



บทที่ 2

ความรู้และทฤษฎีพื้นฐานสำหรับการวิจัย

ในการวิจัยเรื่องนี้ สถิติที่ใช้ศึกษาคือ สัมประสิทธิ์สหสัมพันธ์ และการทดสอบแบบ ไคล์แควร์ สำหรับสัมประสิทธิ์สหสัมพันธ์ ตัวแปรที่นำมาศึกษาจะต้องมีการแจกแจงแบบปกติสองตัวแปร สำหรับการทดสอบแบบไคล์แควร์ในการศึกษาตัวแปรแบบไม่ต่อเนื่อง ตัวแปรทั้งสองจะมีการแจกแจงร่วมแบบพหุนามที่มาจากประชากรเดียวกัน ในการจำลองข้อมูลทั้งสองแบบทำได้โดยวิธีการจำลองแบบ ดังนั้น ในบทนี้จะกล่าวถึงรายละเอียดของสถิติแต่ละวิธี พร้อมทั้งความรู้ที่เกี่ยวข้องอย่างละเอียด ดังนี้

2.1 การแจกแจงแบบปกติสองตัวแปร (The bivariate normal distribution)

ถ้า (x, y) เป็นตัวแปรที่มีฟังก์ชันการแจกแจงความน่าจะเป็นร่วม

$$f(x, y) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp \left[-\frac{1}{2(1-\rho^2)} \left\{ \left(\frac{x-\mu_1}{\sigma_1} \right)^2 - 2\rho \left(\frac{x-\mu_1}{\sigma_1} \right) \left(\frac{y-\mu_2}{\sigma_2} \right) + \left(\frac{y-\mu_2}{\sigma_2} \right)^2 \right\} \right]$$

ภายใต้เงื่อนไข

1. $-\infty < x < \infty$, $-\infty < y < \infty$
 2. $\sigma_1, \sigma_2, \mu_1, \mu_2$ และ ρ เป็นค่าพารามิเตอร์คงที่
- โดยที่ ρ คือค่าสัมประสิทธิ์สหสัมพันธ์ระหว่างตัวแปร x, y ($-1 < \rho < 1$)

μ_1 คือค่าเฉลี่ยของตัวแปร x ($-\infty < \mu_1 < \infty$)

μ_2 คือค่าเฉลี่ยของตัวแปร y ($-\infty < \mu_2 < \infty$)

σ_1^2 คือความแปรปรวนของตัวแปร x ($\sigma_1^2 > 0$)

σ_2^2 คือความแปรปรวนของตัวแปร y ($\sigma_2^2 > 0$)

เราจะเรียก (x, y) ว่าเป็นตัวแปรที่มีการแจกแจงแบบปกติสองตัวแปร

ที่มีค่าเฉลี่ย $\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$

และเมทริกซ์ความแปรปรวนร่วม $\Sigma = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$

2.2 การแจกแจงแบบพหุนามร่วม (Joint multinomial distribution)

ในการศึกษาคุณลักษณะทางประชากร 2 ลักษณะที่สนใจคือ A และ B และจากผลการศึกษาคุณลักษณะ A สามารถแบ่งออกได้ r ระดับ คือ A_1, A_2, \dots, A_r และคุณลักษณะ B สามารถแบ่งได้ c ระดับ คือ B_1, B_2, \dots, B_c สุ่มทำการทดลอง n ครั้งอย่างเป็นอิสระต่อกัน ปรากฏผลการทดลองดังนี้คือ A_1 และ B_1 จะเกิดขึ้น x_{11} ครั้ง A_1 และ B_2 เกิดขึ้น x_{12} ครั้ง ... A_i และ B_j จะเกิดขึ้น x_{ij} ครั้ง ด้วยความน่าจะเป็น $P_{11}, P_{12}, \dots, P_{ij}$ ($i = 1, 2, \dots, r$ และ $j = 1, 2, \dots, c$) ตามลำดับ ดังนั้นฟังก์ชันการแจกแจงความน่าจะเป็นร่วมคือ $f(x_{11}, x_{12}, \dots, x_{rc}; P_{11}, P_{12}, \dots, P_{rc}) = \binom{n}{x_{11}, x_{12}, \dots, x_{rc}} \prod_{i,j} P_{ij}^{x_{ij}}$

เมื่อ x_{ij} = จำนวนหน่วยตัวอย่าง (ความถี่) ที่สอดคล้องกับระดับที่ i ของคุณลักษณะ A และสอดคล้องกับระดับที่ j ของคุณลักษณะ B (จำนวนความถี่ที่ตกในเซลล์ (i, j))

P_{ij} = ความน่าจะเป็นที่หน่วยตัวอย่างใด ๆ จะตกในเซลล์ (i, j) หรือนัยหนึ่งคือ ความน่าจะเป็นที่หน่วยตัวอย่างใด ๆ จะมีคุณลักษณะ A ระดับที่ i และคุณลักษณะ B ระดับที่ j

และ $\sum_{i=1}^r \sum_{j=1}^c x_{ij} = n$ (จำนวนหน่วยตัวอย่างทั้งหมด)

2.3 สัมประสิทธิ์สหสัมพันธ์¹ (Correlation coefficient)

