

บทที่ 1

บทนำ

1. ความสำคัญและที่มาของปัญหา

ในการประมาณค่าตัวแปรที่ไม่ทราบค่าที่นิยมใช้คือการวิเคราะห์จากตัวแปรอื่นที่มีความสัมพันธ์กับตัวแปรที่ต้องการประมาณค่า เพื่อสร้างสมการความสัมพันธ์ระหว่างตัวแปรที่มีข้อมูลกับตัวแปรที่ต้องการประมาณค่าโดยให้เป็นตัวแปรอิสระ (independent variable) และตัวแปรตาม (dependent variable) ตามลำดับ ในทางปฏิบัติการประมาณค่าตัวแปรตามให้มีประสิทธิภาพสูงโดยใช้ตัวแปรอิสระที่เหมาะสมเพียงตัวเดียวอาจไม่เหมาะสมเนื่องจากอธิบายตัวแปรตามได้ไม่มากพอ ผู้วิจัยจึงควรใช้ตัวแปรอิสระตั้งแต่ 2 ตัวขึ้นไปเพื่อนำไปประมาณค่าของตัวแปรตามให้มีความถูกต้องมากขึ้น ซึ่งเราเรียกวิธีการนี้ว่า การวิเคราะห์ความถดถอยพหุคูณ (multiple regression analysis) เราสามารถเขียนตัวแบบทั่วไป (general model) ของความสัมพันธ์เชิงเส้นระหว่างตัวแปรอิสระและตัวแปรตามได้ในรูปแบบข้างดังนี้

$$\tilde{y} = X\tilde{\beta} + \tilde{\epsilon}$$

เมื่อ \tilde{y} คือ เวกเตอร์ของตัวแปรตามขนาด $n \times 1$

X คือ เมทริกซ์ของตัวแปรอิสระขนาด $n \times (p+1)$ และ $X'X$ มีค่าลำดับชั้นเท่ากับ $p+1$

$\tilde{\beta}$ คือ เวกเตอร์ของสัมประสิทธิ์การถดถอยขนาด $(p+1) \times 1$

$\tilde{\epsilon}$ คือ เวกเตอร์ของความคลาดเคลื่อนขนาด $n \times 1$ โดยที่ $E(\tilde{\epsilon}) = 0$ และ $Cov(\tilde{\epsilon}) = \sigma^2 I_n$

n เป็นขนาดตัวอย่างของตัวแปรแต่ละตัว

และ p เป็นจำนวนตัวแปรอิสระ

การประมาณค่าสัมประสิทธิ์การถดถอยพหุคูณจากตัวแบบดังกล่าวนี้ วิธีที่นิยมใช้มากที่สุดคือวิธีกำลังสองน้อยสุด (least square method) ซึ่งค่าประมาณสัมประสิทธิ์การถดถอยพหุคูณอยู่ในรูปของ

$$(1.1) \quad \hat{\tilde{\beta}} = (X'X)^{-1}(X'y)$$

ตัวประมาณในสมการที่ (1.1) มีคุณสมบัติเป็นตัวประมาณไม่เอนเอียงและให้ค่าเฉลี่ยความคลาดเคลื่อนกำลังสอง (mean square error) ต่ำสุดในบรรดาตัวประมาณไม่เอนเอียงเชิงเส้น แต่การประมาณค่าสัมประสิทธิ์การถดถอยพหุคูณด้วยวิธีกำลังสองน้อยสุดมีสมมุติฐานที่สำคัญข้อหนึ่งคือตัวแปรอิสระแต่ละตัวต้องไม่มีความสัมพันธ์กับตัวแปรอิสระตัวอื่นซึ่งในทางปฏิบัติเกิดขึ้นได้น้อยมาก เมื่อตัวแปรอิสระมีพหุสัมพันธ์กันสูงจะทำให้เมทริกซ์ $X'X$ เกิดเงื่อนไขไม่ดี (ill-condition) คือทำให้ $|X'X|$ มีค่าถึงลงเข้าใกล้ศูนย์ เนื่องจากเมทริกซ์ความแปรปรวนร่วมของค่าประมาณสัมประสิทธิ์การถดถอยอยู่ในรูปของ $\hat{Cov}(\beta) = \sigma^2 (X'X)^{-1}$ จึงส่งผลให้ความแปรปรวนของค่าประมาณสัมประสิทธิ์การถดถอยมีค่ามากและเกิดความสัมพันธ์กันสูงระหว่างสัมประสิทธิ์การถดถอยที่ใช้ประมาณค่า ดังนั้นถ้าตัวแปรอิสระมีความสัมพันธ์กันสูงเราอาจแก้ไขโดยตัดตัวแปรอิสระบางตัวออกจากตัวแบบ แต่บางกรณีความสัมพันธ์ระหว่างตัวแปรอิสระไม่ชัดเจนเราจึงไม่ตัดตัวแปรอิสระตัวใดตัวหนึ่งออกได้เพราะถือว่าตัวแปรอิสระทุกตัวมีผลในการอธิบายตัวแปรตามได้พอสมควร

