าวไม่สามารถแปลความหมายกลับเป็นเชิงคุณภาพได้ชัดเจนหรือรัดกุมเพียงพอ ด้วยเหตุดังกล่าวจึงได้มีความพยายามนำเอาหลักเกณฑ์ทางตรรกวิทยามาใช้เป็นแนวทางในการบรรณาธิกรณและสร้างข้อมูลทดแทนสำหรับกรณีของข้อมูลเชิงคุณภาพดังกล่าว วิธีการนี้เรียกว่า การบรรณาธิกรณเชิงตรรกวิทยา (Logical Editing)

หลักการโดยสรุปของการบรรณาธิกรณเชิงตรรกวิทยาประกอบด้วยขั้นตอนดำเนินการกว้าง ๆ 4 ขั้นตอนดังนี้คือ

1. สร้างเงื่อนไขการตรวจสอบเรียกว่าเงื่อนไขขัดแย้ง (fail edit, e) ขึ้นโดยพิจารณาจากระหัสคำตอบ (Code) ที่เป็นไปได้ทั้งหมด ของแบบสำรวจ (Record) แล้วประมวล Combination ของระหัสที่ขัดแย้งกับความเป็นจริงขึ้นมาเป็นเซต
2. นำ Record ที่ต้องการมาตรวจสอบตามเงื่อนไข
3. สร้างข้อมูลทดแทนให้แก่รายการข้อมูล (field) ที่ขัดแย้งกับความเป็นจริง หรือสอดคล้องกับเงื่อนไขขัดแย้ง (fail edit)

ข. ตัวอย่าง

เพื่อความสะดวกต่อการศึกษาวิธีการตลอดจนทฤษฎีของวิธีการเชิงตรรกวิทยา (logical editing) ซึ่งจะกล่าวถึงต่อไป ควรจะได้ทำความเข้าใจในตัวอย่างต่อไปนี้ซึ่งจักเป็นแนวทางให้เข้าใจสัญลักษณ์และนิยามบางประการที่เกี่ยวข้อง

สมมุติแบบสอบถามชุดหนึ่ง Record ประกอบด้วยรายการข้อมูล (field) 3 รายการคือ เพศ อายุและความสัมพันธ์กับหัวหน้าครอบครัว

จากการพิจารณาตามหลักแห่งเหตุผลและความเป็นไปได้ เราสามารถจำแนกระหัสคำตอบของ

แต่ละรายการข้อมูลได้ดังนี้

เพศ	อายุ	ความสัมพันธ์กับหัวหน้าครอบครัว
ชาย	0-14	ภรรยา
หญิง	15 +	สามี
		ธิดา
		ญาติ
		อื่น ๆ

ให้ A_i = เซตของรหัสคำตอบของรายการข้อมูลที่ i ; $i = 1, 2, 3$

n_i = จำนวนรหัสหรือจำนวนสมาชิกของเซต A_i

ดังนั้น

$A_1 = \{\text{ชาย, หญิง}\} ; n_1 = 2$

$A_2 = \{0-14 \text{ ปี, } 15+ \text{ ปี}\} ; n_2 = 2$

$A_3 = \{\text{ภรรยา, สามี, บุตร, ธิดา, ญาติ, อื่น ๆ}\} ; n_3 = 6$

ดังนั้น คำตอบทั้งหมด (all possible code) ของแบบสอบถาม (Record) ชุดนี้ จึงประกอบด้วยคำตอบทั้งหมด $2 \times 2 \times 6 = 24$ คำตอบคือ

