



บทที่ 1

บทนำ

1. ความสำคัญและความเป็นมาของปัญหา

ในการประมาณค่าเพื่อคาดคะเนเหตุการณ์ล่วงหน้าหรือการพยากรณ์ ผู้วิจัยมักเลือกวิธีการวิเคราะห์ความถดถอยพหุ (Multiple regression analysis) เมื่อตัวแปรตามที่ใช้ในการศึกษามีความสัมพันธ์กับปัจจัยอื่น ๆ เราเรียกว่าตัวแปรอิสระ มากกว่า 1 ตัว ซึ่งสามารถเขียนในรูปของตัวแบบทั่วไปของการวิเคราะห์ความถดถอยพหุเชิงเส้นดังนี้

$$(1.1) \quad \tilde{y} = X\tilde{\beta} + \tilde{\varepsilon}$$

เมื่อ \tilde{y} เป็นเวกเตอร์ของตัวแปรตามขนาด $n \times 1$ โดยที่ n เป็นจำนวนค่าสังเกต

X เป็นเมทริกซ์ของตัวแปรอิสระขนาด $n \times p$ ($p < n$) และมี full rank p

$\tilde{\beta}$ เป็นเวกเตอร์ของพารามิเตอร์ที่ไม่ทราบค่า ขนาด $p \times 1$

และ $\tilde{\varepsilon}$ เป็นเวกเตอร์ของความผิดพลาด ขนาด $n \times 1$

โดยมีข้อตกลงเบื้องต้น (Assumption) ของความผิดพลาดดังนี้

$$E(\tilde{\varepsilon}) = \underline{0}$$

$$E(\tilde{\varepsilon}\tilde{\varepsilon}') = \sigma^2 I_n$$

วิธีที่นิยมใช้ในการประมาณค่าสัมประสิทธิ์การถดถอยพหุกันมาก คือ วิธีกำลังสองน้อยที่สุด (least square method) ซึ่งตัวประมาณกำลังสองน้อยที่สุดอยู่ในรูปของ

$$(1.2) \quad \hat{\tilde{\beta}} = (X'X)^{-1}X'\tilde{y}$$

โดยจะเป็นตัวประมาณที่ไม่เอนเอียง (unbiased estimator) และมีความแปรปรวนต่ำสุดในบรรดาตัวประมาณที่ไม่เอนเอียง

ในกรณีที่มีค่าลุ่มตัวอย่างที่ได้มาจากประชากรที่มีการแจกแจงแบบเบ้ และการแจกแจงแบบหางยาวกว่าปกติ (skewed distribution and long tailed distribution) วิธีการข้างสองน้อยที่สุดอาจจะไม่เหมาะสม เพราะวิธีนี้มีความไวต่อข้อมูลที่ผิดปกติ (Outliers) และสูญเสียประสิทธิภาพไปเมื่อการแจกแจงของความผิดพลาดไม่เป็นแบบปกติ ซึ่งได้มีผู้ค้นหาวิธีการแก้ปัญหาที่เกิดขึ้นนี้กันหลายวิธี แต่ที่รู้จักกันดี คือ วิธีการทางนอนพาราเมตริก โดยวิธีการประมาณค่าความผิดพลาดมาตรฐาน (Standard error) ในกรณีที่ไม่ทราบลักษณะการแจกแจงของประชากรและไม่สามารถหาได้จากสูตรทั่วไป โดยใช้เทคนิคการลุ่มตัวอย่างซ้ำ (Resampling) ซึ่งมีอยู่หลายวิธี ได้แก่