ในการศึกษาระดับความสัมพันธ์ระหว่างตัวแปรสองตัว หรือลักษณะที่สนใจศึกษาสองลักษณะ โดยที่ไม่มีตัวแปรใดหรือลักษณะใดถูกกำหนดค่าไว้ล่วงหน้า จะพิจารณาได้จากค่าสัมประสิทธิ์สหสัมพันธ์ (Correlation coefficient) สัมประสิทธิ์สหสัมพันธ์ประชากร สัญลักษณ์ที่ใช้คือ ρ (rho) ค่าประมาณของสัมประสิทธิ์สหสัมพันธ์ ρ คือ สัมประสิทธิ์สหสัมพันธ์ตัวอย่าง (sample correlation) ใช้สัญลักษณ์ r หรือ r_{xy}

¹ในที่นี้จะหมายถึงสัมประสิทธิ์สหสัมพันธ์เชิงเส้นเท่านั้น

สัมประสิทธิ์สหสัมพันธ์ระหว่างสองตัวแปร ที่รู้จักกันดีที่สุดคือ สัมประสิทธิ์สหสัมพันธ์เชิงเส้นแบบเพียร์สัน ซึ่ง คาร์ล เพียร์สัน (Karl Pearson) เป็นผู้เริ่มไปมาตั้งแต่ปี ค.ศ. 1895 บางทีจึงเรียกสัมประสิทธิ์สหสัมพันธ์นี้ว่า สัมประสิทธิ์สหสัมพันธ์ผลคูณโมเมนต์ของเพียร์สัน (Pearson Product-Moment Correlation Coefficient) สูตรที่ใช้คำนวณ คือ

$$r = \frac{n\sum XY - (\sum X)(\sum Y)}{\sqrt{\{n\sum X^2 - (\sum X)^2\} \{n\sum Y^2 - (\sum Y)^2\}}}$$

โดยที่ r คือ สัมประสิทธิ์สหสัมพันธ์ตัวอย่างระหว่างตัวแปร X กับตัวแปร Y

$\sum X$ คือ ผลรวมของค่าข้อมูลของตัวแปร X

$\sum Y$ คือ ผลรวมของค่าข้อมูลของตัวแปร Y

$\sum XY$ คือ ผลรวมของผลคูณระหว่างค่าข้อมูลของตัวแปร X และตัวแปร Y

$\sum X^2$ คือ ผลรวมของกำลังสองของค่าข้อมูลของตัวแปร X

$\sum Y^2$ คือ ผลรวมของกำลังสองของค่าข้อมูลของตัวแปร Y

n คือ จำนวนคู่ของข้อมูล

และค่าของ X และ Y จะต้องเป็นค่าที่วัดในลักษณะควบคู่กัน (bivariate (x, y) data)

ข้อตกลงเบื้องต้น (Assumption)

1. ตัวแปรทั้งสองต้องเป็นค่าต่อเนื่อง และมีการแจกแจงปกติสองตัวแปร

(bivariate normal distribution)

2. ความสัมพันธ์ระหว่างตัวแปรทั้งสองเป็นแบบเส้นตรง (linear relationship)

3. ข้อมูล (x_i, y_i) ($i = 1, 2, \dots, n$) เป็นตัวอย่างเชิงสุ่ม

สัมประสิทธิ์สหสัมพันธ์มีค่าอยู่ระหว่าง -1 ถึง $+1$ ถ้าสัมประสิทธิ์สหสัมพันธ์มีค่าสมบูรณ์

