

## บทที่ 2

### ระเบียบวิธีที่ใช้ในการวิจัย

การศึกษาค้นคว้าวิจัยเชิงเส้นพหุเมื่อตัวแปรตามบางค่าถูกตัดทิ้งทางขวาประเภทที่ 1 โดยในการศึกษาค้นคว้าครั้งนี้ได้ทำการศึกษาเปรียบเทียบด้วยวิธีการประมาณพารามิเตอร์ 4 วิธี คือ วิธีกำลังสองค่าสุด วิธีการของแมทเทอร์จีและแมคกีธ วิธีการของบักเลย์และเจมส์ และวิธีภาวะน่าจะเป็นสูงสุดด้วยขั้นตอนวิธีอีเอ็ม ซึ่งในบทนี้กล่าวถึงความรู้ ทฤษฎีพื้นฐาน และรายละเอียดของวิธีการประมาณพารามิเตอร์แต่ละวิธี รายละเอียดต่างๆ เป็นดังนี้

#### 2.1 ทฤษฎีพื้นฐาน

##### 2.1.1 ประเภทของการถูกตัด (Type of Censoring)

###### ๑. การถูกตัดประเภทที่ 1 (Type I Censoring)

ข้อมูลที่ถูกตัดประเภทที่ 1 เกิดขึ้นเนื่องจากการกำหนดค่าสูงสุดของข้อมูลล่วงหน้าด้วยค่าคงที่  $T_c$  ซึ่งจะเรียกว่า Fixed Censoring Time ตัวอย่างเช่นการทดลองเกี่ยวกับหนูทดลองที่กำหนดระยะเวลาของการทดลองไว้เท่ากับ 3 เดือนแล้วเริ่มทำการทดลอง โดยการฉีดยาที่ต้องการทดลองให้กับหนู แล้วเริ่มบันทึกเวลาตั้งแต่เริ่มทำการทดลองจนกระทั่งหนูทดลองตาย ถ้าหนูตัวใดตายภายใน 3 เดือนจะถือว่าเป็นข้อมูลที่ไม่ถูกตัด และจำนวนวันตั้งแต่เริ่มทดลองจนหนูตายเรียกว่าเป็น Survival Time แต่เมื่อครบระยะเวลาทดลอง 3 เดือนแล้วหนูทดลองตัวใดที่ไม่ตาย จะถือว่าค่าสังเกตของหนูตัวนั้นถูกตัดที่ Fixed Censoring Time เท่ากับ 3 เดือน สำหรับข้อมูลทางค่านประกันภัยที่กรมธรรม์จะกำหนดความเสียหายสูงสุดที่บริษัทรับผิดชอบไว้ เช่นการรับประกันสุขภาพที่บริษัทจะจ่ายค่ารักษาพยาบาลให้ตามที่ผู้เอาประกันภัยจ่ายไปจริงแต่ไม่เกิน 10,000 บาท ถ้าผู้เอาประกันภัยรายใดต้องเสียค่ารักษาจริงมากกว่า 10,000 บาท บริษัทจะจ่ายชดเชยให้เพียง 10,000 บาทเท่านั้น และถือว่าข้อมูลนี้ถูกตัดที่  $T_c$  เท่ากับ 10,000 บาท

ถ้า  $T_1, T_2, \dots, T_N$  เป็นค่าสังเกตที่ไม่ถูกตัดที่มีการแจกแจงเหมือนกันและเป็นอิสระกัน  $T_c$  เป็นเวลาหรือค่าสูงสุดที่กำหนดไว้ล่วงหน้า จะได้ตัวแปรสุ่มของค่าสังเกต  $Y_1, Y_2, \dots, Y_N$  ซึ่ง

$$Y_i = \begin{cases} T_i & ; T_i \leq T_c \\ T_c & ; T_i > T_c \end{cases}$$

โดยมีฟังก์ชันภาวะน่าจะเป็น (Likelihood Function) ดังนี้

$$L(y_i) = \begin{cases} f(y_i) & \text{ถ้าค่าสังเกต ไม่ถูกตัด} \\ P(T_i > T_c) - S(T_c) & \text{ถ้าค่าสังเกตถูกตัด} \end{cases}$$

และมีฟังก์ชันภาวะน่าจะเป็นรวมดังนี้

$$L = \prod_{i \in u} f(y_i) \prod_{i \in c} S(T_c)$$

$i \in u$  หมายถึงค่าสังเกตที่ไม่ถูกตัดทั้ง

$i \in c$  หมายถึงค่าสังเกตที่ถูกตัดทั้ง

## ข. การถูกตัดประเภทที่ 2 (Type II Censoring)

ในบางกรณีที่ไม่อาจกำหนดเวลา หรือค่าสูงสุดของการตัดข้อมูลที่เหมาะสมได้ ดังนั้นจะกำหนดจำนวนค่าสังเกตที่ไม่ถูกตัดแทน นั่นคือเมื่อค่าสังเกตที่ไม่ถูกตัดเกิดขึ้นครบตามจำนวนที่กำหนดแล้วก็จะหยุดการทดลอง เช่นการทดสอบอายุการใช้งานของหลอดไฟ จะกำหนดจำนวนหลอดไฟที่เสื่อมสภาพไว้ล่วงหน้า เริ่มทดลองโดยเปิดให้หลอดไฟทำงานทั้งหมด เริ่มบันทึกเวลาและนับจำนวนหลอดไฟที่เสื่อมสภาพ เมื่อได้จำนวนหลอดไฟที่เสื่อมสภาพครบแล้วก็จะหยุดทำการทดสอบ

ถ้า  $N$  คือจำนวนข้อมูลทั้งหมดและกำหนด  $n$  คือจำนวนค่าสังเกตที่ไม่ถูกตัด  $n \leq N$  ให้  $T_1 \leq T_2 \leq \dots \leq T_n$  เป็นค่าสังเกตที่ไม่ถูกตัด และ  $T_{n+1} \leq T_{n+2} \leq \dots \leq T_N$  เป็นค่าสังเกตที่ถูกตัด ซึ่ง  $T_i \geq T_n$  ;  $i = n+1, n+2, \dots, N$  ไม่ทราบค่าที่แท้จริงของค่าสังเกต ดังนั้น  $Y_i$  เป็นตัวแปรสุ่มของค่าสังเกตซึ่ง

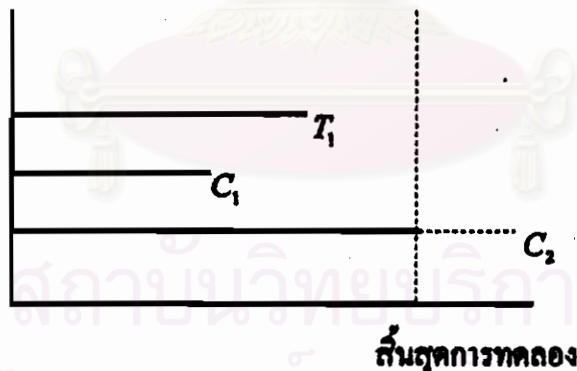
$$Y_i = \begin{cases} T_i & \text{เมื่อ } i \leq n \\ T_n & \text{เมื่อ } i = n+1, n+2, \dots, N \end{cases}$$

ฟังก์ชันความหนาแน่นร่วมของค่าสังเกตคือ

$$\frac{N!}{(N-n)!} f(y_1) f(y_2) \dots f(y_n) [S(y_n)]^{N-n}$$

### ก. การตัดแบบสุ่ม (Random Censoring)

ข้อมูลลักษณะนี้เกิดขึ้นคล้ายกับการตัดประเภทที่ 1 คือมีการกำหนดค่าสูงสุดหรือกำหนดเวลาไว้ล่วงหน้า แต่การตัดข้อมูลอาจจะเกิดขึ้นก่อนนั้นได้ เช่นการทดลองทางการแพทย์ที่คนไข้ถอนตัวออกจากการศึกษาการทดลองก่อนสิ้นสุดการทดลอง หรือคนไข้ยังมีชีวิตอยู่รอดเมื่อสิ้นสุดการทดลอง หรือคนไข้เสียชีวิตเนื่องจากสาเหตุอื่นที่ไม่เกี่ยวข้องกับสิ่งที่กำลังศึกษาทดลอง เป็นต้น จึงไม่สามารถทราบค่าที่แน่นอนของค่าสังเกตนั้นได้ รูปที่ 2.1 แสดงลักษณะความเป็นไปได้ของข้อมูลถูกตัดแบบสุ่ม



