



บทที่ 1

บทนำ

1.1 ความสำคัญและความเป็นมาของปัญหา

ในการศึกษาว่าตัวแปรอิสระ (Independent Variables) มีผลทำให้ตัวแปรตามมีการเปลี่ยนแปลงอย่างไรนั้น อาจารย์วิธีการทางสถิติมาใช้ ซึ่งวิธีหนึ่งที่ใช้ในการวิเคราะห์ความสำคัญของตัวแปรอิสระ คือ การวิเคราะห์ความถดถอยพหุ (Multiple regression analysis) โดยถือว่าการใช้ตัวแปรอิสระที่เหมาะสมมากกว่าหนึ่งตัว จะทำให้ผลของการประมาณค่าตัวแปรตามมีความถูกต้องมากกว่าการใช้ตัวแปรอิสระเพียงตัวเดียว แต่อย่างไรก็ตามจำเป็นต้องคำนึงถึงกรณีที่ตัวแปรอิสระมีพหุสัมพันธ์ (Multicollinearity) กัน ซึ่งหมายถึงการที่ตัวแปรอิสระมีความสัมพันธ์กันเอง ปัญหาดังกล่าวข้างต้นจะส่งผลกระทบต่อค่าสัมประสิทธิ์ความถดถอยพหุ ( $\beta$ )

สำหรับตัวแบบทั่วไป (General Model) ที่แสดงความสัมพันธ์ระหว่างตัวแปรอิสระและตัวแปรตามแบบเชิงเส้น (Linear Relationship) จะมีรูปแบบดังนี้

$$\tilde{y} = \tilde{x}\beta + \tilde{\epsilon} \tag{1}$$

- เมื่อ  $\tilde{y}$  คือ เมตริกซ์ของตัวแปรตามขนาด  $n \times 1$
- $\tilde{x}$  คือ เมตริกซ์ของตัวแปรอิสระขนาด  $n \times p$
- $\beta$  คือ เมตริกซ์ของสัมประสิทธิ์ความถดถอยพหุขนาด  $p \times 1$
- $\tilde{\epsilon}$  คือ เมตริกซ์ของความคลาดเคลื่อนขนาด  $n \times 1$
- $n$  คือ ขนาดตัวอย่าง
- $p$  คือ จำนวนตัวแปรอิสระ

ในการประมาณค่าสัมประสิทธิ์ความถดถอยพหุ โดยทั่วไปจะใช้วิธีกำลังสองน้อยที่สุด (Least squares method) ซึ่งเป็นวิธีที่นิยมใช้กันมาก ตัวประมาณ ( $\beta$ ) ที่ได้เท่ากับ  $(xx')^{-1}xy'$  โดยตัวประมาณนี้จะเป็นตัวประมาณที่ไม่เอนเอียงและให้ค่าเฉลี่ยความคลาดเคลื่อนกำลังสอง (Mean square error) น้อยที่สุด ในบรรดาตัวประมาณที่ไม่เอนเอียงทั้งหลาย แต่ในการประมาณค่าสัมประสิทธิ์ความถดถอยพหุด้วยวิธีกำลังสองน้อยที่สุดนั้น ตัวแปรอิสระจะต้องไม่มีความสัมพันธ์ในลักษณะเชิงเส้น ซึ่งในทางปฏิบัติเป็นไปได้ไม่น้อยมาก เนื่องจากตัวแปรต่าง ๆ ที่นำมาศึกษาอาจมีความสัมพันธ์กัน ตัวแปรอิสระบางตัวอาจเป็นฟังก์ชันของตัวแปรอิสระตัวอื่น ๆ กรณีเช่นนี้กล่าวได้ว่าตัวแปรสัมพันธ์กัน จะมีผลทำให้ตัวประมาณกำลังสองน้อยที่สุดที่ได้เอนเอียง และค่าเฉลี่ยความคลาดเคลื่อนกำลังสองไม่เป็นค่าที่ต่ำที่สุด นั่นคือ ค่าประมาณสัมประสิทธิ์การถดถอยพหุที่ได้ขาดความเที่ยงตรง (Accuracy) ดังนั้นวิธีการกำลังสองน้อยที่สุดจะแก้ไขปัญหานี้โดยการขจัดตัวแปรอิสระที่มีความสัมพันธ์กันออกไป แต่ในบางครั้งการตัดสินใจว่าตัวแปรอิสระตัวใดควรถูกคัดออกป็นนั้นทำได้ยาก เนื่องจากลักษณะความสัมพันธ์ระหว่างตัวแปรอิสระที่เกิดขึ้นไม่ชัดเจนพอ และถือว่าตัวแปรอิสระทุกตัวต่างมีผลต่อการเปลี่ยนแปลงของตัวแปรตามมากพอสมควรซึ่งจากเหตุผลดังกล่าววิธีกำลังสองน้อยที่สุด จึงไม่เหมาะสมที่จะนำมาใช้ในกรณีที่เกิดความสัมพันธ์ระหว่างตัวแปรอิสระ