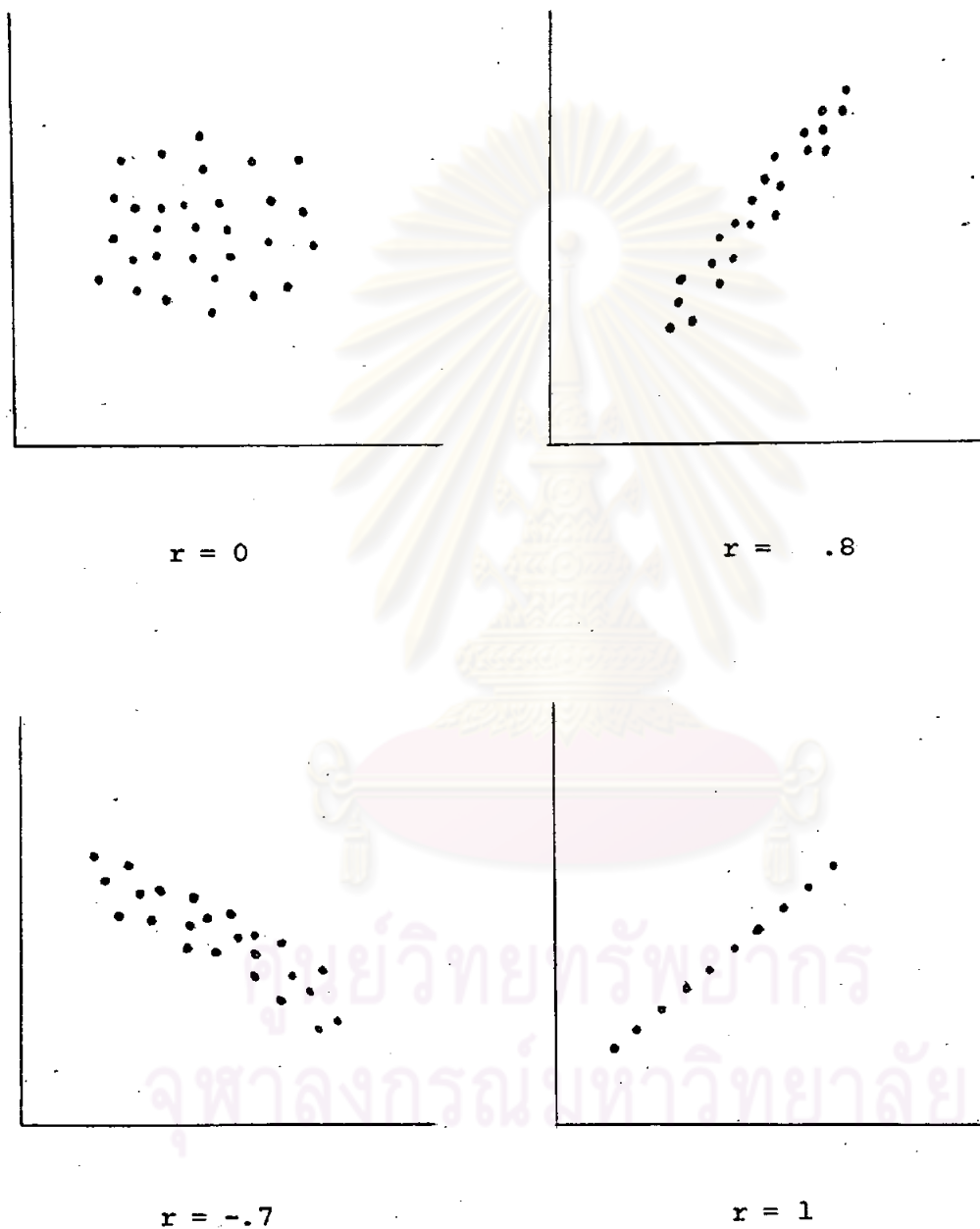
เข้าใกล้ 1 แสดงว่าตัวแปรคู่ที่กำลังศึกษามีความสัมพันธ์เชิงเส้นในระดับสูง ถ้าสัมประสิทธิ์สหสัมพันธ์มีค่าสมบูรณ์ไม่มากนัก (ประมาณ $0.3 - 0.7$) แสดงว่ามีความสัมพันธ์เชิงเส้นในระดับปานกลาง แต่ถ้าสัมประสิทธิ์สหสัมพันธ์มีค่าเข้าใกล้ 0 แสดงว่ามีความสัมพันธ์กันน้อย และถ้าสัมประสิทธิ์สหสัมพันธ์เท่ากับ 0 แสดงว่าตัวแปรทั้งสองไม่มีความสัมพันธ์เชิงเส้น

สำหรับทิศทางของความสัมพันธ์ จะพิจารณาได้จากเครื่องหมายของสัมประสิทธิ์สหสัมพันธ์ ถ้าสัมประสิทธิ์สหสัมพันธ์มีค่าเป็นบวก แสดงว่าเมื่อตัวแปรตัวหนึ่งมีค่าเพิ่มขึ้น ตัวแปรอีกตัวหนึ่งจะมีค่าเพิ่มขึ้นตาม แต่ถ้าสัมประสิทธิ์สหสัมพันธ์มีค่าเป็นลบ แสดงว่าเมื่อตัวแปรตัวหนึ่งมีค่าเพิ่มขึ้น ตัวแปรอีกตัวหนึ่งจะมีค่าลดลง นั่นคือเป็นไปในทางตรงกันข้าม อย่างไรก็ตามผลที่ได้จากการคำนวณค่า r^2 แล้วตีความในรูปของค่าสัมประสิทธิ์สหสัมพันธ์ ยังไม่สามารถบอกทิศทางของความสัมพันธ์ได้ ต้องอาศัยวิธีการอื่น ๆ ประกอบ จึงจะสามารถบ่งชี้ได้



ศูนย์วิทยทรัพยากร
จุฬาลงกรณ์มหาวิทยาลัย

รูปที่ 2.1 แสดง ค่าสัมประสิทธิ์สหสัมพันธ์ในระดับความสัมพันธ์ต่าง ๆ



สัมประสิทธิ์สหสัมพันธ์เป็นตัวเลขไม่มีหน่วย ไม่ขึ้นอยู่กับขนาดมากน้อยของค่าตัวแปร ไม่จำเป็นที่ตัวแปรทั้งสองจะต้องมีหน่วยเดียวกัน และความสัมพันธ์ที่ได้เป็นการ แสดงแต่เพียงว่า การเกิดของสองปรากฏการณ์นั้นมีความสัมพันธ์เชิงบวกหรือเชิงลบเท่านั้น มิได้หมายความว่าอะไร เป็นเหตุหรือเป็นผลซึ่งกันและกันแต่อย่างใด

2.4 การจำลองแบบ (Simulation)

ในการศึกษาปัญหาซึ่งมีข้อจำกัดว่าข้อมูลที่ผ่านมาวิเคราะห์นั้น จะต้องมีความสัมพันธ์ตามข้อสมมติเบื้องต้น เช่น สามารถระบุค่าเฉลี่ย ความแปรปรวน หรือลักษณะการแจกแจงของประชากรได้ไหม วิธีการสร้างข้อมูลที่มีความสัมพันธ์ตามต้องการ กระทำได้โดยวิธีการที่เรียกว่า เทคนิคการจำลองแบบ (Simulation technique) เป็นเทคนิคการทดลองพื้นฐานที่มีความรวดเร็ว และเสียค่าใช้จ่ายถูกกว่าในการทดลองจริง ซึ่งอาจทำได้ยาก หรือไม่สามารทำได้ เพราะมีปัจจัยหรือข้อจำกัดมากมาย