ในปี ค.ศ. 1970 โฮเอิร์ต (Hoerl) และเคนนาร์ด์ (Kennard) ได้เสนอวิธีประมาณค่าสัมประสิทธิ์การถดถอยพหุคูณเมื่อเกิดพหุสัมพันธ์ระหว่างตัวแปรอิสระ โดยใช้หลักการนำค่าคงที่ที่เหมาะสมค่าหนึ่งมาบวกกับสมาชิกในแนวทแยงมุมของเมทริกซ์ $X'X$ เพื่อลดค่าเฉลี่ยความคลาดเคลื่อนกำลังสองให้ต่ำกว่าวิธีกำลังสองน้อยสุด โดยเราเรียกวิธีดังกล่าวว่าวิธีรีดจ์รีเกรสชัน (ridge regression method) ซึ่งค่าประมาณสัมประสิทธิ์การถดถอยพหุคูณอยู่ในรูปแบบดังนี้

$$(1.2) \quad \hat{\beta}_R(k) = (X'X + kI)^{-1} X'y ; k > 0$$

ในปี ค.ศ. 1976 สวินเดล (Swindell) ได้มีแนวความคิดว่าการนำข้อมูลจากอดีตหรือข้อมูลที่มีประโยชน์มาเป็นส่วนประกอบเสริมวิธีรีดจ์รีเกรสชันจะทำให้ค่าประมาณสัมประสิทธิ์การถดถอยพหุคูณมีประสิทธิภาพและความเที่ยงตรงมากขึ้น นอกจากนี้ตัวประมาณสัมประสิทธิ์การถดถอยพหุคูณที่ได้มีคุณสมบัติเป็นตัวประมาณไม่เอนเอียงและให้ค่าเฉลี่ยความคลาดเคลื่อนกำลังสองต่ำกว่าวิธีกำลังสองน้อยสุดและวิธีรีดจ์รีเกรสชัน โดยเราเรียกวิธีดังกล่าวว่าวิธีรีดจ์รีเกรสชันที่ใช้ข้อมูลพิเศษโดยหลักเกณฑ์ (ridge regression with prior information method) ซึ่งค่าประมาณสัมประสิทธิ์การถดถอยพหุคูณอยู่ในรูปแบบดังนี้

$$(1.3) \quad \hat{\beta}_R(kI, j) = (X'X + kI)^{-1} (X'y + kJ) ; k > 0$$

เมื่อ j เป็นเวกเตอร์ขนาด $p \times 1$ จากการประมาณค่า β ในอดีต

จากสมการที่ (1.2) และ (1.3) เราจะเห็นได้ว่า $\hat{\beta}_r(k, j) = \hat{\beta}_r(k)$ เมื่อ $j = 0$ ดังนั้นเราจึงกล่าวได้ว่าวิธีรีดจ์รีเกรสชันเป็นกรณีเฉพาะของวิธีรีดจ์รีเกรสชันที่ใช้ข้อสนเทศโดยหลักเกณฑ์

ในปี ค.ศ.1993 ลิว คีเจียน (Liu Kejian) ได้เสนอวิธีประมาณค่าสัมประสิทธิ์การถดถอยพหุคูณในกรณีที่ตัวแปรอิสระมีพหุสัมพันธ์กัน ซึ่งค่าประมาณสัมประสิทธิ์การถดถอยพหุคูณวิธีของลิว คีเจียนอยู่ในรูปแบบดังนี้

$$(1.4) \quad \hat{\beta}_L(d) = (X'X + I)^{-1}(X'y + d\hat{\beta}) \quad ; \quad 0 < d < 1$$