จากปัญหาข้างต้น นักสถิติได้ทำการศึกษาวิธีการประมาณค่าสัมประสิทธิ์ความถดถอยที่ให้ค่าเฉลี่ยความคลาดเคลื่อนของผลต่างกำลังสองต่ำกว่าวิธีกำลังสองน้อยที่สุด เมื่อเกิดความสัมพันธ์ระหว่างตัวแปรอิสระ โดยให้ชื่อว่า ริดจ์รีเกรสชัน (Ridge regression) ซึ่งทำการศึกษาโดย Hoerl and Kennard (1970) วิธีการนี้อาศัยหลักการที่ว่า ถ้า  $|xx'|$  มีค่าเข้าใกล้ศูนย์ซึ่งทำให้มีปัญหาในการคำนวณ  $(xx')^{-1}$  และผลบวกกำลังสองของ  $\beta$  ซึ่งเท่ากับ  $\beta\beta'$  มีค่ามากกว่าความเป็นจริง ดังนั้นการที่จะทำให้ผลบวกกำลังสองของ  $\beta$  มีค่าลดลง โดยการทำให้  $(xx')^{-1}$  มีค่าเพิ่มขึ้นโดยการบวกค่าคงที่  $k$  ที่มากกว่าศูนย์เข้ากับสมาชิกทุกตัวบนเส้นทแยงมุมของ  $(xx')$  ซึ่งมีผลทำให้ค่าเรคเตอร์ริสติก (Characteristic Root) ของ  $(xx')$  มีค่ามากขึ้นและผลบวกกำลังสองของ  $\beta$  มีค่าลดลง ตัวประมาณค่าสัมประสิทธิ์ความถดถอยพหุ โดยวิธีริดจ์รีเกรสชันเป็นได้ดังนี้

$$\hat{\beta}_k = (xx' + kI)^{-1}xy' \quad ; k > 0 \quad (2)$$

สมการ (2) ใช้หลักการคำนวณเหมือนกับวิธีกำลังสองน้อยที่สุด แต่จะแตกต่างกันที่เมทริกซ์ของตัวแปรอิสระ  $(xx')$  Hoerl and Kennard ได้กล่าวว่าตัวประมาณค่าสัมประสิทธิ์ความถดถอยพหุที่ได้จากวิธีริดจ์รีเกรสชันนี้ จะมีลักษณะค่อนข้างคงที่ ค่าสัมบูรณ์ของตัวประมาณมีค่าสมเหตุสมผลและเครื่องหมายของค่าประมาณสัมประสิทธิ์ความถดถอยพหุจะถูกต้อง



นอกจากวิธีวิธีรีเกรสชันนี้แล้วยังมีอีกวิธีการหนึ่งที่น่าสนใจ คือ วิธีลาเท้นรูทรีเกรสชัน (Latent root regression) ซึ่งคิดค้นโดย Gunst , Webster และ Mason (1978) วิธีการนี้จะกำหนดให้เมตริกซ์  $A = [y : x]$  และทำการคำนวณค่าค่าแตรเคเตอร์วิสดิกรูท หรือ ค่าลาเท้นรูท (Latent root) จากเมตริกซ์  $AA'$  นำค่าลาเท้นรูท และ ลาเท้นเวกเตอร์ (Latent vectors) ที่ได้มาเปรียบเทียบกับเกณฑ์ที่กำหนดไว้ ถ้าค่าลาเท้นรูท และ ลาเท้นเวกเตอร์ มีค่าต่ำกว่าเกณฑ์จะถูกคัดออก ซึ่งมีผลทำให้ตัวประมาณสัมประสิทธิ์ความถดถอยที่มีความถูกต้องและใกล้เคียงกว่าตัวประมาณสัมประสิทธิ์ความถดถอยที่ได้จากวิธีกำลังสองน้อยที่สุด