$A = \{(\text{ชาย, อายุไม่เกิน 14 ปี, เป็นภรรยา}), (\text{ชาย, อายุตั้งแต่ 15 ปี ขึ้นไป, เป็นภรรยา}), (\text{ชาย, อายุไม่เกิน 14 ปี, เป็นสามี}), (\text{ชาย, อายุตั้งแต่ 15 ปีขึ้นไป, เป็นสามี}), (\text{ชาย, อายุไม่เกิน 14 ปี, เป็นบุตร}), (\text{ชาย, อายุตั้งแต่ 15 ปีขึ้นไป, เป็นธิดา}), (\text{ชาย, อายุไม่เกิน 14 ปี, เป็นญาติ}), (\text{ชาย, อายุตั้งแต่ 15 ปีขึ้นไป, เป็นญาติ}), (\text{ชาย, อายุไม่เกิน 14 ปี, เกี่ยวข้องในฐานะอื่น}), (\text{ชาย, อายุตั้งแต่ 15 ปีขึ้นไป, เกี่ยวข้องในฐานะอื่น}), (\text{หญิง, อายุไม่เกิน 14 ปี, เป็นภรรยา}), (\text{หญิง, อายุตั้งแต่ 15 ปีขึ้นไป, เป็นภรรยา}), (\text{หญิง, อายุไม่เกิน 14 ปี, เป็นสามี}), (\text{หญิง, อายุตั้งแต่ 15 ปีขึ้นไป, เป็นสามี}), (\text{หญิง, อายุไม่เกิน 14 ปี, เป็นบุตร}), (\text{หญิง, อายุตั้งแต่ 15 ปีขึ้นไป, เป็นธิดา}), (\text{หญิง, อายุไม่เกิน 14 ปี, เป็นญาติ}), (\text{หญิง, อายุไม่เกิน 14 ปี, เกี่ยวข้องในฐานะอื่น}), (\text{หญิงอายุตั้งแต่ 15 ปีขึ้นไป เป็นบุตร}), (\text{ชายอายุตั้งแต่ 15 ปีขึ้นไป, เป็นบุตร}), (\text{ชาย, อายุไม่เกิน 14 ปี, เป็นธิดา})\}$

A คือ Cartesian Product ของ field ใน record นั่นคือ

$$A = A_1 \times A_2 \times A_3$$

ให้ A คือเซตของรหัสคำตอบทั้งหมดของ record a ใด ๆ ที่ field i ของ a เป็นอนุเซตของ

สมมติให้ $i = 2$ นั่นคือ a มี field ที่ 2 เป็นอนุเซตของ A_2 สมมติว่ารหัสใน field ที่ 2 ของ a คือ อายุ 15 ปีขึ้นไป นั่นคือ $= \{15 + \text{ปี}\}$

ดังนั้นรหัสคำตอบทั้งหมด (all possible code) ของ record a จึงมีอยู่ทั้งสิ้น

$$2 \times 1 \times 6 = 12 \text{ รหัส หรือ}$$

คือ $\{(\text{ชาย, อายุตั้งแต่ 15 ปีขึ้นไป, เป็นภรรยา}), (\text{ชาย, อายุตั้งแต่ 15 ปีขึ้นไป, เป็น$

สามี), (ชาย, อายุตั้งแต่ 15 ปีขึ้นไป, เป็นบุตร), (ชาย, อายุตั้งแต่ 15 ปีขึ้นไป, เป็นธิดา), (ชาย, อายุตั้งแต่ 15 ปีขึ้นไป, เป็นญาติ), (ชาย, อายุตั้งแต่ 15 ปีขึ้นไป, เกี่ยวข้องในฐานะอื่น ๆ), (หญิง, อายุตั้งแต่ 15 ปีขึ้นไป, เป็นสามี), (หญิง, อายุตั้งแต่ 15 ปีขึ้นไป, เป็นภรรยา), (หญิง, อายุตั้งแต่ 15 ปีขึ้นไป, เป็นบุตร), (หญิง, ตั้งแต่ 15 ปีขึ้นไป, เป็นธิดา), (หญิง, อายุตั้งแต่ 15 ปีขึ้นไป, เป็นญาติ), (หญิง, อายุตั้งแต่ 15 ปีขึ้นไป, เกี่ยวข้องในฐานะอื่น ๆ) }

จึงเห็นได้ว่า ถ้าเรานิยามว่า $A_i^o = A_1 \times A_2 \times \dots \times A_{i-1} \times A_i^o \times \dots \times A_n$

ซึ่งในที่นี้คือ $A_i^o = A_1 \times A_2^o \times A_3$ แล้ว A_i^o จะเป็นอนุเซตของ A