เมื่อ $\hat{\beta}$ เป็นเวกเตอร์ของสัมประสิทธิ์การถดถอยพหุคูณวิธีกำลังสองน้อยสุด

ตัวประมาณสัมประสิทธิ์การถดถอยพหุคูณวิธีของลิว คีเจียนได้นำข้อดีวิธีรีดจ์รีเกรสชันและวิธีของสไตน์ (Stein) มาผสมผสานกัน กล่าวคือการประมาณค่าสัมประสิทธิ์การถดถอยพหุคูณด้วยวิธีรีดจ์รีเกรสชัน (สมการที่ (1.2)) มีประสิทธิภาพในทางปฏิบัติแต่ประสบความยุ่งยากในการคำนวณหาค่า k ที่เหมาะสม ส่วนค่าประมาณสัมประสิทธิ์การถดถอยพหุคูณวิธีของสไตน์ที่อยู่ในรูปของ $\hat{\beta}_s = c\hat{\beta}$; $0 < c < 1$ เป็นฟังก์ชันเชิงเส้นของค่าคงที่ c ดังนั้นการคำนวณหาค่า c ที่เหมาะสมจึงไม่ยุ่งยากมากนัก แต่อัตราส่วนการลดลงของค่าสัมประสิทธิ์การถดถอยแต่ละตัวเท่ากันซึ่งไม่เหมาะสมในทางปฏิบัติ ดังนั้นข้อดีของการประมาณค่าสัมประสิทธิ์การถดถอยพหุคูณด้วยวิธีของลิว คีเจียนคือตัวประมาณสัมประสิทธิ์การถดถอยพหุคูณวิธีของลิว คีเจียนเป็นฟังก์ชันเชิงเส้นของ d ทำให้การคำนวณหาค่า d ที่เหมาะสมสะดวกกว่าการคำนวณหาค่า k ที่เหมาะสมจากวิธีรีดจ์รีเกรสชัน นอกจากนี้ตัวประมาณสัมประสิทธิ์การถดถอยพหุคูณวิธีของลิว คีเจียนมีประสิทธิภาพในทางปฏิบัติเมื่อตัวแปรอิสระมีพหุสัมพันธ์กัน

ผู้อ่านจะเห็นได้ว่าวิธีของลิว คีเจียน (สมการที่ (1.4)) จะนำค่าคงที่ที่เท่ากันค่าหนึ่งคูณกับ $\hat{\beta}$ ซึ่งการประมาณค่าสัมประสิทธิ์การถดถอยพหุคูณอาจให้ผลไม่ดีเมื่อเทียบกับการนำค่าคงที่ที่เหมาะสมคูณกับแต่ละสมาชิกของ $\hat{\beta}$ โดยค่าคงที่แต่ละค่าอาจไม่เท่ากัน ดังนั้นในปี ค.ศ.1995 Fikri Akdeniz และ Selabattin Kaciranlar ได้เสนอการประมาณค่าสัมประสิทธิ์การถดถอยพหุคูณเป็นวิธีของลิว คีเจียนที่อยู่ในรูปทั่วไปโดยผู้วิจัยเรียกวิธีการนี้ว่าวิธีลิว คีเจียนทั่วไป (generalized Liu Kejian method) ซึ่งค่าประมาณสัมประสิทธิ์การถดถอยพหุคูณอยู่ในรูปแบบดังนี้

$$(1.5) \quad \hat{\beta}_L(D) = (X'X + I)^{-1}(X'y + D\hat{\beta})$$

โดยที่ D เป็นเมทริกซ์ขนาด $p \times p$ ซึ่งสมาชิกในแนวทแยงมุมเป็นค่าคงที่และสมาชิกนอกแนวทแยงมุมเป็น 0 กล่าวคือ $D = \text{diag}(d_1, d_2, \dots, d_p)$; $0 < d_i < 1$ ดังนั้นค่าประมาณสัมประสิทธิ์การถดถอยพหุคูณสมการที่ (1.4) เป็นกรณีเฉพาะของค่าประมาณสัมประสิทธิ์การถดถอยพหุคูณสมการที่ (1.5) เมื่อ $d_1 = d_2 = \dots = d_p = d$