สำหรับการศึกษาในเรื่องนี้ ได้ใช้เทคนิคการจำลองแบบสร้างข้อมูลที่มีการแจกแจงแบบปกติสองตัวแปร และพหุนามแบบสองตัวแปร เครื่องมือที่สำคัญในเทคนิคการจำลองแบบ คือเทคนิคการผลิตเลขสุ่ม (Technique of random number generation) ซึ่งเป็นพื้นฐานในการสร้างลักษณะการแจกแจงแบบต่าง ๆ เช่น การแจกแจงแบบทวินาม การแจกแจงแบบปัวซอง การแจกแจงแบบแกมมา เป็นต้น การจำลองแบบจะใกล้เคียงความเป็นจริงเพียงใดนั้น ขึ้นอยู่กับคุณภาพของตัวเลขสุ่มเป็นสำคัญ ตัวเลขสุ่มที่ได้จะถูกสร้างให้มีลักษณะสุ่ม (random) และให้แต่ละหมายเลขมีโอกาสเกิดขึ้นเท่าเทียมกัน

การสร้างข้อมูลที่มีการแจกแจงแบบปกติสองตัวแปร จะสร้างตัวเลขสุ่มที่มีการแจกแจงแบบสม่ำเสมอ (Uniform distribution) ขึ้นมาก่อน โดยใช้ซับรูทีนแรนดู (Subroutine Randu) จากตัวเลขสุ่มที่มีการแจกแจงแบบสม่ำเสมอในช่วง $(0, 1)$ นี้ จะนำมาสร้างตัวแปรปกติมาตรฐานโดยใช้ซับรูทีนเกาส์ (Subroutine Gauss) เรียกค่าที่ได้ครั้งแรกนี้ว่า x ต่อจากนั้นจะสร้างตัวเลขสุ่มที่มีการแจกแจงแบบสม่ำเสมอขึ้นมาอีกชุดหนึ่ง ซึ่งไม่ขึ้นกับตัวเลขสุ่มที่มีการแจกแจงแบบสม่ำเสมอในชุดแรก โดยใช้ซับรูทีนแรนดูเช่นกัน จากนั้นจะนำตัวเลขสุ่มนี้มาสร้างตัวแปรปกติมาตรฐาน และเรียกค่าที่ได้นี้ว่า y โดยค่า y จะมีลักษณะการแจกแจงที่ขึ้นอยู่กับค่า x ที่ได้จากครั้งแรก ผลลัพธ์ที่ได้คือค่า (x, y) จะเป็นตัวแปรสุ่มปกติสองตัวแปร

การสร้างตัวแปรปกติโดยใช้สัจพจน์เกาส์ อาศัยทฤษฎีลิมิต คือ ทฤษฎีการโน้ม
สู่ค่ากลาง (Central limit theorem)

ทฤษฎี 2.4.1¹ กำหนดให้ y_1, y_2, \dots, y_n ซึ่งเป็นอิสระต่อกัน เป็นตัวอย่างสุ่มขนาด n
จากประชากรที่มีค่าเฉลี่ย μ และค่าความแปรปรวน σ^2 การแจกแจงของค่าสถิติ

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

จะประมาณได้ด้วยการแจกแจงแบบปกติที่มีค่าเฉลี่ย μ และค่าความแปรปรวน $\frac{\sigma^2}{n}$

เนื่องจากเกาส์ ใช้ตัวแปรที่มีการแจกแจงแบบสม่ำเสมอในช่วง $(0, 1)$

12 ตัว ($n = 12$) \bar{y} จะมีค่าเฉลี่ย $\frac{1}{2}$ และค่าความแปรปรวน $\frac{1}{144}$

$$\text{จะได้ } \bar{y} = \frac{1}{12} \sum_{i=1}^{12} y_i \sim N\left(\frac{1}{2}, \frac{1}{144}\right)$$

$$\therefore x = \frac{\bar{y} - \frac{1}{2}}{\sqrt{1/144}}$$

$$= \frac{1}{12} \sum_{i=1}^{12} y_i - 6.0$$

จะมีการแจกแจงโดยประมาณแบบปกติมาตรฐาน

2.5 ตารางการถักร (Contingency table)

การทดสอบความเป็นอิสระ หรือทดสอบความสัมพันธ์ระหว่างตัวแปร 2 ตัว
โดยใช้การทดสอบแบบไคสแควร์ จะต้องนำเสนอบริการข้อมูลให้อยู่ในรูปตารางแจกแจงความถี่ของ
ตัวแปรทั้งสอง (Bivariate frequency table) หรือที่เรียกกันว่า ตารางการถักร
(Contingency table)

ตารางการถักร² เป็นตารางแสดงถึงการแจกแจงความถี่ของข้อมูล ซึ่งจำแนก
ตามตัวแปร 2 ตัวพร้อม ๆ กัน สำหรับตัวแปรแต่ละตัวอาจจำแนกออกเป็น 2, 3, 4... ลักษณะ

¹ ดูพิสัยจาก Robert V. Hoag and Allen T. Craig "Introduction to
Mathematical Statistics หน้า 182-183