คือ $A_i^o \subseteq A$ และขอให้สังเกตไว้เป็น 5 ประการคือ

1. A_i^o คือเซตของรหัสคำตอบทั้งหมดของ record a เมื่อ a คือ record ใด ๆ ที่มี field ที่ 2 ที่เราให้ความสนใจเป็นการเฉพาะเจาะจง
2. ทั้ง A_i^o และ A (ซึ่งกรณีทั่วไปคือ A_i^o และ A) เป็นเซตของเซต หมายความว่าสมาชิกของ A_i^o และ A คือเซตของคำตอบที่ฟังเป็นไปได้อันหนึ่งของ record a และ record ใด ๆ ตามลำดับ
3. เมื่อสนใจผลลัพธ์หรือคำตอบของเฉพาะ record a จะพบว่าคำตอบของ record a อาจจะเป็นสมาชิกใดสมาชิกหนึ่ง A_i^o (ในที่นี้คือ A_i^o) ก็ได้ แสดงว่า aca_i^o
4. A และ A_i^o จะประกอบไปด้วยทั้งเซตของคำตอบที่น่าเป็นไปได้และไม่อาจเป็นไปได้ เช่น (ชาย, อายุไม่เกิน 14 ปี, เป็นภรรยา) เป็นตัวอย่างของคำตอบที่ไม่อาจเป็นไปได้ เรียกว่าเงื่อนไขขัดแย้ง (fail edit)
5. กลุ่มของ Combination ของรหัสคำตอบที่ไม่อาจเป็นไปได้ (Unacceptable) เป็นอนุเซตของ Code Space

ค. สัญลักษณ์และทฤษฎีที่เกี่ยวข้อง

ค. 1 สัญลักษณ์

A_i = เซตของคำตอบที่พึงเป็นไปได้ทั้งหมดของรายการข้อมูลที่ $i; i = 1, 2, \dots, n$

n = จำนวนรายการข้อมูลทั้งหมด (จำนวน field)

n_i = จำนวนสมาชิกของเซต $A_i; i = 1, 2, \dots, n$ n_i อาจเป็นจำนวนนับได้ (finite) หรือนับไม่ได้ (infinite) ก็ได้

A = Cartesian Product ของทุก field ใน record หรือเซตของรหัสคำตอบที่เป็นไปได้ทั้งหมดของ record

ดังนั้น $A = A_1 \times A_2 \times A_3 \times \dots \times A_n$

a = record ใด ๆ = เวกเตอร์ที่สมาชิกของ a คือรหัสจากทุก field ของ a

A_i^0 = เซตของรหัสทั้งหมดของ record a ใด ๆ ที่เขาสนใจเฉพาะ field i ของ a ที่เป็นอนุเซตของ A_i หรือ

$A_i^0 = A_1 \times A_2 \times A_3 \times \dots \times A_{i-1} \times A_i^0 \times \dots \times A_n$ ($i = 1, 2, \dots, n$)

A_i^r, A_j^r = field ที่ i ของเงื่อนไขที่ r , field ที่ j ของเงื่อนไขที่ r

ดังนั้น เราสามารถยังผลสรุปได้เป็น 2 ประการดังนี้

1. $a \in A_i^0$ และรหัสคำตอบของ field ที่ i ของ a จะเป็น Component ของทุกสมาชิกของ A_i^0

2. $A_i^0 \subseteq A_i$ เมื่อนิยามว่า $A_i = A_1 \times A_2 \times \dots \times A_{i-1} \times A_i \times A_{i+1} \times \dots \times A_n$

และ $A_i \subseteq A$ เรียก A ว่า Code Space

ค.2 แนวคิดในการสร้างเงื่อนไขการบรรณาธิการ (Normal form of Edit)

การตรวจสอบหรือบรรณาธิการข้อมูลโดยนัยแห่งวิธีการบรรณาธิการเชิงตรรกวิทยาก็คือ การตรวจสอบโดยสร้างเงื่อนไขการบรรณาธิการขึ้นมาใช้, เงื่อนไขจะถูก

1. super script "0" หมายความว่า A_i คือ field ที่เรากำลังให้ความสนใจ superscript "i" ชี้ให้เห็นว่าเรากำลังสนใจ field ใดโดยที่ $i = 1, 2, \dots, n$