รูปที่ 2.1 แผนภาพแสดงการเกิดค่าถูกตัดแบบสุ่ม

คนไข้คนที่ 1 เข้าทำการทดลองตั้งแต่เริ่มการทดลอง และเสียชีวิตที่เวลา  $T_1$  จะถือว่าเป็นค่าสังเกตที่ไม่ถูกตัดทิ้ง

คนไข้คนที่ 2 เข้าทำการทดลองตั้งแต่เริ่มการทดลอง และถอนตัวจากการทดลองที่เวลา  $C_1$  จะถือว่าเป็นค่าสังเกตที่ถูกตัดทิ้ง

คนไข้คนที่ 3 เข้าทำการทดลองตั้งแต่เริ่มการทดลอง และมีชีวิตอยู่รอดจนสิ้นสุดการทดลอง จะถือว่าเป็นค่าสังเกตที่ถูกตัดทิ้ง  $C_2$

ถ้า  $T_1, T_2, \dots, T_N$  เป็นตัวแปรสุ่มของค่าสังเกตที่ไม่ถูกตัดที่มีการแจกแจงเหมือนกันและเป็นอิสระกันมีฟังก์ชันการอยู่รอดและฟังก์ชันความหนาแน่น  $S$  และ  $f$  ตามลำดับ และ  $C_1, C_2, \dots, C_N$  เป็นตัวแปรสุ่มของค่าสังเกตที่ถูกตัดทั้งที่มีการแจกแจงเหมือนกันและเป็นอิสระกันมีฟังก์ชันการอยู่รอดและฟังก์ชันความหนาแน่น  $G$  และ  $g$  ตามลำดับ

ดังนั้น  $T_i$  และ  $C_i$ ,  $i=1, 2, \dots, N$  เป็นอิสระกัน จากนิยามการตัดทั้งแบบสุ่ม นิยามให้  $Y_i = \min(T_i, C_i)$  จะได้ค่าสังเกตสุ่ม  $Y_1, Y_2, \dots, Y_N$  ดังนี้

$$Y_i = \begin{cases} T_i & \text{ถ้า } T_i \leq C_i \quad (\text{ไม่ถูกตัด}) \\ C_i & \text{ถ้า } T_i > C_i \quad (\text{ถูกตัด}) \end{cases}$$

$$\delta_i = \begin{cases} 1 & \text{ถ้า } T_i \leq C_i \quad (\text{ไม่ถูกตัด}) \\ 0 & \text{ถ้า } T_i > C_i \quad (\text{ถูกตัด}) \end{cases}$$

จะมีฟังก์ชันภาวะน่าจะเป็นดังนี้

$$L(y_i, \delta_i) = \begin{cases} f(y_i)G(y_i) & \text{ถ้า } \delta_i = 1 \\ g(y_i)S(y_i) & \text{ถ้า } \delta_i = 0 \end{cases}$$

และมีฟังก์ชันภาวะน่าจะเป็นรวมดังนี้

$$L = \left[ \prod_{i \in u} f(y_i) \prod_{i \in u} G(y_i) \right] \left[ \prod_{i \in c} g(y_i) \prod_{i \in c} S(y_i) \right]$$

เนื่องจาก  $G(y_i)$  และ  $g(y_i)$  ไม่เกี่ยวข้องกับพารามิเตอร์ที่สนใจ จึงละไว้โดยใช้ฟังก์ชันภาวะน่าจะเป็นดังนี้

$$L = \left[ \prod_{i \in u} f(y_i) \prod_{i \in c} S(y_i) \right]$$

$i \in u$  หมายถึงค่าสังเกตที่ไม่ถูกตัดทั้ง

$i \in c$  หมายถึงค่าสังเกตที่ถูกตัดทั้ง

## ๑. อารคณแบบอื่ๆ (Other Type of Censoring)

นอกจากการคักระเภทที่ 1 การคัคประเภทที่ 2 และการคัคแบบตุม ซึ่งที่กล่าวมาเป็นการคัคทางขวา ยังมีการคัคทางซ้าย (Left Censoring) และการคัคทางซ้ายและทางขวา (Left and Right Censoring) ข้อมูลที่ถูกคัคทางซ้าย เช่นกรรมธรรม์ประกันภัชวตอนคัคกำหนดความเสียหายต่วนแรก (Deductible) ที่ผู้เอาประกันภัชวจะคัคองรับคัคชองเอง คัคนัันบริษัทประกันภัชวจึงมีเฉพาะข้อมูลการจ่ายความชชวเสียหายที่มากกว่าความเสียหายต่วนแรกที่กำหนดไว้เท่านั้น และถ้าบริษัทรับประกันภัชวไม่รับประกันภัชวคัคกับบริษัทรับประกันภัชวคัคคัค ซึ่งบริษัทรับประกันภัชวจะกำหนดจำนวนความรับคัคชองสูงสุดของบริษัคเอาไว้ ความชชวเสียหายต่วนที่เกินจากที่กำหนดนี้บริษัทรับประกันภัชวคัคจะเป็นผู้รับคัคชอง คัคนัันข้อมูลทีบริษัทรับประกันภัชวมีอยู่จึงเป็นความชชวเสียหายที่มากกว่าความเสียหายต่วนแรกและไม่เกินจำนวนเงินสูงสุดที่รับคัคชอง

### 2.1.2 ฟัคชันการอยู่รอดและฟัคชันการชชวเสียหาย (Survival Function and Hazard Function)

ให้  $T$  เป็นตัวแปรตุมคัคเนื่อง

$f(t)$  เป็นฟัคชันความหนาแน่นของ  $T$  (Probability density function)

$F(t)$  เป็นฟัคชันการแจกแจงสะสมของ  $T$  (Distribution Function)

$S(t)$  เป็นฟัคชันการอยู่รอดของ  $T$  (Survival Function)

$h(t)$  เป็นฟัคชันความชชวเสียหายหรืออัตราการชชวเสียหาย (Hazard Function

or Hazard Rate)

นิยามฟัคชัน  $S(t)$  คัคความน่าจะเป็นที่ตัวแปรตุม  $T$  จะมีค่ามากกว่าหรือเท่ากับ  $t$

$$\begin{aligned} S(t) &= P(T > t) \\ &= \int_t^{\infty} f(t) dt \end{aligned}$$

นิยามฟัคชัน  $h(t)$  แทนฟัคชันการชชวเสียหายมีค่าเท่ากับลิมิตของความน่าจะเป็นที่ตัวแปรตุม  $T$  จะมีค่าอยู่ในช่วงเวลาคัคๆ  $(t, t + \Delta t)$  คัคหน่วยเวลา  $\Delta t$  เมื่อกำหนดค่า  $T > t$  และ  $h(t)$  กำหนดได้คัคนี้

$$\begin{aligned} h(t) &= \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T < t + \Delta t / T \geq t)}{\Delta t} \\ &= \lim_{\Delta t \rightarrow 0} \frac{[F(t + \Delta t) - F(t)]}{(1 - F(t)) \Delta t} \end{aligned}$$

$$\begin{aligned}
 &= \frac{d}{dt} F(t) \cdot \frac{1}{1-F(t)} \\
 &= \frac{f(t)}{1-F(t)} = \frac{f(t)}{S(t)}, \quad t > 0
 \end{aligned}$$

กรณีเมื่อ  $T$  เป็นตัวแปรสุ่มที่ไม่ต่อเนื่อง และมีค่าเป็น  $t_1, t_2, t_3, \dots$  โดยที่  $0 \leq t_1 < t_2 < t_3 < \dots$  ดังนั้นฟังก์ชันความน่าจะเป็น  $p(t_j)$  (Probability Function) กำหนดได้ดังนี้