The Jackknife

The Bootstrap

Half-sampling

Subsampling

Balanced repeated replication

Influence function techniques

The delta method

แต่ละวิธีมาจากแนวความคิดพื้นฐานคล้ายกัน คือ การหาค่าประมาณของค่าความผิดพลาดมาตรฐาน โดยการลุ่มซ้ำจากข้อมูลที่เก็บรวบรวมมา ซึ่งพบว่าวิธีบูตสเตรป (The Bootstrap Method) เป็นวิธีที่ให้ผลดีที่สุด เพราะว่าการหาค่าประมาณโดยวิธีนี้เป็น การประมาณค่าภาวะน่าจะเป็นสูงสุดของนอนพาราเมตริก ทำให้ตัวประมาณที่ได้เป็นตัวประมาณแบบภาวะน่าจะเป็นสูงสุด (Maximum Likelihood Estimator (MLE)) และจากการศึกษาของ Bradley Efron (1982; 27) ได้ศึกษาวิธีบูตสเตรป พบว่าสามารถนำมาใช้ในการประมาณค่าคร่าว ๆ ที่ไม่สามารถหาค่าได้โดยตรงทางพาราเมตริก เช่น ไม่ทราบลักษณะการแจกแจงของประชากร ไม่ทราบลักษณะการแจกแจงของความผิดพลาด

อีกวิธีหนึ่ง คือ ตัวประมาณ M ซึ่งอยู่ในรูปของค่าน้อยที่สุดของฟังก์ชัน ค่าผิดพลาดนี้สามารถเขียนได้ดังนี้

$$\min_{\beta} \sum_{i=1}^n \rho(\epsilon_i/s) = \min_{\beta} \sum_{i=1}^n \rho[(y_i - X_i\beta)/s]$$

เมื่อ ϵ_i เป็นค่าผิดพลาดของค่าสังเกตที่ i และ s เป็นค่าที่เหมาะสมสำหรับการกระจายของ ϵ_i ซึ่ง Huber (1981) ได้เสนอข้อคิดเห็นเมื่อค่าสังเกตมีการกระจายแบบปโลมปน และค่าสังเกตที่ได้จากกรณีที่มีการแจกแจงแบบสมมาตร ว่าการเลือกตัวประมาณ s ที่เหมาะสมขึ้นกับ

1. ขอบเขตของ breakdown point มีขนาดใหญ่
2. มีความเอนเอียงไม่มาก
3. มีความแปรปรวนไม่มาก

ในการวิจัยครั้งนี้ ผู้วิจัยจึงสนใจที่จะศึกษาเปรียบเทียบการประมาณสัมประสิทธิ์การถดถอยของ 3 วิธี คือ วิธีกำลังสองน้อยที่สุด วิธีบูตสเตรป และวิธีตัวประมาณชนิด M ภายใต้สภาวะการกระจายของลักษณะข้อมูลที่มีการแจกแจงแบบต่าง ๆ โดยใช้ค่าเฉลี่ยความผิดพลาดกำลังสอง (MSE) ของตัวประมาณแต่ละตัวเป็นตัวเปรียบเทียบ

1.2 วัตถุประสงค์ของการวิจัย

1. เพื่อศึกษาวิธีการประมาณค่าสัมประสิทธิ์การถดถอยพหุ และคุณสมบัติของตัวประมาณของวิธีกำลังสองน้อยที่สุดวิธีบูตสเตรป และวิธีตัวประมาณชนิด M เมื่อค่าความผิดพลาดมีการแจกแจงแบบต่างๆ
2. เพื่อศึกษาวิธีการประมาณค่าสัมประสิทธิ์การถดถอยพหุวิธีกำลังสองน้อยที่สุด วิธีบูตสเตรป และวิธีตัวประมาณชนิด M เมื่อกำหนดขนาดของตัวอย่างต่างๆ กัน
3. เพื่อศึกษาเปรียบเทียบประสิทธิภาพของการประมาณค่าสัมประสิทธิ์การถดถอยพหุด้วยวิธีกำลังสองน้อยที่สุด วิธีบูตสเตรป และวิธีตัวประมาณชนิด M

1.3 ข้อตกลงเบื้องต้น

1. ค่าความผิดพลาด (ϵ_i) เป็นตัวแปรสุ่มที่มีการแจกแจงเหมือนกัน และเป็นอิสระซึ่งกันและกัน