จากที่กล่าวมาข้างต้นผู้วิจัยจึงสนใจเปรียบเทียบวิธีประมาณค่าสัมประสิทธิ์การถดถอยพหุคูณตามวิธีดังนี้ วิธีการแรกคือวิธีกำลังสองน้อยสุด (สมการที่ (1.1)) ซึ่งเป็นวิธีที่นิยมใช้กันส่วนใหญ่และตัวประมาณที่ได้มีคุณสมบัติเป็นตัวประมาณไม่เอนเอียงและให้ความแปรปรวนต่ำสุดในบรรดาตัวประมาณไม่เอนเอียงเชิงเส้น วิธีการที่สองคือวิธีรีดจ์รีเกรสชันที่ใช้ข้อมูลสนเทศโดยหลักเกณฑ์ (สมการที่ (1.3)) ซึ่งปรับปรุงจากวิธีรีดจ์รีเกรสชันทำให้การประมาณค่าสัมประสิทธิ์การถดถอยพหุคูณมีความเชื่อถือได้มากขึ้นและวิธีการที่สามคือวิธีลิว คีเจียนทั่วไป (สมการที่ (1.5)) ซึ่งเป็นวิธีที่ผสมผสานข้อดีของวิธีรีดจ์รีเกรสชันและวิธีของส ไคน์ ผู้วิจัยจึงเปรียบเทียบตามวิธีการดังกล่าวเพื่อศึกษาว่าการประมาณค่าสัมประสิทธิ์การถดถอยพหุคูณด้วยวิธีใดมีประสิทธิภาพที่สุดเมื่อเกิดพหุสัมพันธ์ระหว่างตัวแปรอิสระโดยใช้ค่าเฉลี่ยความคลาดเคลื่อนกำลังสองเป็นเกณฑ์ในการพิจารณา

2. วัตถุประสงค์ของการวิจัย

การวิจัยครั้งนี้มีวัตถุประสงค์เพื่อเปรียบเทียบวิธีประมาณค่าสัมประสิทธิ์การถดถอยพหุคูณเมื่อเกิดพหุสัมพันธ์ระหว่างตัวแปรอิสระ ซึ่งในที่นี้มี 3 วิธี คือ

1. วิธีกำลังสองน้อยสุด (Least Square Method (LS))
2. วิธีรีดจ์รีเกรสชันที่ใช้ข้อมูลสนเทศโดยหลักเกณฑ์ (Ridge Regression with Prior Information Method (RP))
3. วิธีลิว คีเจียนทั่วไป (Generalized Lin Kajian Method (LK))

3. ขอบเขตงานของการวิจัย

เมื่อตัวแปรอิสระเกิดพหุสัมพันธ์ระดับสูงวิธี RP และ LK ให้ค่าประมาณสัมประสิทธิ์การถดถอยพหุคูณมีความถูกต้องและเชื่อถือได้มากกว่าวิธี LS ภายใต้จำนวนตัวแปรอิสระ ขนาดตัวอย่างและระดับสัมประสิทธิ์ความแปรผันของความคลาดเคลื่อนเดียวกัน และวิธี LS จะมีประสิทธิภาพดีขึ้นเมื่อการแจกแจงของความคลาดเคลื่อนอยู่เข้าสู่การแจกแจงแบบปกติ

4. ข้อตกลงเบื้องต้น

เนื่องจากเราสามารถเขียนค่าประมาณสัมประสิทธิ์การถดถอยพหุคูณวิธีวิเศษวิธีเกรตส์ที่ใช้ข้อสมมติ โดยหักกันภายในอีกรูปแบบหนึ่งดังนี้

$$\hat{\beta}_j(kL, j) = (I + k(X'X)^{-1})^{-1}(\hat{\beta} - j) + j$$

ดังนั้นเราจะได้ว่า $\hat{\beta}_j(kL, j) = \hat{\beta}$ เมื่อ $j = \hat{\beta}$ ซึ่งค่าประมาณสัมประสิทธิ์การถดถอยพหุคูณที่ได้ คือค่าประมาณสัมประสิทธิ์การถดถอยพหุคูณวิธีกำลังสองน้อยสุด ดังนั้นสมมุติฐานที่กำหนดขึ้นคือ $j \neq \hat{\beta}$ สำหรับการวิจัยครั้งนี้ผู้วิจัยกำหนด j เป็นเวกเตอร์ที่ทุกสมาชิกในเวกเตอร์คือค่าเฉลี่ยของ $\hat{\beta}$ ที่ได้จากวิธีกำลังสองน้อยสุด