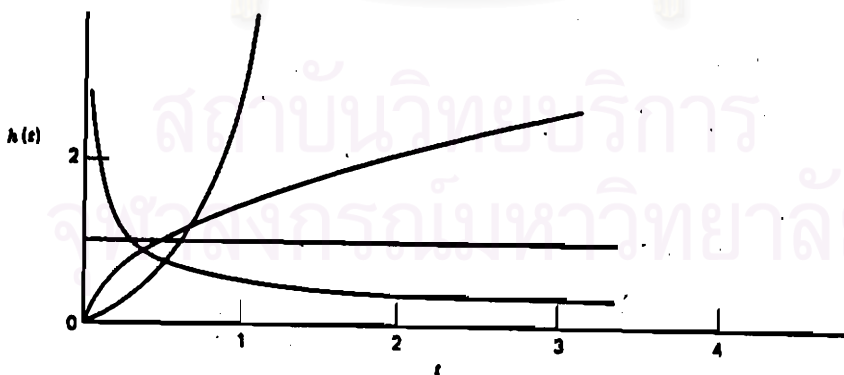
$$p(t_j) = P(T = t_j) \quad j = 1, 2, 3, \dots$$

และฟังก์ชันการอยู่รอด  $S(t)$  คือ

$$S(t) = P(T \geq t_j) = \sum_{j \geq t} p(t_j)$$

ดังนั้นฟังก์ชันความสูญเสีย  $h(t)$  แสดงได้ดังนี้

$$\begin{aligned}
 h(t) &= P(T = t_j / T \geq t_j) \\
 &= \frac{P(t_j)}{S(t_j)}
 \end{aligned}$$

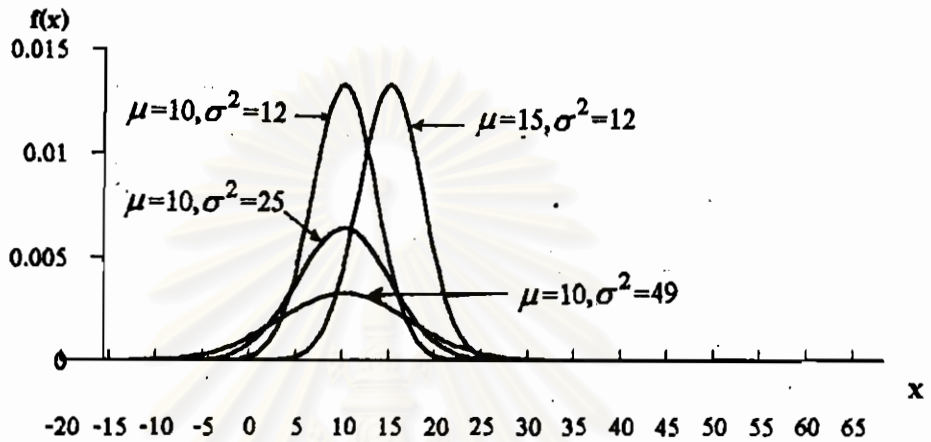


รูปที่ 2.2 แสดงตัวอย่างลักษณะต่างๆ ของฟังก์ชันการสูญเสีย  $h(t)$  ในรูปแบบต่าง ๆ

2.1.3 **ฟังก์ชันความหนาแน่น ค่าคาดหวัง และความแปรปรวน** ของการแจกแจงภายใต้การศึกษานี้

2.1.3.1 การแจกแจงปกติ (Normal Distribution)

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-(x-\mu)^2/2\sigma^2\right] \quad ; \quad -\infty < x < \infty$$



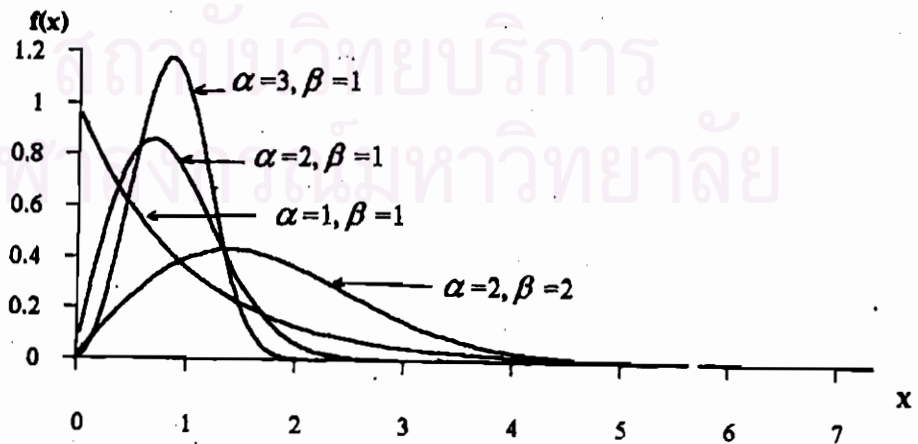
$$E(x) = \mu$$

$$V(x) = \sigma^2$$

รูปที่ 2.3 แสดงการแจกแจงปกติ

2.1.3.2 การแจกแจงแบบไวบูลต์ (Weibull Distribution)

$$f(x) = \alpha\beta^{-\alpha} x^{\alpha-1} \exp\left(-x/\beta\right)^\alpha \quad ; \quad x > 0$$



$$E(x) = \frac{\beta}{\alpha} \Gamma(1/\alpha)$$

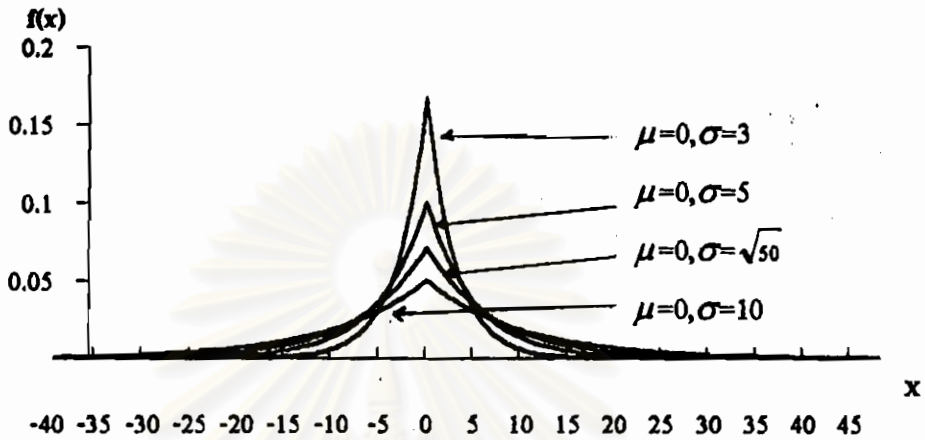
$$V(x) = \frac{\beta^2}{\alpha} \left\{ 2\Gamma\left(\frac{2}{\alpha}\right) - \frac{1}{\alpha} \left(\Gamma\left(\frac{1}{\alpha}\right)\right)^2 \right\}$$

รูปที่ 2.4 แสดงการแจกแจงไวบูลต์

2.1.3.3 การแจกแจงแบบคัมเบิลเอกซ์โพเนนเชียล (Double Exponential

Distribution)

$$f(x) = \frac{1}{2\sigma} \exp(-|x - \mu|/\sigma) \quad ; \quad -\infty < x < \infty$$



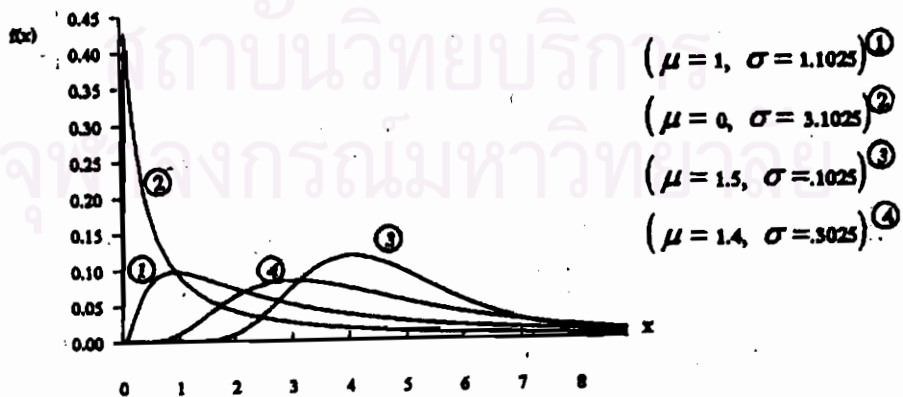
$$E(x) = \mu$$

$$V(x) = 2\sigma^2$$

รูปที่ 2.5 แสดงการแจกแจงแบบคัมเบิลเอกซ์โพเนนเชียล

2.1.3.4 การแจกแจงแบบล็อกนอร์มอล (Lognormal Distribution)