² วิบูลย์ พิชาลัย และ สมจิต วัฒนาชยากุล "สถิติสำหรับนักสังคมศาสตร์" พิมพ์ครั้งที่
ที่ 5 กรุงเทพมหานคร สำนักพิมพ์ประกายพรึก, 2527

ฉะนั้น ตารางการแจกแจงอาจเป็นชนิด 2×2 2×3 3×2 3×3 ... ฯลฯ ตารางการแจกแจงที่ได้จากการจำแนกประเภทของหน่วยตัวอย่าง โดยอาศัยค่าของตัวแปร 2 ตัวนี้ จะเรียกได้อีกชื่อหนึ่งว่า ตารางการแจกแจงสองทาง (Two-way contingency table)

ถ้ากำหนดสัญลักษณ์สำหรับตารางการแจกแจงว่า ให้ตารางมีแถวนอน r แถว และแถวตั้ง c แถว นั่นคือจะพิจารณาตารางการแจกแจงชนิด $r \times c$ และให้ O_{ij} เป็นค่าของข้อมูลที่เก็บรวบรวมได้ในแถวนอนที่ i และแถวตั้งที่ j ซึ่งตารางการแจกแจงชนิด $r \times c$ แสดงได้ดังนี้

ตารางที่ 2.1 แสดงตารางการแจกแจงชนิด $r \times c$

		จำแนกแถวตั้ง						
		1	2	...	j	...	c	รวมแถวนอน
จำแนกแถวนอน	1	O_{11}	O_{12}	...	O_{1j}	...	O_{1c}	$O_{1.}$
	2	O_{21}	O_{22}	...	O_{2j}	...	O_{2c}	$O_{2.}$
	.							
	.							
	i	O_{i1}	O_{i2}	...	O_{ij}	...	O_{ic}	$O_{i.}$
.								
.								
r	O_{r1}	O_{r2}	...	O_{rj}	...	O_{rc}	$O_{r.}$	
รวมแถวตั้ง	$O_{.1}$	$O_{.2}$...	$O_{.j}$...	$O_{.c}$	n	

ผลรวมของแต่ละแถวนอนคือ $O_{1.}, O_{2.}, \dots, O_{r.}$

ผลรวมของแต่ละแถวตั้งคือ $O_{.1}, O_{.2}, \dots, O_{.c}$

นั่นคือ $O_{i.} = \sum_{j=1}^c O_{ij}$, $O_{.j} = \sum_{i=1}^r O_{ij}$

n คือจำนวนความถี่ทั้งหมดที่เก็บรวบรวมมา หรือคือขนาดตัวอย่าง

$$\text{โดยที่ } \sum_{i=1}^r \sum_{j=1}^c O_{ij} = \sum_{i=1}^r O_{i.} = \sum_{j=1}^c O_{.j} = n$$

ในการจำแนกข้อมูลตามแถวอนและตามแถวตั้งนี้ ข้อมูลในแต่ละช่องจะไม่ซ้ำซ้อนกัน (mutually exclusive and exhaustive class) ดังนั้นข้อมูลแต่ละค่าจะต้องตกในแถวอนแถวใดแถวหนึ่ง และแถวตั้งแถวใดแถวหนึ่ง

2.6 การทดสอบแบบไคส์แควร์ (Chi-square test)

ในการทดสอบความเป็นอิสระระหว่างตัวแปรสองตัว ในกรณีข้อมูลเชิงคุณภาพหรือแม้แต่ข้อมูลเชิงปริมาณ นักวิจัยมักจะนิยมใช้การทดสอบแบบไคส์แควร์ทดสอบ ซึ่งตัวสถิติของการทดสอบนี้คือ

$$W = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}}, \quad W \sim \chi^2$$

$$\text{โดยที่ } E_{ij} = \frac{O_{i.} \cdot O_{.j}}{n}$$

ตัวสถิตินี้จะประมาณได้ด้วยการแจกแจงแบบไคส์แควร์ ดังใหม่ ในส่วนนี้จะกล่าวถึงการแจกแจงแบบไคส์แควร์ก่อน แล้วจึงจะกล่าวถึงรายละเอียดของการทดสอบแบบไคส์แควร์ดังนี้

2.6.1 การแจกแจงแบบไคส์แควร์ (Chi-square distribution)

ตัวแปรสุ่มไคส์แควร์ คือตัวแปรสุ่มที่พัฒนาขึ้นมาจากฟังก์ชันของตัวแปรสุ่มปกติ กล่าวคือ ถ้า x_1, x_2, \dots, x_n เป็นตัวแปรสุ่มที่มีการแจกแจงปกติโดยมี $E(x_i) = \mu$ และ $V(x_i) = \sigma^2$, $i = 1, 2, \dots, n$ ถ้าสุ่มตัวอย่างขึ้นมา 1 ตัวแล้วแปลงเป็นคะแนนปกติมาตรฐาน Z ดังนี้

$$Z = \frac{x - \mu}{\sigma} \sim N(0, 1)$$

$$\text{ถ้ายกกำลังสอง } Z^2 = \left(\frac{x - \mu}{\sigma}\right)^2 = W, \quad W \sim \chi^2_{(1)}$$