2. การวิจัยครั้งนี้ใช้เกณฑ์การประมาณค่าสัมประสิทธิ์การถดถอยพหุภาคย์ได้ค่าความผิดพลาดที่มีการแจกแจงแบบต่าง ๆ จากวิธีใดที่ทำให้ค่าเฉลี่ยความผิดพลาดกำลังสอง (Mean Square Error) น้อยที่สุด จะเป็นวิธีที่เหมาะสมของแต่ละสถานการณ์

1.4 ขอบเขตการวิจัย

1.4.1 ลักษณะการแจกแจงของความผิดพลาดที่ศึกษามีดังนี้

1.4.1.1 การแจกแจงแบบปกติปลอมปน (scale-contaminated normal distribution)

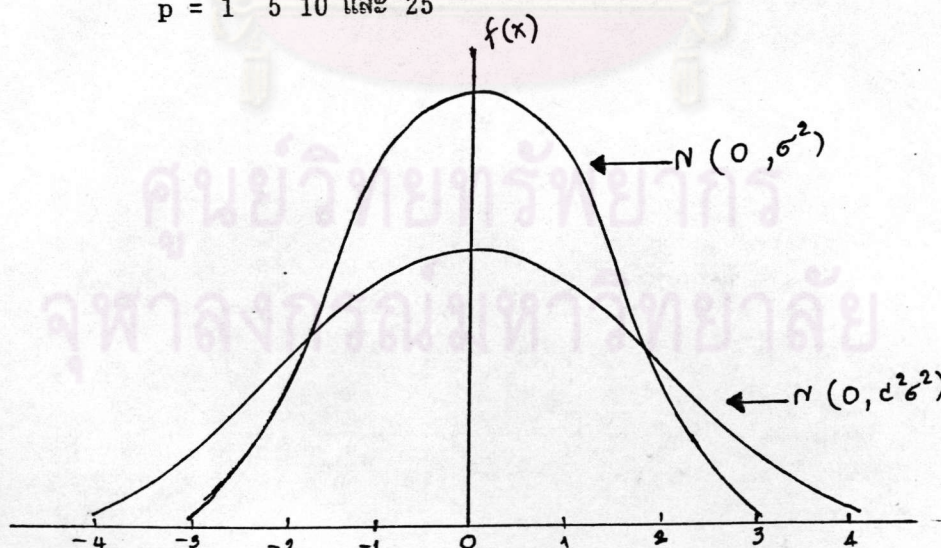
ฟังก์ชันการแจกแจงอยู่ในรูปของ

$$F = (1 - p) N(0, \sigma^2) + p N(0, c^2 \sigma^2)$$

เมื่อ c คือสเกลแฟกเตอร์ (scale factor) ซึ่งถ้าค่าสูงจะทำให้ค่าสังเกตที่ผิดปกติมีค่าสูงด้วย ในที่นี้ใช้ $c = 3$ และ 10

และ p คือเปอร์เซ็นต์การปลอมปน (percent of contamination) ในที่นี้ใช้

$p = 1 \quad 5 \quad 10$ และ 25



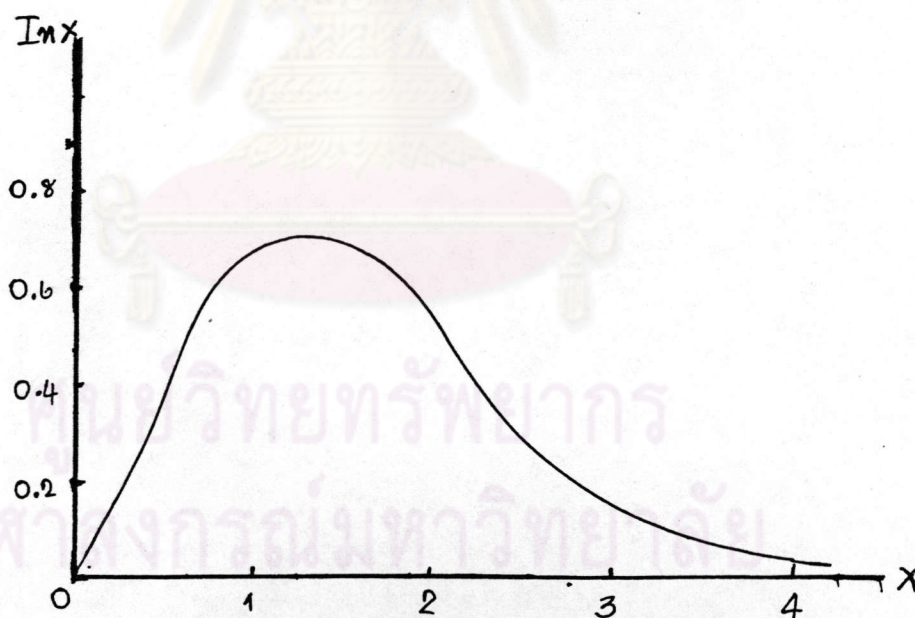
รูปที่ 1.4.1 แสดงเส้นโค้งของการแจกแจงแบบปกติปลอมปน ณ p และ c