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



$$E(x) = e^{\mu + \sigma^2/2}$$

$$V(x) = e^{2\mu + \sigma^2} (e^{\sigma^2} - 1)$$

รูปที่ 2.6 แสดงการแจกแจงแบบล็อกนอร์มอล



### 2.1.4 ตัวประมาณที่แยก (Product Limit Estimator : PL Estimator)<sup>1</sup>

ตัวประมาณที่แยกพัฒนาขึ้นโดย แคพแลน และ ไมเออร์ เป็นวิธีการประมาณฟังก์ชันการอยู่รอดที่เป็นนอนพารามตริกซ์ (Nonparametric) เมื่อข้อมูลบางส่วนถูกตัดทิ้ง แนวคิดของฟังก์ชันการอยู่รอด  $t_k$  ปี  $S(t_k)$  คือความน่าจะเป็นที่ระยะเวลาของอายุที่รอด (Survival Time) มากกว่า  $t_k$  ปี ดังนั้นจะได้

$$\begin{aligned} S(t_k) &= P(T > t_k) \\ &= P(T > t_1)P(T > t_2 / T > t_1) \dots P(T > t_k / T > t_{k-1}) \\ &= p_1 p_2 p_3 \dots p_k \end{aligned}$$

เมื่อ  $p_i$  คือความน่าจะเป็นที่จะมีชีวิตรอด  $t_i$  ปี หลังจากมีชีวิตอยู่รอดมาแล้ว  $t_{i-1}$  ปี

$$p_i = P(T > t_i / T > t_{i-1})$$

ฟังก์ชันการอยู่รอดมีค่าที่  $S(t_0) = 1$  และ  $S(t_N) = 0$  และเป็นฟังก์ชันขั้นบันได

สมมติค่าสังเกตของอายุที่รอด (Survival Time) จำนวน  $N$  คนมีค่าเป็น  $y_1, y_2, y_3, \dots, y_N$  นำอายุที่รอดมาเรียงตามลำดับเป็น  $y(1) < y(2) < \dots < y(N)$  ดังนั้นหาค่าประมาณ  $\hat{p}_i$  โดย

$$\text{ให้ } \hat{q}_i = \frac{1}{n_i}$$

<sup>1</sup> Kaplan E.L. and Meier P., "Nonparametric Estimation from Incomplete Observations" Journal of the American Statistical Association, 53 (June 1958): 475-81.

Rupert G. Miller, Survival Analysis. (New York : John Wiley and Sons Inc., 1981), PP.46-50

$$\hat{p}_i = 1 - \hat{q}_i = \begin{cases} 1 - \frac{1}{n_i} & \text{เมื่อ } \delta_{(i)} = 1 \\ i & \text{เมื่อ } \delta_{(i)} = 0 \end{cases}$$

$$\delta_i = \begin{cases} 1 & \text{เมื่อข้อมูลลำดับที่ } i \text{ ไม่ถูกตัดทิ้ง} \\ 0 & \text{เมื่อข้อมูลลำดับที่ } i \text{ ถูกตัดทิ้ง} \end{cases}$$

ตัวประมาณที่เอนทสามารถแสดงได้ดังนี้

$$\hat{S}(t) = \prod_{x(t) \leq t} \left( \frac{N-i}{N-i+1} \right)^{\delta_{(i)}}$$

และตัวประมาณที่เอนทเมื่อมีจำนวนค่าสังเกตที่ถูกตัด ณ  $x(t)$  เท่ากับ  $d_i$  ค่า สามารถแสดงได้ดังนี้

$$\hat{S}(t) = \prod_{x(t) \leq t} \left( 1 - \frac{d_i}{n_i} \right)^{\delta_{(i)}}$$

- เมื่อ
- $x(t)$  เป็นค่าสังเกตที่ไม่ถูกตัดทิ้ง ลำดับที่  $i$
  - $i$  เป็นลำดับที่ของข้อมูล
  - $N$  เป็นจำนวนข้อมูลทั้งหมดที่ไม่ถูกตัด และที่ถูกตัด
  - $n_i$  คือจำนวนคนที่มิ่วชีวิตรอด ณ เวลา  $x(t)$ .
  - $d_i$  คือจำนวนคนที่เสียชีวิตรอด ณ เวลา  $x(t)$

ตัวอย่างการหาตัวประมาณที่เอนท (PL Estimator) จากค่าสังเกตต่อไปนี้ 9, 12, 12\*, 18, 23, 28\*, 30, 34, 45, 47, 160\* เมื่อ + หมายถึงข้อมูลถูกตัดทิ้ง จะได้ PL Estimator ดังนี้

$$\begin{aligned} \hat{S}(0) &= 1 \\ \hat{S}(9) &= \hat{S}(0) \times \frac{10}{11} = .91 \\ \hat{S}(12) &= \hat{S}(9) \times \frac{9}{10} = .82 \\ \hat{S}(18) &= \hat{S}(12) \times \frac{7}{8} = .72 \end{aligned}$$

$$\hat{S}(23) - \hat{S}(18) \times \frac{6}{7} = .61$$

$$\hat{S}(30) - \hat{S}(23) \times \frac{4}{5} = .49$$

$$\hat{S}(34) - \hat{S}(30) \times \frac{3}{4} = .37$$

$$\hat{S}(47) - \hat{S}(34) \times \frac{1}{2} = .18$$

### 2.1.3 อีเอ็ม อัลกอริทึม (EM Algorithm)<sup>2</sup>

เนื่องจากสถิติที่เพียงพอ (Sufficient Statistics) เป็นสถิติที่เพียงพอสำหรับพารามิเตอร์ของข้อมูลที่สมบูรณ์เท่านั้น (Complete Data) สำหรับข้อมูลที่ถูกคัดบางส่วนถือว่าเป็นข้อมูลที่ไม่สมบูรณ์ (Incomplete Data) อีเอ็ม อัลกอริทึม เป็นการประมาณค่าที่ประกอบไปด้วย 2 ขั้นตอนคือ ก่อขึ้นการประมาณค่าสังเกตที่ถูกคัดด้วยค่าคาดหวังอย่างมีเงื่อนไขซึ่งเรียกขั้นนี้ว่า ขั้นตอนการหาค่าคาดหวัง (Expectation Step : E Step) ซึ่งเมื่อรวมค่าประมาณของค่าสังเกตที่ถูกคัดและค่าสังเกตที่ไม่ถูกคัดแล้วจะทำให้ได้ข้อมูลที่สมบูรณ์ขึ้น จากขั้นนี้จะทำให้ได้สถิติที่เพียงพอสำหรับพารามิเตอร์ ขั้นต่อมาจะนำข้อมูลทั้งหมดประมาณค่าสูงสุดของพารามิเตอร์ ซึ่งเรียกขั้นนี้ว่า ขั้นตอนการประมาณค่าสูงสุด (Maximization Step : M Step) ซึ่งในการศึกษาครั้งนี้วิธีการน่าจะเป็นสูงสุดใช้ขั้นตอนวิธีอีเอ็มในการประมาณพารามิเตอร์

## 2.2 การประมาณค่าพารามิเตอร์

### 2.2.1 วิธีกำลังสองต่ำสุด

การหาค่าประมาณของพารามิเตอร์วิธีนี้ มีรากฐานมาจากทฤษฎีการประมาณเชิงเส้น (Theory of Linear Estimation) ซึ่งคิดขึ้นโดย คาร์ล ฟูรควิก เกาส์

---

<sup>2</sup>A.P. Dempster, N.M. Laird and D.B. Rubin , "Maximum Likelihood from Incomplete Data via the EM Algorithm " Journal of the Royal Statistical Society, B,39 ( 1977) : 1-38.