จะได้ว่า Z^2 จะมีการแจกแจงแบบไคส์แควร์ที่มีอันหนึ่งความถี่เป็นอิสระเท่ากับ 1 ในทำนองเดียวกัน ถ้าสุ่มกลุ่มตัวอย่างขนาด n ซึ่งเป็นอิสระต่อกันจะได้ว่าผลรวมของ Z^2 จะมีการแจกแจงแบบไคส์แควร์ที่มีอันหนึ่งความถี่เท่ากับ n ดังนี้

$$W = Z_1^2 + Z_2^2 + \dots + Z_n^2, \quad W \sim \chi^2_{(n)}$$

สำหรับฟังก์ชันความน่าจะเป็นของการแจกแจงไคส์แควร์แสดงได้ดังนี้

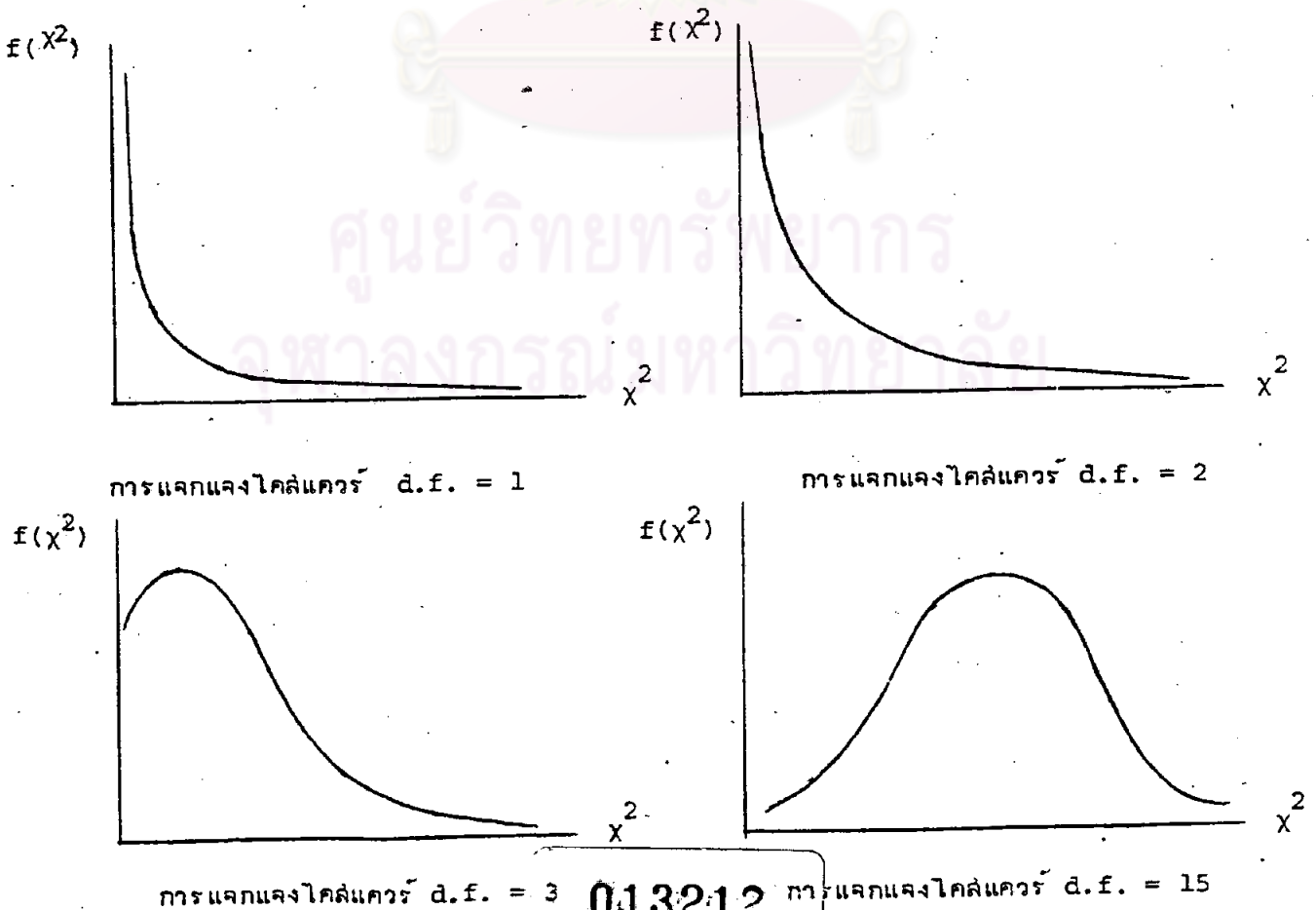
$$f(x) = \frac{1}{2^{n/2} \left(\frac{n}{2} - 1\right)!} x^{\frac{n}{2} - 1} e^{-x/2} \quad \text{เมื่อ } x > 0$$

n เป็นจำนวนขั้นแห่งความเป็นอิสระ

e เป็นจำนวนอตรรกยะ (irrational) มีค่าเท่ากับ 2.71828 (โดยประมาณ)