1.4.1.2 การแจกแจงแบบลอการิทึม (lognormal distribution)

ฟังก์ชันความหนาแน่นอยู่ในรูปของ

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

เมื่อ μ และ σ^2 เป็นค่าเฉลี่ยและความแปรปรวนของ Y โดยที่ $Y = \ln X$
 และ Y มีการแจกแจงแบบปกติ ในที่นี้พิจารณาใช้ $\mu = 0, \sigma^2 = 1$



รูปที่ 1.4.3 แสดงเส้นโค้งของการแจกแจงแบบลอการิทึม

1.4.1.3 การแจกแจงแบบแกมมา (gamma-distribution)

ฟังก์ชันความหนาแน่นอยู่ในรูปของ

$$f(x) = \begin{cases} \frac{x^{\alpha-1} \exp\{-x/\beta\}}{\beta^{\alpha} \Gamma(\alpha)} & ; x > 0, \alpha > 0, \beta > 0 \\ 0 & ; \text{อื่น ๆ} \end{cases}$$

เมื่อ β เป็น scale parameter

และ α เป็น shape parameter

$$\text{จะได้ว่า } E(X) = \beta \alpha$$

$$\text{Var}(X) = \beta^2 \alpha$$

$$\text{สัมประสิทธิ์ความแปรปรวน (C.V.)} = 1/\alpha$$

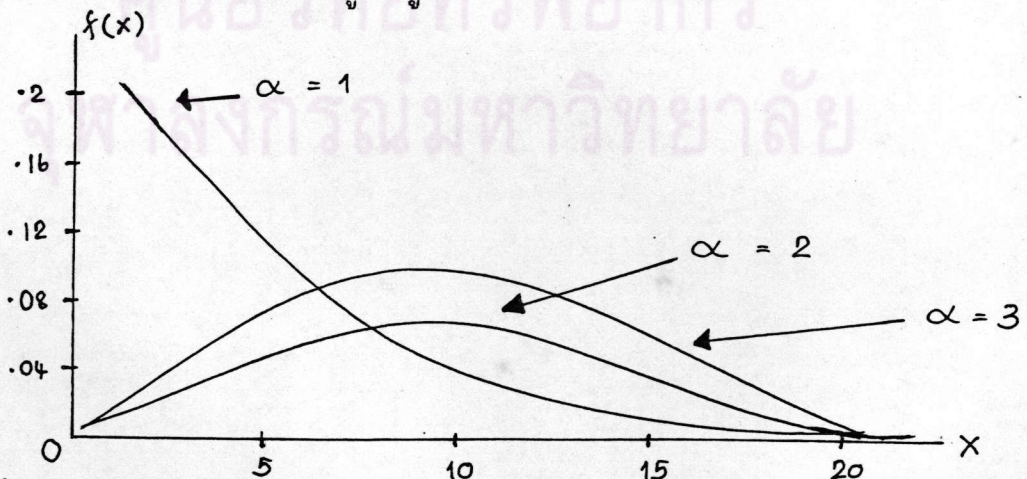
ในการวิจัยครั้งนี้จะศึกษา ณ ค่า $\alpha = 5$ 10 และ 150 เมื่อ $\beta = 1$ และ 2
กล่าวคือ