และอังเดร แอนควีวีร มาร์คอฟ<sup>3</sup> โดยมีหลักเกณฑ์ว่าค่าประมาณพารามิเตอร์ที่ทำให้ผลบวกกำลังสองของผลต่างระหว่างค่าสังเกตกับค่าประมาณมีค่าต่ำสุด จะเป็นตัวประมาณจริงเต็มที่ดีที่สุดและไม่เอนเอียง (Best Linear Unbiased Estimation, BLUE) เมื่อข้อมูลเป็นไปตามข้อตกลงเบื้องต้นของการวิเคราะห์ความถดถอย คือ

1. ค่าความคลาดเคลื่อนต้องมีการแจกแจงปกติ ที่มีค่าเฉลี่ย 0 และความแปรปรวนเป็น  $\sigma^2$
2. ค่าความคลาดเคลื่อนจะต้องเป็นอิสระต่อกัน คือ  $\varepsilon_i$  และ  $\varepsilon_j$  ไม่มีความสัมพันธ์ต่อกันเมื่อ  $i \neq j, i = 1, 2, 3, \dots, N, j = 1, 2, 3, \dots, N$  เมื่อ  $N$  คือขนาดตัวอย่าง
3. ค่าความคลาดเคลื่อน  $\varepsilon_i$  จะต้องเป็นอิสระกับตัวแปรอิสระ  $X$  หรือ  $Cov(\varepsilon_i, X_i) = 0$  เมื่อ  $i = 1, 2, 3, \dots, N$  เมื่อ  $N$  คือขนาดตัวอย่าง

แต่เนื่องจากข้อมูลที่ถูกคัดเลือกรวมข้อมูลที่ไม่ได้ถูกคัดทิ้ง ดังนั้น วิธีกำลังสองต่ำสุดจะทำให้ได้ตัวประมาณที่เอนเอียง โดยเฉลี่ยจะน้อยกว่าความเป็นจริง

หลักเกณฑ์ในการหาตัวประมาณกำลังสองต่ำสุด จากความสัมพันธ์ระหว่างตัวแปรตาม  $Y$  และตัวแปรอิสระ  $X$  คือ

$$\underline{Y} = \underline{X} \underline{\beta} + \underline{\varepsilon}_i$$

ให้  $\underline{\hat{\beta}} = (\hat{\alpha}, \hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k)$  เป็นเวกเตอร์ของตัวประมาณพารามิเตอร์จะได้ความสัมพันธ์ค่าคาดหวังคือ

$$\underline{\hat{Y}} = \underline{X} \underline{\hat{\beta}}$$

เมื่อ  $\underline{\varepsilon}$  คือความคลาดเคลื่อนระหว่างค่าสังเกตของ  $Y$  กับค่าประมาณ  $\hat{Y}$  ดังนั้น

$$\underline{\varepsilon} = \underline{Y} - \underline{X} \underline{\hat{\beta}}$$

พิจารณาผลบวกกำลังสองของความคลาดเคลื่อน (Sum of Squared Error)

$$\begin{aligned} \underline{\varepsilon}' \underline{\varepsilon} &= (\underline{Y} - \underline{X} \underline{\hat{\beta}})' (\underline{Y} - \underline{X} \underline{\hat{\beta}}) \\ &= (\underline{Y}' - \underline{\hat{\beta}}' \underline{X}') (\underline{Y} - \underline{X} \underline{\hat{\beta}}) \\ &= \underline{Y}' \underline{Y} - 2 \underline{\hat{\beta}}' \underline{X}' \underline{Y} + \underline{\hat{\beta}}' \underline{X} \underline{X}' \underline{\hat{\beta}} \end{aligned}$$

<sup>3</sup> ประจวบ สุวัคฉี, ดร., ทฤษฎีการอนุมานเชิงสถิติ (กรุงเทพมหานคร: 2527), หน้า 158

ตัวประมาณกำลังสองต่ำสุด คือตัวประมาณที่ได้จากการทำให้  $\underline{\underline{\epsilon}}' \underline{\underline{\epsilon}}$  มีค่าต่ำสุด หาได้โดยการหาอนุพันธ์ (Differentiate) เทียบกับ  $\underline{\underline{\beta}}$  แล้วทำให้เท่ากับ 0

$$\frac{\partial}{\partial \underline{\underline{\beta}}} (\underline{\underline{Y}}' \underline{\underline{Y}} - 2 \underline{\underline{\beta}}' \underline{\underline{X}}' \underline{\underline{Y}} + \underline{\underline{\beta}}' \underline{\underline{X}} \underline{\underline{X}}' \underline{\underline{\beta}}) = 0$$

จะได้  $\underline{\underline{\hat{\beta}}} = (\underline{\underline{X}} \underline{\underline{X}}')^{-1} \underline{\underline{X}}' \underline{\underline{Y}}$

### 2.2.2 วิธีการของแชตเตอร์จีและแมคเลิช<sup>4</sup>

เป็นวิธีการประมาณพารามิเตอร์ซึ่งเสนอโดยแชตเตอร์จีและแมคเลิช (S. Chatterjee and D.L. McLeish) ได้พัฒนาขึ้นจากวิธีการประมาณเชิงเส้นเมื่อค่าสังเกตบางค่าสูญหาย ซึ่งคิดขึ้นโดย เยตส์ (Yates) และ ซิเบอร์ (Seber) โดย เยตส์และซิเบอร์ประมาณค่าที่สูญหายด้วยค่าคาดหวังที่ทำให้ผลรวมกำลังสองของความคลาดเคลื่อนมีค่าต่ำสุด วิธีการของแชตเตอร์จีและแมคเลิชเป็นวิธีการนอนพารามตริกซ์ (Nonparametric Method) โดยมีข้อสมมติว่าความคลาดเคลื่อน  $\epsilon_i$  มีการแจกแจงที่เหมือนกันและเป็นอิสระกันด้วยค่าเฉลี่ย เท่ากับ 0 และความแปรปรวน  $\sigma^2$  วิธีการนี้ใช้การวนซ้ำซึ่งในแต่ละรอบจะประมาณค่าคาดหวังของค่าสังเกตที่ถูกตัด  $\hat{Y}_i$  โดยมีหลักการว่าค่าคาดหวังที่ได้จะมีค่ามากกว่าค่าที่ถูกตัด  $T_c$  ดังนี้

$$\hat{Y}_i \begin{cases} T_c & \text{เมื่อ } T_c - \sum_j \hat{\beta}_j X_j^c y_j \geq 0 \\ \sum_j \hat{\beta}_j X_j^c y_j & \text{เมื่อ } T_c - \sum_j \hat{\beta}_j X_j^c y_j < 0 \end{cases}$$

และข้อมูลที่ปามาวิเคราะห์จะถือว่าค่าสังเกตที่ถูกตัดเสมือนเป็นค่าสังเกตที่ไม่ถูกตัด แล้วประมาณพารามิเตอร์  $\underline{\underline{\beta}}$  ด้วยวิธีกำลังสองต่ำสุด

<sup>4</sup> S. Chatterjee and D.L. McLeish. Fitting Linear Regression to Censored Data by Least Squares and Maximum Likelihood Methods. Communications Statistics, (1986), A 15, 11, pp. 3227-3243.

### ขั้นตอนการหาค่าประมาณ $\hat{\beta}$ ด้วยวิธีของเบคเกอร์และแมกลิช

ขั้นที่ 1 หาค่าประมาณ  $\hat{\beta}$  เริ่มต้น โดยการใช้วิธีการกำลังสองต่ำสุดกับเฉพาะค่าสังเกตที่ไม่ถูกตัดทิ้ง

$$\hat{\beta} = (X'X)^{-1}X'Y$$

เมื่อ  $X$  เป็นเมตริกซ์ของตัวแปรอิสระของข้อมูลที่ไม่ถูกตัด  
 $Y$  เป็นเวกเตอร์ของตัวแปรตามของข้อมูลที่ไม่ถูกตัด

ขั้นที่ 2 สำหรับข้อมูลที่ถูกตัด คำนวณค่าประมาณ  $\hat{Y}_i$  ดังนี้

$$\hat{Y}_i = \sum_{j=0}^2 \hat{\beta}_j X_{ij} \quad \text{เมื่อ } X_{i0} = 1$$

เปรียบเทียบค่าประมาณ  $\hat{Y}_i$  ที่ได้กับค่าที่ถูกตัด  $T_c$  ถ้าค่า  $\hat{Y}_i$  มากกว่าค่า  $T_c$  ให้แทนค่าที่ถูกตัดด้วย  $\hat{Y}_i$  ถ้าค่า  $\hat{Y}_i$  น้อยกว่าค่า  $T_c$  ให้คงค่าที่ถูกตัด  $T_c$  ไว้อย่างเดิม

ขั้นที่ 3 นำค่าประมาณที่คำนวณได้จากขั้นที่ 2 และข้อมูลที่ไม่ได้ถูกตัดคำนวณหาตัวประมาณ  $\hat{\beta}$  ด้วยวิธีกำลังสองต่ำสุด