จากฟังก์ชันความน่าจะเป็นของการแจกแจงไคส์แควร์ จะเห็นว่า การแจกแจงไคส์แควร์ขึ้นอยู่กับจำนวนขั้นแห่งความเป็นอิสระ ส่วนโค้งการแจกแจงจะมีค่าไคส์แควร์เป็นแอบซิสซา (Absissa) และมีค่าฟังก์ชันความน่าจะเป็น ($f(x)$) เป็นออร์ดิเนต (Ordinate) การแจกแจงไคส์แควร์จะเบ้ขวา เมื่อจำนวนขั้นแห่งความเป็นอิสระมีค่าน้อย และความเบ้จะลดลงไปเมื่อจำนวนขั้นแห่งความเป็นอิสระมีค่าเพิ่มขึ้น และถ้าจำนวนขั้นแห่งความเป็นอิสระมีค่ามาก (ประมาณ 20) การแจกแจงค่อนข้างจะลู่มาตรง ซึ่งอาจประมาณการแจกแจงไคส์แควร์ด้วยการแจกแจงปกติได้

รูปที่ 2.2 แสดงการแจกแจงไคส์แควร์เมื่อจำนวนขั้นแห่งความเป็นอิสระแตกต่างกัน



คุณสมบัติของการแจกแจงไคส์แควร์

1. การแจกแจงไคส์แควร์เป็นการแจกแจงแบบต่อเนื่อง
2. ค่าไคส์แควร์เป็นได้ตั้งแต่ 0 ไปจนถึง ∞ นั่นคือค่าไคส์แควร์

เป็นค่าบวกเสมอ

3. มี moment generating function (m.g.f.) = $\left(\frac{1}{1-2t}\right)^{n/2}$

มีค่าเฉลี่ยเท่ากับ n ค่าความแปรปรวนเท่ากับ $2n$ และมี mode ที่ $n-1$ ($n = d.f.$)

4. ลักษณะของส่วนโค้งจะเบ้ไปทางขวามือ
5. เป็นส่วนโค้งแบบเพียร์สันชนิดที่สาม (Pearson Type III Curve)
6. การกระจายเป็นไปตามกฎของการแจกแจงแบบมัลติโนเมียล

(Multinomial Distribution)

ค่าวิกฤติของการแจกแจงไคส์แควร์

เนื่องจากการแจกแจงไคส์แควร์แตกต่างกันไปตามค่าของจำนวนขั้นแห่งความเป็นอิสระ ดังนั้นในการสร้างตารางไคส์แควร์จึงกำหนดค่าวิกฤติ (Critical Value) ในรูปของจำนวนขั้นแห่งความเป็นอิสระ และระดับความสำคัญต่าง ๆ กันตั้งแต่ระดับ 0.99 ลงถึงระดับ 0.001

จำนวนขั้นแห่งความเป็นอิสระ (degrees of freedom)

ในการทดสอบความเป็นอิสระระหว่างสองตัวแปร จะต้องนำเสนอบรรยากาศข้อมูลในรูปของตารางการถักรขณนร $r \times c$ ดังนั้นจำนวนขั้นแห่งความเป็นอิสระจะขึ้นอยู่กับขนาดของตาราง ดังนี้

จำนวนพารามิเตอร์ที่จำเป็นต้องประมาณจากตาราง $= r-1 + c-1 = r+c-2$

สำหรับตารางการถักรขณนร $r \times c$ ข้อมูลที่เก็บรวบรวมมาคือ rc ตัว

$$\begin{aligned} \therefore \text{จำนวนขั้นแห่งความเป็นอิสระ} &= rc - 1 - (r+c-2) \\ &= rc - r - c + 1 \\ &= (r-1)(c-1) \end{aligned}$$

r คือ จำนวนแถวบน (row)

c คือ จำนวนแถวตั้ง (column)

ความสำคัญของ การแจกแจงไคส์แควร์

การแจกแจงไคส์แควร์เป็นการแจกแจงชนิดต่อเนื่อง แต่มักใช้กับข้อมูลชนิดไม่ต่อเนื่อง โดยเฉพาะเกี่ยวกับความถี่ของข้อมูลหรือช่วงของข้อมูล ใช้ได้กับการทดสอบที่ใช้พารามิเตอร์ (parameter) และที่ไม่ใช้พารามิเตอร์ (non-parameter) เช่น ใช้ในการหาช่วงความเชื่อมั่น (confidence interval) ของความแปรปรวนประชากร การทดสอบภาวะล้ารูปสัณฐาน (test of goodness of fit) การทดสอบความเป็นอิสระ (test of independent) เป็นต้น

2.6.2 การทดสอบความเป็นอิสระ (Test of independent)

การทดสอบความเป็นอิสระ เป็นวิธีการทางสถิติที่ใช้ทดสอบว่าตัวแปรสองตัวมีความเป็นอิสระต่อกันหรือไม่ โดยมีสมมติฐาน คือ

H_0 : ตัวแปรทั้งสอง เป็นอิสระต่อกัน

H_a : ตัวแปรทั้งสอง ไม่เป็นอิสระต่อกัน

สถิติทดสอบ

การทดสอบแบบไคส์แควร์ตัวสถิติที่ใช้ในการทดสอบ คือ

$$W = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

W จะมีการแจกแจงแบบไคส์แควร์

เมื่อ O_{ij} = ความถี่ของข้อมูลที่ได้จากการสังเกตหรือทดลองระดับที่ i ของตัวแปรตัวที่ 1 และระดับที่ j ของตัวแปรตัวที่ 2

E_{ij} = ความถี่ที่คาดหวังของระดับที่ i ของตัวแปรตัวที่ 1 และระดับที่ j ของตัวแปรตัวที่ 2 ภายใต้สมมติฐานว่าง (H_0)