$$\text{C.V.}(X) = 100\% (\alpha = 1, \beta = 1)$$

$$\text{C.V.}(X) = 70\% (\alpha = 2, \beta = 1)$$

$$\text{C.V.}(X) = 58\% (\alpha = 3, \beta = 1)$$

และสาเหตุที่เลือกใช้ค่า C.V. = 100% 70% และ 58% โดยที่ไม่เลือกค่า C.V. ที่มีค่าต่ำกว่านี้ เพราะจากการพิจารณากราฟที่แสดงเส้นโค้งของการแจกแจงแบบแกมมา ถ้าค่า C.V. ที่มีค่าต่ำกว่านี้ กราฟของการแจกแจงจะลู่เข้าสู่การแจกแจงแบบปกติมากขึ้น

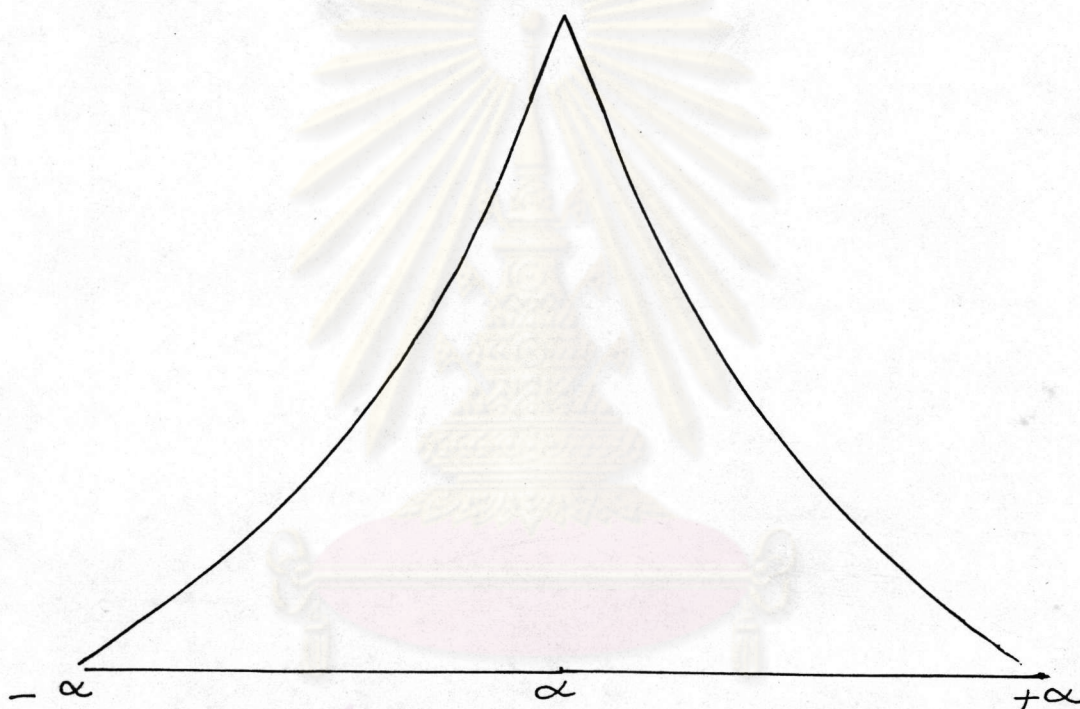


รูปที่ 1.4.4 แสดงเส้นโค้งของการแจกแจงแบบแกมมา ณ $\beta = 1$ และ $\alpha = 1, 2$ และ 3

1.4.1.4 การแจกแจงแบบคัมเบิ้ลเอ็กซ์โปเนนเชียล (Double Exponential Distribution)

ฟังก์ชันความน่าจะเป็น คือ

$$f(x) = \frac{1}{2\beta} \cdot e^{-\left| \frac{x - \alpha}{\beta} \right|} \quad ; -\alpha < x < \alpha$$