5. ขอบเขตของการวิจัย

1. ตัวแปรอิสระที่ใช้ในการวิจัยมี 2 จำนวน คือ 3 และ 5 ตัวแปร
2. ขนาดตัวอย่าง (n) มี 4 ขนาด คือ 12, 30, 50 และ 100
3. การแจกแจงของความคลาดเคลื่อนมี 3 การแจกแจง คือ

3.1 เมื่อความคลาดเคลื่อนมีการแจกแจงแบบปกติ (Normal Distribution)

ฟังก์ชันความหนาแน่นของ X อยู่ในรูปของ

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\} ; \sigma^2 > 0$$

การวิจัยครั้งนี้ศึกษาที่สัมประสิทธิ์ความแปรผัน (C.V.) 3 ระดับ คือ 5%, 10% และ 15% ผู้วิจัยศึกษา C.V. 3 ระดับนี้เนื่องจากผู้วิจัยพิจารณากราฟที่ C.V. มีค่าสูงการแจกแจงของความคลาดเคลื่อนเป็นการแจกแจงอื่นที่ไม่ใช่การแจกแจงแบบปกติ (จิรายุส พุ่มนศรี, "การเปรียบเทียบตัวประมาณวิเศษสำหรับการวิเคราะห์การถดถอยแบบวิเศษ", วิทยานิพนธ์ปริญญาโทบริหารศาสตราจารย์มหาบัณฑิต สาขาวิชาสถิติ บัณฑิตวิทยาลัย จุฬาลงกรณ์มหาวิทยาลัย, 2534) หน้า 6.

3.2 เมื่อความคลาดเคลื่อนมีการแจกแจงแบบปกติปน (Scale-Contaminated Normal Distribution)

ฟังก์ชันการแจกแจงอยู่ในรูปของ

$$F = (1-p) N(\mu, \sigma^2) + p N(\mu, c^2 \sigma^2)$$

เมื่อ c คือ สเกลแฟกเตอร์ (scale factor) (ถ้าสเกลแฟกเตอร์มีค่าสูงทำให้ค่าสังเกตที่มีค่าน้อยมีค่าสูง)
 และ p คือ เปอร์เซ็นต์การปลอมปน (percent of contamination) (ถ้าเปอร์เซ็นต์การปลอมปนมีค่าสูงทำให้โอกาสที่เกิดค่าสังเกตที่มีค่าน้อยมีค่าสูง)

ผู้วิจัยศึกษาในกรณีที่ $c=3, 10$ และ $p=5, 10$ เมื่อกำหนด C.V. = 5%, 10%, 15% ตามลำดับ

3.3 เมื่อความคลาดเคลื่อนมีการแจกแจงแบบลอการิทึม (Lognormal Distribution)

ฟังก์ชันความหนาแน่นของ X อยู่ในรูปของ

$$f(x) = \frac{1}{\sqrt{2\pi} \sigma x} \exp\left\{-\frac{1}{2\sigma^2}(\ln x - \mu)^2\right\}; x > 0$$

สัมประสิทธิ์ความแปรผันของการแจกแจงแบบลอการิทึมอยู่ในรูปของ

$$C.V. = \sqrt{\exp(\sigma^2) - 1}$$

ผู้วิจัยศึกษาในกรณีที่ $\sigma^2 = 0.05, 0.30$ และ 0.70 ซึ่งทำให้ได้ C.V. = 23%, 59% และ 100% ตามลำดับ สำหรับสาเหตุที่ผู้วิจัยไม่เลือก C.V. ที่มีค่าต่ำกว่านี้เนื่องจากผู้วิจัยพิจารณากราฟที่ C.V. มีค่าต่ำกว่านี้กราฟของการแจกแจงจะดูเข้าสู่การแจกแจงแบบปกติมากขึ้น (เจษฎาพร ตูทรนวิบูลย์ชัย, "การศึกษาเปรียบเทียบตัวประมาณวิธี", วิทยานิพนธ์ปริญญาโท สาขาวิชาสถิติ บัณฑิตวิทยาลัย จุฬาลงกรณ์มหาวิทยาลัย, 2533) หน้า 7.