ขั้นที่ 4 แทนค่า  $\hat{\beta}$  ที่ได้จากขั้นที่ 3 ลงในขั้นที่ 2 แล้วทำการวนซ้ำขั้นที่ 2 ถึงขั้นที่ 3 ทำไปจนกระทั่งค่าของ  $\hat{\beta}$  ที่ได้เท่ากับค่าของ  $\hat{\beta}$  ในรอบที่ผ่านมา ก็จะหยุดการวนซ้ำจะได้  $\hat{\beta}$  เป็นค่าประมาณของพารามิเตอร์

ในบางครั้งซึ่งค่า  $\hat{\beta}$  จะแกว่งอยู่ระหว่าง 2 ค่า ในกรณีนี้จะใช้ค่าเฉลี่ยระหว่างสองค่าเป็นค่าประมาณของ  $\hat{\beta}$

### 2.2.3 วิธีการประมาณของบัคเลย์และเจมส์<sup>5</sup>

วิธีการของบัคเลย์และเจมส์เสนอโดยโจนาธานบัคเลย์และไออันเจมส์ (Jonathan Buckley and Ian James) เป็นวิธีการประมาณพารามิเตอร์ในสมการถดถอยเชิงเส้นที่เป็นนอนพารามตริก (Nonparametric Method) ซึ่งวิธีการของบัคเลย์และเจมส์มีข้อสมมติว่า ค่าภาคเคลื่อน  $\varepsilon_i$  เป็นอิสระ และมีฟังก์ชันการแจกแจง  $F$  ที่ไม่มีรูปแบบเฉพาะ มีค่าเฉลี่ยเป็น 0 ความแปรปรวน  $\sigma^2$  และมีฟังก์ชันการอยู่รอดเป็น  $S = 1 - F$

บัคเลย์และเจมส์ได้ทำการประมาณค่าของข้อมูลที่ถูกลบด้วยค่าภาคหวังอย่างมีเงื่อนไข  $E(Y_i/Y_i > T_c, X_i, \beta)$  เมื่อ  $T_c$  เป็นค่าที่ถูกลบ แต่เนื่องจากไม่ทราบฟังก์ชันการแจกแจง  $F$  จึงไม่สามารถหาค่า  $E(Y_i/Y_i > T_c, X_i, \beta)$  จึงประมาณค่า  $F$  ด้วยตัวประมาณที่แอล

ดังนั้นข้อมูลที่ถูกลบจะถูกแทนด้วย

$$\bar{Y}_i(\hat{\beta}) = (\hat{\beta}_1 X_{i1} + \hat{\beta}_2 X_{i2}) + \frac{\sum W_k(\hat{\beta})(Y_k - (\hat{\beta}_1 X_{k1} + \hat{\beta}_2 X_{k2}))}{1 - \hat{F}(T_c - (\hat{\beta}_1 X_{i1} + \hat{\beta}_2 X_{i2}))}$$

$$\text{เมื่อ } \delta_{(i)} = \begin{cases} 1 & \text{เมื่อข้อมูลลำดับที่ } i \text{ ไม่ถูกลบ} \\ 0 & \text{เมื่อข้อมูลลำดับที่ } i \text{ ถูกลบ} \end{cases}$$

ขั้นตอนการหาตัวประมาณ  $\hat{\beta}$  สำหรับวิธีของบัคเลย์และเจมส์

ขั้นที่ 1 หาค่าประมาณ  $\hat{\beta}$  เริ่มต้น โดยการใช้วิธีการกำลังสองค่าสุดท้ายเฉพาะค่าสังเกตที่ไม่ถูกลบทิ้ง

<sup>5</sup> Jonathan Buckley and Ian James. Linear Regression with Censored Data. *Biometrika*. (1979), 66, 3, pp. 429-436

Rupert G. Miller. *Survival Analysis*. (New York : John Wiley, 1981), pp. 150-154

$$\hat{\beta} = (X'X)^{-1}X'Y$$

เมื่อ  $X$  เป็นเมทริกซ์ของตัวแปรอิสระของข้อมูลที่ไม่ถูกตัด  
 $Y$  เป็นเวกเตอร์ของตัวแปรตามของข้อมูลที่ไม่ถูกตัด

ขั้นที่ 2 จากข้อมูลทั้งหมดคือรวมทั้งข้อมูลที่ถูกตัดและข้อมูลที่ไม่ถูกตัด ให้หาความคลาดเคลื่อน  $e$  จาก

$$e_i = Y_i - \hat{\beta}_1 X_{i1} - \hat{\beta}_2 X_{i2} \quad ; \quad i = 1, 2, \dots, N$$

ขั้นที่ 3 ให้เรียงลำดับที่ (Rank) ความคลาดเคลื่อน จากน้อยไปมาก จะได้  $e_{(1)}, e_{(2)}, \dots, e_{(N)}$  ในกรณีที่ลำดับที่ของค่าที่ถูกตัดและค่าที่ไม่ถูกตัดมีค่าเท่ากันให้ลำดับที่ของค่าที่ไม่ถูกตัดนำหน้าค่าที่ถูกตัด

ขั้นที่ 4 ให้เปลี่ยนค่าลำดับที่ของค่าที่ถูกตัดเป็น 0 ส่วนค่าลำดับที่ของค่าที่ไม่ถูกตัดให้คงไว้

ขั้นที่ 5 หาค่า PL Estimator  $\hat{S}$  จาก

$$\hat{S}(e_i, \hat{\beta}) = \prod_{i: e_{(i)} \leq e_1} \left[ \frac{N-i}{N-i+1} \right]^{s(i)}$$

เมื่อ  $i$  คือลำดับที่ของความคลาดเคลื่อน

$N$  คือจำนวนตัวอย่างทั้งหมด

PL Estimator  $\hat{S}$  นี้คือ ค่า Survival Time

ขั้นที่ 6 คำนวณหาฟังก์ชัน  $\hat{F}(e_i, \hat{\beta})$  จาก

$$\hat{F}(e_i, \hat{\beta}) = 1 - \prod_{i: e_{(i)} \leq e_1} \left[ \frac{N-i}{N-i+1} \right]^{s(i)}$$



ขั้นที่ 7 คำนวณค่าถ่วงน้ำหนัก  $w_i(\hat{\beta})$  จาก

$$\begin{aligned} w_1(\hat{\beta}) &= \hat{R}_1 \\ w_2(\hat{\beta}) &= \hat{R}_2 - \hat{R}_1 \\ &\vdots \\ w_N(\hat{\beta}) &= \hat{R}_N - \hat{R}_{N-1} \end{aligned}$$

ในกรณีลำดับที่สูงสุดของความคลาดเคลื่อนเป็นค่าที่ถูกตัดให้ปรับตัวถ่วงน้ำหนัก  $w_i(\hat{\beta})$  ใหม่ดังนี้

$$w^*_i(\hat{\beta}) = \frac{w_i(\hat{\beta})}{\sum_i w_i(\hat{\beta})}$$

เมื่อ  $\sum_i$  เป็นค่าผลบวกเฉพาะกับข้อมูลที่ไม่ถูกตัด

ขั้นที่ 8 ให้เปลี่ยนค่าลำดับที่ของค่าที่ไม่ถูกตัดเป็น 0 ส่วนค่าลำดับที่ของค่าที่ถูกตัดให้คงไว้

ขั้นที่ 9 หากค่า PL Estimator  $\hat{S}$

$$\hat{S}(e_1, \hat{\beta}) = \prod_{i: e_{(i)} \leq e_1} \left[ \frac{N-i}{N-i+1} \right]^{S(i)}$$

ขั้นที่ 10 หาฟังก์ชัน  $\hat{F}(e_1, \hat{\beta})$  จาก

$$\begin{aligned} \hat{F}(e_1, \hat{\beta}) &= 1 - \hat{S}(e_1, \hat{\beta}) \\ \hat{F}(e_1, \hat{\beta}) &= 1 - \prod_{i: e_{(i)} \leq e_1} \left[ \frac{N-i}{N-i+1} \right]^{S(i)} \end{aligned}$$

ขั้นที่ 11 ประมาณค่าที่ถูกตัดด้วย

$$\hat{Y}_i(\hat{\beta}) = (\hat{\beta}_1 X_{i1} + \hat{\beta}_2 X_{i2}) + \frac{\sum_k w_k(\hat{\beta})(Y_k - (\hat{\beta}_1 X_{k1} + \hat{\beta}_2 X_{k2}))}{1 - \hat{F}(T_c - (\hat{\beta}_1 X_{i1} + \hat{\beta}_2 X_{i2}))}$$