โดย
$$E_{ij} = \frac{O_{i.} \times O_{.j}}{n}, \quad i = 1, 2, \dots, r \text{ และ } j = 1, 2, \dots, c$$

$$= \frac{(\text{ผลรวมของ แถวนอนที่ } i)(\text{ผลรวมของ แถวดิ่งที่ } j)}{n}$$

r = จำนวนแถวของแถวของตัวแปรตัวที่ 1

c = จำนวนแถวตั้งของตัวแปรตัวที่ 2

N = จำนวนข้อมูลทั้งหมด

เกณฑ์การตัดสินใจ

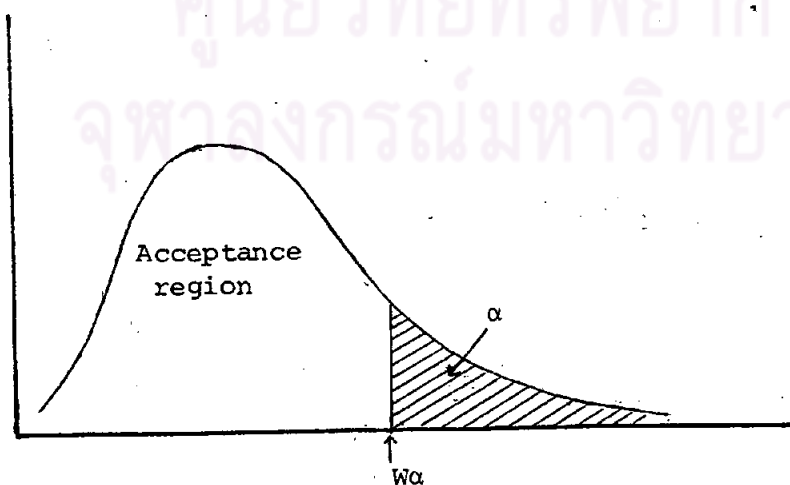
ในการตัดสินใจว่าจะยอมรับ H_0 หรือไม่นั้น อาจทำได้ 2 แบบ คือ

1. ถ้าให้ W ที่คำนวณได้เป็น W_0 เมื่อทราบค่า W_0 และจำนวนชั้นแห่งความเป็นอิสระ (d.f.) แล้ว สามารถอ่านค่าของความน่าจะเป็นได้จากตาราง ถ้าค่าความน่าจะเป็นที่อ่านได้มีค่าน้อยกว่าระดับความมีนัยสำคัญ (α) ที่กำหนด แสดงว่า โอกาสที่สมมติฐานจะเป็นจริงมีน้อยมาก จึงปฏิเสธ H_0 ถ้าความน่าจะเป็นที่อ่านได้จากตารางสูงกว่า α แสดงว่า ความแตกต่างที่เกิดขึ้นไม่มีนัยสำคัญ จึงยอมรับ H_0 ได้

2. อีกแบบคือ เมื่อกำหนดค่า α และทราบค่า d.f. แล้วก็อ่านค่า W_α จากตาราง ถ้า $W_0 \geq W_\alpha$ ความแตกต่างมีนัยสำคัญ นั่นคือ จะปฏิเสธ H_0 แต่ถ้า $W_0 < W_\alpha$ ความแตกต่างไม่มีนัยสำคัญ จะยอมรับ H_0

การทดสอบไคสแควร์เป็นการทดสอบแบบทางเดียวเสมอ ทางด้านมากกว่า ทั้งนี้เพราะว่า O_{ij} จะต่างกับ E_{ij} อย่างมีนัยสำคัญ ถ้า $(O_{ij} - E_{ij})^2$ มีค่ามาก ซึ่งค่าที่ได้จะมีเฉพาะค่าบวกอย่างเดียว

รูปที่ 2.3 แสดงบริเวณปฏิเสธและยอมรับสมมติฐาน



การแจกแจงไคส์แควร์ เป็นการแจกแจงของตัวแปรสุ่มชนิดต่อเนื่อง แต่ความถี่ที่เก็บรวบรวมได้เป็นตัวแปรสุ่มชนิดไม่ต่อเนื่อง อย่างไรก็ตาม การแจกแจงไคส์แควร์ ยังคงเป็นการแจกแจงโดยประมาณที่ดีของค่า χ^2 ที่คำนวณได้เมื่อจำนวนชั้นแห่งความเป็นอิสระมากกว่า 1 และค่าคาดหวังมีค่ามากพอ แต่ในกรณีที่จำนวนชั้นแห่งความเป็นอิสระเท่ากับ 1 นั่นคือ ตารางการถัวชนิด 2×2 จำเป็นต้องมีการแก้ไขโดยใช้ Yates' correction for continuity ตั้งชื่อเพื่อเป็นเกียรติแก่ Dr. Frank Yates ผู้คิด มีสูตรดังนี้

$$\chi^2 (\text{yates-corrrection}) = \sum_{i=1}^2 \sum_{j=1}^2 \frac{(|O_{ij} - E_{ij}| - 0.5)^2}{E_{ij}}$$

โดยทั่วไป เมื่อขนาดตัวอย่างมากกว่าหรือเท่ากับ 50 ไม่จำเป็นต้องปรับค่า χ^2

ต้องปรับค่า χ^2



ศูนย์บริการสุขภาพ
จุฬาลงกรณ์มหาวิทยาลัย