วิธีการคำนวณค่าสัมประสิทธิ์ความถดถอยทั้งสองวิธีที่กล่าวมานี้ไม่จำเป็นต้องคัดตัวแปรอิสระออกจากตัวแบบ แม้ว่าตัวแปรอิสระนั้นจะมีพหุสัมพันธ์กันก็ตาม แต่ตัวประมาณที่ได้จากทั้ง 2 วิธี จะเป็นตัวประมาณที่เอนเอียง

ถ้าค่าสังเกตสุ่มมาจากประชากรที่มีการแจกแจงแบบเบ้ ตัวประมาณวิธีรีเกรสชันและวิธีลาเท้นรูทรีเกรสชันอาจแตกต่างไปจากกรณีที่ค่าสังเกตของประชากรมีการแจกแจงแบบปกติ เนื่องจากทั้งสองวิธีการเป็นตัวประมาณที่ถูกดัดแปลงมาจากตัวประมาณกำลังสองน้อยที่สุด และวิธีการกำลังสองน้อยที่สุดเป็นวิธีการที่ไวต่อข้อมูลที่ผิดปกติ ดังนั้นจึงเป็นสิ่งที่น่าสนใจที่จะศึกษาว่าวิธีการรีเกรสชันและวิธีลาเท้นรูทรีเกรสชัน วิธีใดให้ค่าเฉลี่ยความคลาดเคลื่อนกำลังสองของตัวประมาณสัมประสิทธิ์การถดถอยพหุคูณน้อยที่สุด

## 1.2 วัตถุประสงค์ของการวิจัย

1. ศึกษาและเปรียบเทียบ ตัวประมาณกำลังสองน้อยที่สุด ตัวประมาณวิธีรีเกรสชันและตัวประมาณลาเท้นรูทรีเกรสชัน เมื่อความคลาดเคลื่อนมีการแจกแจงแบบปกติ การแจกแจงแบบปกติปลอมปน และการแจกแจงแบบลอกนอร์มอล

2. เพื่อเป็นแนวทางในการเลือกใช้ตัวประมาณสัมประสิทธิ์ความถดถอยพหุ เมื่อข้อมูลเกิดพหุสัมพันธ์ระหว่างตัวแปรอิสระ

### 1.3 สมมติฐานของการวิจัย

โดยทั่วไป ตัวประมาณค่าสัมประสิทธิ์ความถดถอยพหุ โดยวิธีลาเท็นรุทซ์ จะให้ค่าตัวประมาณได้ถูกต้องและแม่นยำกว่า วิธีรีดจ์รีเกรสชัน

### 1.4 ข้อตกลงเบื้องต้น

ในการวิจัยครั้งนี้เกณฑ์ที่ใช้ในการเปรียบเทียบตัวประมาณ คือ ค่าเฉลี่ยความคลาดเคลื่อนกำลังสอง (Mean Square Error) ของการประมาณสัมประสิทธิ์การถดถอยพหุด้วยวิธีกำลังสองน้อยที่สุด วิธีการรีดจ์รีเกรสชันและวิธีการลาเท็นรุทซ์รีเกรสชัน

### 1.5 ขอบเขตของการวิจัย

1.5.1 ทำการศึกษาในกรณีที่รูปแบบของการแจกแจงของความคลาดเคลื่อนมีลักษณะการแจกแจง ดังนี้

1.5.1.1 เมื่อความคลาดเคลื่อนมีการแจกแจงแบบปกติ (Normal distribution) ฟังก์ชันความหนาแน่นอยู่ในรูปของ

$$f(x) = \begin{cases} \frac{1}{\sigma\sqrt{2\pi}} \exp \frac{-(x-\mu)^2}{2\sigma^2} & ; x > 0 \\ 0 & ; \text{อื่น ๆ} \end{cases}$$

ในการวิจัยครั้งนี้ จะศึกษาเมื่อค่าเฉลี่ย  $\mu = 0$  และ  $\sigma^2 = 1$

1.5.1.2 เมื่อความคลาดเคลื่อนมีการแจกแจงแบบปกติปลอมปน (Scale - contaminated distribution) ฟังก์ชันการแจกแจงอยู่ในรูปของ

$$F(x) = (1-p)N(0, \sigma^a) + pN(0, c^a \sigma^a)$$



เมื่อ  $c$  คือ สเกลแฟคเตอร์ (scale factor) โดยที่สเกลมีค่าสูงจะทำให้เกิดค่าสังเกตที่ผิดปกติมีค่าสูงด้วย ในการวิจัยครั้งนี้จะศึกษาที่ค่า  $c = 3$  และ  $10$