ค่าคาดหวัง $E(X) = \alpha$

ค่าความแปรปรวน $V(X) = 2\beta^2$

1.4.8 ขนาดตัวอย่างและตัวแปรอิสระ

จำนวนตัวแปรอิสระและขนาดตัวอย่างที่ใช้ในการวิจัยครั้งนี้แบ่งเป็น 3 กลุ่ม

ดังนี้

1. ตัวแปรอิสระ $p = 5$ จะใช้ขนาดตัวอย่าง 50 และ 100
2. ตัวแปรอิสระ $p = 3$ จะใช้ขนาดตัวอย่าง 5 10 และ 20
3. ตัวแปรอิสระ $p = 2$ จะใช้ขนาดตัวอย่าง 4

ทั้งนี้ เพื่อให้ได้ผลของการวิจัยได้ครอบคลุมทั้งขนาดตัวอย่างและจำนวนตัวแปรอิสระ ที่มีปริมาณน้อย ปานกลาง และใหญ่ คือจะใช้ที่จำนวนตัวแปรอิสระ = 5 ขนาดตัวอย่าง 50 และ 100 เป็นตัวแทนของกลุ่มตัวอย่างขนาดใหญ่ สำหรับที่จำนวนตัวแปรอิสระ = 3 ขนาดตัวอย่าง = 5, 10, 20 เป็นตัวแทนของกลุ่มตัวอย่างขนาดปานกลาง และที่จำนวนตัวแปรอิสระ = 2 ขนาดตัวอย่าง = 4 เป็นตัวแทนของกลุ่มตัวอย่างที่มีขนาดเล็ก

1.5 ประโยชน์ที่คาดว่าจะได้รับ

ผลการศึกษาวិธีการประมาณและคุณสมบัติของตัวประมาณจะทำให้สามารถบอกได้ว่าวิธีการประมาณค่าสัมประสิทธิ์ถดถอยพหุทั้ง 3 มีคุณสมบัติและวิธีการที่จะนำไปใช้ในการประมาณค่าสัมประสิทธิ์ถดถอยพหุได้อย่างไร

1.6 คำจำกัดความ

1.6.1 ค่าเฉลี่ยความผิดพลาดกำลังสองเฉลี่ย (Mean Square Error) หรือ MSE ของตัวประมาณ คือ ถ้า $\hat{\theta}$ เป็นตัวประมาณของพารามิเตอร์ แล้ว ความผิดพลาดกำลังสองเฉลี่ยของ θ คือ $E(\hat{\theta} - \theta)^2$

ในการคำนวณหาความผิดพลาดกำลังสองเฉลี่ยของฟังก์ชันเชิงเส้นของตัวประมาณ θ ในที่นี้คือ $1^T\beta$; เมื่อ 1 เป็นเมตริกซ์ขนาด $p \times 1$ ที่สมาชิกทุกตัวมีค่าเป็น 1 นั่นก็คือ การหาความผิดพลาดกำลังสองเฉลี่ยในรูปของผลบวกของตัวประมาณนั่นเอง

1.6.2 ความแปรปรวน (Variance) ของตัวประมาณ คือ ถ้า $\hat{\theta}$ เป็นตัวประมาณของพารามิเตอร์ θ แล้ว ความแปรปรวนของ $\hat{\theta}$ คือ $E(\hat{\theta} - E(\hat{\theta}))^2$