ขั้นที่ 12 จากข้อมูลที่ไม่วถูกและค่าประมาณของข้อมูลที่ถูกต้องคำนวณหา  $\hat{\beta}$  จาก

$$\hat{\beta} = (X'X)^{-1}X'Y$$

ขั้นที่ 13 แทนค่า  $\hat{\beta}$  ที่ได้จากขั้นที่ 12 ลงในขั้นที่ 2 แล้วคำนวณจากขั้นที่ 2 จนถึงขั้นที่ 12 คำนวณซ้ำคังนี้จนกระทั่งค่า  $\hat{\beta}$  ที่ได้เท่ากับค่า  $\hat{\beta}$  ในรอบที่ผ่านมาจะหยุด ก็จะได้ค่าประมาณของ  $\hat{\beta}$  ในบางครั้งการประมาณค่า  $\hat{\beta}$  จะแกว่งระหว่าง 2 ค่า ในกรณีนี้จะใช้ค่าเฉลี่ยระหว่าง 2 ค่านั้น เป็นค่าประมาณ  $\hat{\beta}$

ขั้นที่ 14 คำนวณ  $d$  จาก

$$d = \frac{1}{N} \left\{ \sum_u Y_i + \sum_c \bar{Y}_i(\hat{\beta}) \right\} - (\hat{\beta}_1 \bar{X}_1 + \hat{\beta}_2 \bar{X}_2)$$

$\sum_u$  เป็นค่าผลบวกเฉพาะข้อมูลที่ไม่วถูกคัก

$\sum_c$  เป็นค่าผลบวกเฉพาะข้อมูลที่ถูกต้อง

$\bar{X}_1$  เป็นค่าเฉลี่ยของตัวแปรอิสระที่ 1

$\bar{X}_2$  เป็นค่าเฉลี่ยของตัวแปรอิสระที่ 2

#### 2.2.4 วิธีการภาว่น่าจะเป็นสูงสุดด้วยขั้นตอนวิธีเอ็ม<sup>6</sup>

การประมาณพารามิเตอร์ในสมการถดถอยเชิงเส้นเมื่อตัวแปรตามบางค่าถูกคักทางขวาค้วภาว่น่าจะเป็นสูงสุด ซึ่งเสนอโดยเมอร์เรย์ ไอท์กิน (Murray Aitkin) เป็นวิธีการที่อักษ อีเอ็ม อัลกอริทึม (EM Algorithm) ซึ่งเสนอโดยคัมสเตอร์ ลายด์ และ รูบิน

<sup>6</sup> Murry Aitkin, A Note on The Regression Analysis of Censored Data. *Technometrics*. (1981), 23, 2, pp. 161-163.

(Dempster Laird and Rubin) เป็นวิธีที่วนซ้ำ และในแต่ละรอบของการทำซ้ำจะประกอบไปด้วย 2 ขั้นตอนคือขั้นตอนของการประมาณค่าพารามิเตอร์อย่างมีเงื่อนไขของข้อมูลที่ถูกต้อง และขั้นตอนของการประมาณค่าสูงสุดของพารามิเตอร์โดยอาศัยทั้งข้อมูลที่ไม่ถูกต้องและถูกต้อง

จากรูปแบบความถัมพันธ์

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \varepsilon_i$$

เมื่อ  $\varepsilon_i \sim N(0, \sigma^2)$

กำหนดให้  $\mu_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2}$

$$z_i = (y_i - \mu_i) / \sigma$$

$$f(z_i) = \frac{1}{\sqrt{2\pi}} \exp(-z^2/2)$$

$$S(z_i) = 1 - F(z_i) = \int_{z_i}^{\infty} f(t) dt$$

$$h(z_i) = \frac{f(z_i)}{1 - F(z_i)}$$

$$\begin{aligned} \phi(y_i) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(y_i - \mu_i)^2\right) \\ &= \frac{1}{\sigma} f(z_i) \end{aligned}$$

ดังนั้นฟังก์ชันภาวะน่าจะเป็นของค่าที่ไม่ถูกต้องและค่าที่ถูกต้องคือ

$$\begin{aligned} L(\beta, \sigma) &= \prod_{i=1}^n f(y_i) \prod_{i=n+1}^{n+m} S(y_i) \\ &= \frac{1}{\sigma^n} \prod_{i=1}^n f(z_i) \prod_{i=n+1}^{n+m} S(z_i) \end{aligned}$$

เมื่อ  $n$  คือจำนวนตัวอย่างที่ไม่ถูกต้อง

$m$  คือจำนวนตัวอย่างที่ถูกต้อง

และฟังก์ชันลอการิทึมภาวะน่าจะเป็นคือ

$$\begin{aligned} \log L(\beta, \sigma) &= -n \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu_i)^2 + \sum_{i=n+1}^{n+m} \log S\left(\frac{y_i - \mu_i}{\sigma}\right) \\ &= -n \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n \left( y_i^2 - 2y_i \sum_{j=1}^2 \beta_j x_{ij} + \left( \sum_{j=1}^2 \beta_j x_{ij} \right)^2 \right) \\ &\quad + \sum_{i=n+1}^{n+m} \log S\left(\frac{y_i - \mu_i}{\sigma}\right) \end{aligned}$$

หาอนุพันธ์บางส่วน (Partial Derivatives) ของฟังก์ชันน่าจะเป็นเทียบกับ  $\beta_j$  และฟังก์ชันถ้อยภาวะน่าจะเป็นคือ

$$\begin{aligned}\frac{\partial \log L}{\partial \beta_j} &= \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \mu_i) x_{ij} + \frac{1}{\sigma} \sum_{i=n+1}^{n+m} \frac{f(z_i)}{1 - F(z_i)} x_{ij} \\ &= \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \mu_i) x_{ij} + \frac{1}{\sigma} \sum_{i=n+1}^{n+m} h(z_i) x_{ij} \\ &= \frac{1}{\sigma^2} \sum_{i=1}^{n+m} (W_i - \mu_i) x_{ij}\end{aligned}$$

โดยที่  $W_i = \begin{cases} y_i & \text{เมื่อ } i = 1, 2, \dots, n \\ \mu_i + \sigma h(z_i) & \text{เมื่อ } i = n+1, \dots, n+m \end{cases}$

หาอนุพันธ์บางส่วนของฟังก์ชันน่าจะเป็นเทียบกับ  $\sigma$  และให้เท่ากับ 0

$$\frac{\partial \log L}{\partial \sigma} = \frac{-n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (y_i - \mu_i)^2 + \frac{1}{\sigma^2} \sum_{i=n+1}^{n+m} h(z_i) (y_i - \mu_i) = 0$$

$$\begin{aligned}-n + \sum_{i=1}^n h(z_i) z_i + \frac{1}{\sigma^2} \sum_{i=n+1}^{n+m} (y_i - \mu_i)^2 &= 0 \\ \sigma^2 &= \frac{\left( \sum_{i=1}^n (y_i - \mu_i)^2 \right)}{\left( n - \sum_{i=n+1}^{n+m} h(z_i) z_i \right)}\end{aligned}$$

สถิติที่เพียงพอสำหรับพารามิเตอร์ของข้อมูลที่สมบูรณ์คือ  $\sum_{i=1}^{n+m} x_{ij} y_i$  และ  $\sum_{i=1}^{n+m} y_i^2$  และค่าคาดหวังอย่างมีเงื่อนไขคือ

$$E\left(\sum_{i=1}^{n+m} x_{ij} y_i\right) = \sum_{i=1}^n x_{ij} + \sum_{i=n+1}^{n+m} x_{ij} E(y_i / y_i > T_c, \beta, \sigma^2)$$

$$E\left(\sum_{i=1}^{n+m} y_i^2\right) = \sum_{i=1}^n y_i^2 + \sum_{i=n+1}^{n+m} E(y_i^2 / y_i > T_c, \beta, \sigma^2)$$

$$\text{เมื่อ } E(y_i / y_i > T_c, \beta, \sigma^2) = \mu_i + \sigma h(z_i)$$

$E(y_i^2 / y_i > Tc, \beta, \sigma^2) = \mu_i + \sigma^2 + \sigma(y_i + \mu_i)H(z_i)$   
 ดังนั้นสำหรับข้อมูลที่ถูกต้อง หากค่าประมาณ  $\hat{y}_i$  และ  $\hat{y}_i^2$  ได้ดังนี้

$$\hat{y}_i = \hat{\mu}_i + \hat{\sigma}H(z_i) \quad (1)$$

$$\hat{y}_i^2 = \hat{\mu}_i^2 + \hat{\sigma}^2 + \hat{\sigma}(y_i + \hat{\mu}_i)H(z_i) \quad (2)$$