เมื่อ  $p$  คือ เปอร์เซ็นต์การปลอมปน (percent of contamination) ในการวิจัยครั้งนี้จะศึกษาที่ค่า  $p = 5$  และ  $10$  %

1.5.1.2 เมื่อความคลาดเคลื่อนมีการแจกแจงแบบลอกลอนอร์มอล (Lognormal distribution) ฟังก์ชันความหนาแน่นอยู่ในรูปของ

$$f(x) = \begin{cases} \frac{1}{x\sigma\sqrt{2\pi}} \exp \left[ -\frac{(\ln(x-\mu))^2}{2\sigma^2} \right] ; x > 0, \sigma > 0 \\ 0 & ; \text{อื่น ๆ} \end{cases}$$

ในการวิจัยครั้งนี้ จะศึกษาเมื่อค่าเฉลี่ย  $\mu = 0$  และความแปรปรวน เท่ากับ 1

1.5.2 ขนาดตัวอย่างที่ใช้ในการศึกษา กำหนดให้มีขนาดดังนี้ 10, 30 และ 50

1.5.3 จำนวนตัวแปรอิสระ ซึ่งนำมาพิจารณาในการวิจัย เท่ากับ 3 และ 5

1.5.4 กำหนดค่าสหสัมพันธ์ระหว่างตัวแปรอิสระ เพื่อบังคับให้ตัวแปรอิสระเกิดพหุสัมพันธ์กัน ในระดับที่ต้องการ โดยจัดแบ่งระดับของการเกิดพหุสัมพันธ์เป็น 9 ช่วงดังนี้ [0.11-0.20] , [0.21-0.30] , [0.31-0.40] , [0.41-0.50] , [0.51-0.60] , [0.61-0.70] , [0.71-0.80] , [0.81-0.90] และ [0.91-1.00] ทั้งนี้ เพื่อให้ได้ผลลัพธ์ที่มีความละเอียดและมีความถูกต้องมากขึ้น อีกทั้งจะช่วยให้สามารถเห็นแนวโน้มของค่าเฉลี่ยความคลาดเคลื่อนกำลังสองของตัวแปรตามแต่ละวิธี เมื่อระดับความสัมพันธ์เปลี่ยนแปลงไป

1.5.5 ในการวิจัยครั้งนี้ จะสร้างแบบจำลองข้อมูลให้มีสถานการณ์ตามต้องการโดยอาศัยเครื่องคอมพิวเตอร์เขียนด้วยโปรแกรม Fortran 77 และจำลองข้อมูลโดยอาศัยหลักการของ Multinormal ซึ่งจะได้ค่าตัวแปรอิสระที่มีความสัมพันธ์กันในระดับต่าง ๆ ทั้งนี้จะทำการทดลอง 1,000 รอบ ในแต่ละสถานการณ์



1.5.6 ในการคำนวณค่าเฉลี่ยความคลาดเคลื่อนกำลังสองและตัวประมาณลาเท็นรุธวีเกรสซัน เกณฑ์ที่ใช้ในการตัดค่าลาเท็นรุธ และลาเท็นเวกเตอร์ เพื่อคำนวณค่าดังกล่าว คือ เมื่อค่าลาเท็นรุธ มากกว่าหรือเท่ากับ 0.05 และค่าลาเท็นเวกเตอร์ มากกว่าหรือเท่ากับ 0.10

#### 1.6 ประโยชน์ที่คาดว่าจะได้รับ

เพื่อเป็นแนวทางในการเลือกใช้ตัวประมาณสัมประสิทธิ์ความถดถอยที่เหมาะสมในกรณีเกิดพหุสัมพันธ์ ภายใต้การแจกแจงของความคลาดเคลื่อนแบบปกติหรือแบบเบ้ โดยที่ขนาดตัวอย่าง จำนวนตัวแปรอิสระ และค่าสหสัมพันธ์มีค่าต่าง ๆ กัน



ศูนย์วิทยทรัพยากร  
จุฬาลงกรณ์มหาวิทยาลัย