4. การกำหนดค่าสัมประสิทธิ์การถดถอยพหุคูณ (β) เพื่อสร้างค่า \hat{y} ขึ้นจากตัวแบบ $y = X\beta + \epsilon$
 ผู้วิจัยกำหนดค่า β สองชุดดังนี้

- ก) กำหนด β จากเวกเตอร์เฉพาะที่สอดคล้องกับค่าเฉพาะที่มีค่ามากที่สุดของเมทริกซ์ $X'X$
 โดยผู้วิจัยมีเหตุผลว่าถ้ากำหนด β เช่นนี้ทำให้ค่ากำลังสองของความแตกต่างระหว่าง \hat{y} และ y มีค่าน้อยสุด
- ข) กำหนด β จากเวกเตอร์เฉพาะที่สอดคล้องกับค่าเฉพาะที่มีค่าน้อยสุดของเมทริกซ์ $X'X$
 โดยผู้วิจัยมีเหตุผลว่าถ้ากำหนด β เช่นนี้ทำให้ค่ากำลังสองของความแตกต่างระหว่าง \hat{y} และ y มีค่ามากที่สุด

5. ระดับความสัมพันธ์ระหว่างตัวแปรอิสระ

ในกรณีตัวแปรอิสระ 3 ตัว ผู้วิจัยศึกษาที่ทุกสัมพันธ์ระหว่างตัวแปรอิสระ 4 ระดับ คือ

- | | |
|------------------|---------------|
| (1) ระดับต่ำ | $\rho = 0.30$ |
| (2) ระดับปานกลาง | $\rho = 0.60$ |
| (3) ระดับสูง | $\rho = 0.90$ |
| (4) ระดับสูงมาก | $\rho = 0.99$ |

โดยที่ ρ คือ ระดับความสัมพันธ์ระหว่าง x_1 กับ x_2 , x_1 กับ x_3 และ x_2 กับ x_3

ในกรณีตัวแปรอิสระ 5 ตัว ผู้วิจัยศึกษาที่ทุกสัมพันธ์ระหว่างตัวแปรอิสระ 4 ระดับ คือ

- | | |
|------------------|-----------------------|
| (1) ระดับต่ำ | $\rho = (0.30, 0.30)$ |
| (2) ระดับปานกลาง | $\rho = (0.60, 0.60)$ |
| (3) ระดับสูง | $\rho = (0.90, 0.90)$ |
| (4) ระดับสูงมาก | $\rho = (0.99, 0.99)$ |

โดยที่ค่าแรกในวงเล็บของ ρ คือ ระดับความสัมพันธ์ระหว่าง x_1 กับ x_2 , x_1 กับ x_3 และ x_2 กับ x_3 ค่าที่สองในวงเล็บของ ρ คือ ระดับความสัมพันธ์ระหว่าง x_4 กับ x_5

6. ประโยชน์ที่คาดว่าจะได้รับ

1. ผลการศึกษากลายเป็นแนวทางให้ผู้วิจัยเลือกใช้วิธีการประมาณค่าสัมประสิทธิ์การถดถอยพหุคูณเมื่อเกิดทุกสัมพันธ์ระหว่างตัวแปรอิสระได้เหมาะสมกับสถานการณ์และเงื่อนไขที่เกิดขึ้น

2. ผลการศึกษากลายเป็นแนวทางในการแก้ไขและปรับปรุงวิธีการประมาณค่าสัมประสิทธิ์การถดถอยพหุคูณเมื่อเกิดทุกสัมพันธ์ระหว่างตัวแปรอิสระได้ต่อไปในอนาคต

สถาบันวิทยบริการ
จุฬาลงกรณ์มหาวิทยาลัย