การประมาณค่าสูงสุดของพารามิเตอร์ในขั้นเอ็ม (Maximization Step: M Step) ในรอบที่  $k+1$  จะประมาณพารามิเตอร์  $\hat{\beta}$  ด้วยการแทนค่าที่ถูกต้องด้วย  $\hat{\mu}_i + \hat{\sigma}H(z_i)$  จากรอบที่  $k$  แล้ว

คำนวณหา  $\hat{\beta}$  ด้วยวิธีกำลังสองต่ำสุด และ  $\hat{\sigma}^2$  คำนวณหาโดยใน  $\sum_{i=1}^{n+m} (y_i - \hat{\mu}_i^k)^2$  จะแทน  $y_i$  ด้วย  $\hat{\mu}_i + \hat{\sigma}H(z_i)$  และแทน  $y_i^2$  ด้วย  $\hat{\mu}_i^2 + \hat{\sigma}^2 + \hat{\sigma}(y_i + \hat{\mu}_i)H(z_i)$  ดังนั้น

$$\hat{\sigma}^2_{(k+1)} = \frac{1}{(n+m)} \left[ \sum_1^n (y_i - \hat{\mu}_i^{(k)})^2 + \hat{\sigma}^2_{(k)} \sum_{n+1}^{n+m} (1 + z_i^{(k)} H(z_i^{(k)})) \right]$$

ขั้นตอนการประมาณพารามิเตอร์  $\hat{\beta}$  และ  $\hat{\sigma}^2$  โดยวิธีการวนซ้ำจะเป็น มีขั้นตอนดังนี้

ขั้นที่ 1 ประมาณค่าพารามิเตอร์  $\hat{\beta}, \hat{\sigma}^2$  เริ่มต้นด้วยวิธีกำลังสองต่ำสุด โดยข้อมูลที่นำมาวิเคราะห์จะถือว่าข้อมูลที่ถูกต้องทั้งหมดเหมือนข้อมูลที่ไม่ถูกต้องทั้งค่าประมาณพารามิเตอร์ที่ได้คือ  $\hat{\beta}, \hat{\sigma}^2$

ขั้นที่ 2 เฉพาะข้อมูลที่ถูกต้องทั้ง และใช้พารามิเตอร์จากขั้น 1 คำนวณ

$$\hat{\mu}_i = \hat{\beta}_0 + \sum_{j=1}^2 \hat{\beta}_j x_{ij} \quad ; i = n+1, n+2, \dots, n+m$$

$$z_i = (Tc - \hat{\mu}_i) / \hat{\sigma}$$

$$f(z_i) = \frac{1}{\sqrt{2\pi}} \exp(-z_i^2 / 2)$$

$$S(z_i) = 1 - F(z_i) = \int_z^\infty f(t) dt$$

$$H(z_i) = \frac{f(z_i)}{1 - F(z_i)}$$

ขั้นที่ 3 ประมาณข้อมูลที่ถูกต้องทั้งด้วยค่าคาดหวังอย่างมีเงื่อนไข

$$E(y_i / y_i > Tc, \beta, \sigma^2) = W_i$$

$$W_i = \hat{\mu}_i + \hat{\sigma}H(z_i) \quad ; i = n+1, n+2, \dots, n+m$$

แล้วคำนวณค่า  $y_i^*(\beta)$

$$y_i^*(\beta) = y_i \delta_i + E(y_i / y_i > Tc, \beta, \sigma^2)(1 - \delta_i)$$

$$y_i^*(\beta) = y_i \delta_i + W_i(1 - \delta_i)$$

$$\text{เมื่อ } \delta_i = \begin{cases} 1 & \text{เมื่อข้อมูลไม่ถูกตัดทิ้ง} \\ 0 & \text{เมื่อข้อมูลถูกตัดทิ้ง} \end{cases}$$

ขั้นที่ 4 นำค่า  $y_i^*(\beta)$  จากขั้นที่ 3 เพื่อประมาณค่าพารามิเตอร์  $\beta$  ด้วยวิธีกำลังสองต่ำสุดและ  $\hat{\sigma}^2$  จาก

$$\hat{\sigma}^2 = \left[ \sum_1^n (y_i - \hat{\mu}_i) + \hat{\sigma}^2 \sum_{n+1}^{n+m} (1 + z_i h(z_i)) \right] / (n + m)$$

ขั้นที่ 5 เปรียบเทียบค่าประมาณพารามิเตอร์  $\hat{\beta}$ ,  $\hat{\sigma}^2$  จากขั้นที่ 4 กับขั้นที่ 1 ถ้ายังไม่เท่ากันให้คำนวณรอบต่อไป

ขั้นที่ 6 เฉพาะข้อมูลที่ถูกตัด และพารามิเตอร์  $\hat{\beta}$ ,  $\hat{\sigma}^2$  จากขั้นที่ 4 คำนวณ

$$\hat{\mu}_i = \hat{\beta}_0 + \sum \hat{\beta}_j x_{ij}$$

$$z_i = (Tc - \hat{\mu}_i) / \hat{\sigma}$$

$$f(z_i) = \frac{1}{\sqrt{2\pi}} \exp(-z_i^2 / 2)$$

$$S(z_i) = 1 - F(z_i) = \int_{z_i}^{\infty} f(t) dt$$

$$S(z_i) = \frac{f(z_i)}{1 - F(z_i)}$$

ขั้นที่ 7 ประมาณค่าที่ถูกตัดทิ้งด้วยค่าคาดหวังอย่างมีเงื่อนไข

$$E(y_i / y_i > Tc, \beta, \sigma^2) = W_i$$

$$W_i = \hat{\mu}_i + \hat{\sigma} h(z_i)$$

แล้วคำนวณหา

$$y_i^*(\beta) = y_i \delta_i + W_i(1 - \delta_i)$$

$$\text{เมื่อ } \delta_i = \begin{cases} 1 & \text{เมื่อข้อมูลไม่ถูกตัดทิ้ง} \\ 0 & \text{เมื่อข้อมูลถูกตัดทิ้ง} \end{cases}$$

ขั้นที่ 8 นำค่าประมาณพารามิเตอร์  $\tilde{\beta}, \tilde{\sigma}^2$  จากรอบที่ผ่านมา ( $k$ ) และ  $y_i^*(\beta)$  จากขั้นที่ 7 เพื่อประมาณค่าพารามิเตอร์รอบปัจจุบัน ( $k+1$ ) คือ  $\hat{\beta}^{(k+1)}$  ด้วยวิธีกำลังสองต่ำสุดและ  $\hat{\sigma}^2_{(k+1)}$  จาก

$$\hat{\sigma}^2_{(k+1)} = \left[ \sum_1^n (y_i - \hat{\mu}_i^{(k)}) + \hat{\sigma}^2_{(k)} \sum_{n+1}^{n+m} (1 + \hat{z}_i^{(k)} h(\hat{z}_i^{(k)})) \right] / (n+m)$$

ขั้นที่ 9 เปรียบเทียบค่าประมาณพารามิเตอร์จากขั้นที่ 4 กับขั้นที่ 8 คือเปรียบเทียบ  $\tilde{\beta}$  ขั้นที่ 8 กับ  $\tilde{\beta}$  ในขั้นที่ 4 และ  $\tilde{\sigma}^2$  ในขั้นที่ 8 กับ  $\tilde{\sigma}^2$  ในขั้นที่ 4 ถ้าค่าประมาณของพารามิเตอร์ไม่เท่ากันก็จะนำค่าประมาณพารามิเตอร์ในขั้นที่ 8 แทนค่าลงในขั้นที่ 6 และวนทำซ้ำรอบใหม่จากขั้นที่ 6 ถึงขั้นที่ 9 ถ้าค่าประมาณพารามิเตอร์มีค่าอยู่ระหว่างสองค่าแล้วค่าประมาณพารามิเตอร์คือค่าเฉลี่ยของสองค่านั้น

สถาบันวิทยบริการ  
จุฬาลงกรณ์มหาวิทยาลัย