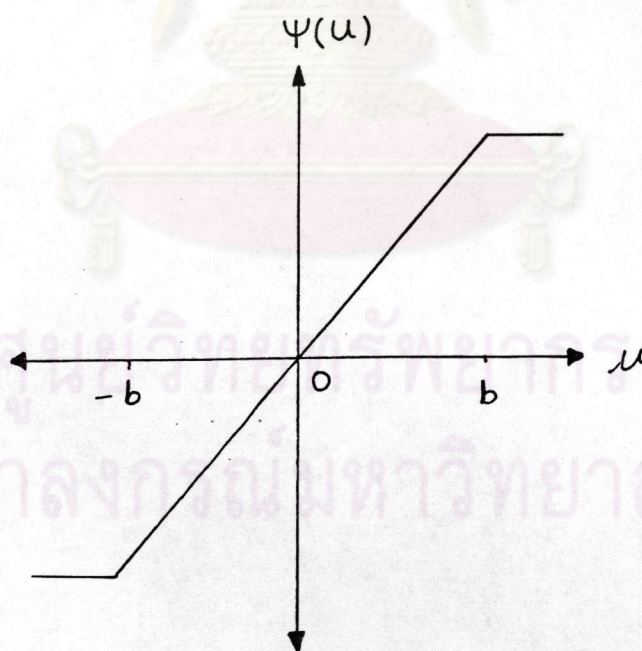
1.6.3 ความไม่เอนเอียง (Unbiased) ของตัวประมาณ คือถ้า $\hat{\theta}$ เป็นตัวประมาณที่ไม่เอนเอียงของ θ ก็ต่อเมื่อ $E(\hat{\theta}) = \theta$

1.6.4 ตัวประมาณเชิงเส้นที่ดีที่สุดและไม่เอนเอียง (Best Linear Unbiased Estimator) หรือ BLUE เป็นคุณสมบัติหนึ่งของตัวประมาณ โดยตัวประมาณ $\hat{\theta}$ จะมีคุณสมบัติเป็น BLUE ของพารามิเตอร์ θ ถ้า $\hat{\theta}$ มีคุณสมบัติครบ 3 ข้อดังต่อไปนี้

- 5.1 เป็นฟังก์ชันเชิงเส้นของตัวแปรสุ่ม
- 5.2 เป็นตัวประมาณที่ไม่เอนเอียง
- 5.3 เป็นตัวประมาณที่มีค่าความแปรปรวนต่ำสุด

1.6.5 ตัวประมาณที่มีค่าความแปรปรวนต่ำสุดเพียงตัวเดียวในบรรดาตัวประมาณที่ไม่เอนเอียง (Uniformly Minimum Variance Unbiased Estimator) หรือ UMVUE เป็นคุณสมบัติของตัวประมาณที่ผู้วิจัยต้องการ θ คือ ถ้ามีตัวประมาณที่ไม่เอนเอียงสำหรับ θ ซึ่งเป็นสถิติที่พอเพียงสำหรับ θ และมีค่าความแปรปรวนต่ำกว่าค่าความแปรปรวนของตัวประมาณอื่นๆ ที่ไม่เอนเอียงสำหรับ θ แล้ว ตัวประมาณดังกล่าวจะมีคุณสมบัติเป็นตัวประมาณที่ดีที่สุดหรือมีค่าความแปรปรวนต่ำสุดเพียงตัวเดียวในบรรดาตัวประมาณที่ไม่เอนเอียง

1.6.6 เกณฑ์ความแกร่งของ Huber คือสมการเชิงเส้นที่จุดกำเนิดและมีค่าคงที่ที่ปลายทั้งสองข้างพิจารณาจากรูปที่ 1.6.7



รูป 1.6.7 แสดงกราฟของเกณฑ์ความแกร่ง Huber ψ เมื่อ b คือจุดเปลี่ยนเว้า

และมีฟังก์ชันความแรงดังนี้

$$(r) = \begin{cases} -r^2/2 & ; \text{ ถ้า } |r| < k \\ k|r| = k^2/2 & ; \text{ อื่น ๆ } \end{cases}$$

ซึ่งค่าประมาณของ Huber เป็นค่าประมาณภาวะน่าจะเป็นสูงสุด (MLE) ที่ดีสำหรับการแจกแจงแบบ \mathcal{E} -contaminate Gaussian และทำให้ค่าประมาณสเกลเป็นเงื่อนไขที่สำคัญของค่าสังเกตในการแปลงข้อมูลจากกำลังสองในรูปสมการเชิงเส้น

ศูนย์วิทยทรัพยากร
จุฬาลงกรณ์มหาวิทยาลัย