



หน้าที่ 1

หน้า

1.1 ความเป็นมาและความสำคัญของปัจจุบัน

การใช้เทคนิคการวิเคราะห์ความถดถอยในงานวิจัยส่วนใหญ่ มักจะเลือกใช้วิธีกำลัง ส่องตัวสุด เนื่องจากเป็นวิธีการที่ใช้ง่าย และเป็นวิธีที่มีประสิทธิภาพภายใต้ข้อตกลงเบื้องต้น ของการวิเคราะห์ความถดถอย แต่ในการวิเคราะห์ข้อมูลเพื่อจะนำไปสู่การสรุปผลที่ถูกต้องนั้น สิ่งสำคัญเบื้องต้น ส่วนหนึ่งคือข้อมูลที่นำมาวิเคราะห์ ถ้าหากข้อมูลที่นำมาวิเคราะห์นั้นเป็น ข้อมูลที่ไม่สมบูรณ์ (Incomplete Data) คือ การที่ไม่ทราบค่าสังเกตของตัวแปรที่สนใจ หรือ ค่าสังเกตนั้นมีลักษณะเป็นข้อมูลที่ถูกตัดกึ่ง (Censored Data) ย่อมทำให้เกิดผลกระทบจากการ นำข้อมูลลักษณะดังกล่าวมาใช้เพื่อวิเคราะห์และสรุปผล เช่น จะทำให้ได้ตัวประมาณที่เออนเอียง ออกจากนี้ยังมีผลกระทบอันเนื่องมาจากจำนวนของข้อมูลที่ถูกตัดกึ่ง คือ เมื่อจำนวนข้อมูลที่ถูกตัด กึ่งมีจำนวนมากขึ้น ความแย่ร้ายจากการประมาณค่าอย่างจะลดน้อยลง ดังนี้จะเกิดความคลาดเคลื่อนจากการประมาณมากขึ้น

กรณีข้อมูลถูกตัดกึ่งบางส่วนนั้น คือการที่ไม่ทราบค่าที่แท้จริงต้องบันทึกค่าเท่าที่สังเกต ได้ ตัวอย่างเช่น ข้อมูลความเสียหายทางด้านประกันภัย สมมติกำหนดจำนวนเงินประกันไว้สูงสุด 100,000 บาท เพราะฉะนั้น ในกรณีค่าเสียหายจริงสูงกว่า 100,000 บาท บริษัทจะจ่าย 100,000 บาท และบันทึกตัวเลข 100,000 บาท เราเรียกตัวเลข 100,000 บาท ที่บันทึกไว้นี้ว่าข้อมูลถูกตัดกึ่ง ซึ่งจะเห็นได้ว่าตัวเลขข้อมูลนี้ไม่ใช่ตัวเลขแสดงค่าเสียหายที่เป็นจริง ข้อมูลที่ถูกตัดกึ่งบางส่วนอาจจะเกิดขึ้นได้จากหลายสาเหตุ ซึ่งนอกจากตัวอย่างที่กล่าวมาแล้วอาจมี กรณีอื่น ๆ อีก เช่น ไม่สามารถบันทึกค่าจริงของข้อมูลได้ อันเนื่องมาจากได้ทำการทดลองเพียง ช่วงระยะเวลาหนึ่ง เมื่อยุดทำการทดลองจึงไม่สามารถทราบค่าจริงของข้อมูลได้ ข้อมูลในลักษณะถูกตัดกึ่งนี้พบได้บ่อยในด้านการแพทย์ การทดลองด้านอุตสาหกรรม และการประกันภัย

ลักษณะข้อมูลที่ถูกตัดทิ้งบางส่วนอาจเกิดได้หลายรูปแบบ เช่น จะเกิดข้อมูลถูกตัดทิ้งทางซ้าย (Left Censoring) ที่ C ถ้าไม่ทราบค่าที่แท้จริงของค่าสังเกตเฉพาะค่าที่น้อยกว่า C หรือข้อมูลถูกตัดทิ้งทางขวา (Right Censoring) ที่ C ถ้าไม่ทราบค่าที่แท้จริงของค่าสังเกตเฉพาะค่าที่มากกว่า C

ในการวิจัยครั้งนี้สนใจที่จะศึกษา กรณีที่ตัวแปรตามในสมการถูกตัดออก เชิงเส้นอย่างง่าย มีค่าถูกตัดทิ้งทางขวาเท่านั้น การทดลองหรือการศึกษาที่ทำให้ตัวแปรตามมีค่าถูกตัดทิ้งทางขวา เช่น การทดลองเกี่ยวกับความทนทาน หรืออายุการใช้งานของจนวนกับความร้อนว่าจะช้าอยู่กับอุณหภูมิหรือความร้อนที่ได้รับหรือไม่ ในการทดลองนี้ตัวแปรตาม คือ อายุการใช้งานของจนวน และตัวแปรอิสระ คือ อุณหภูมิหรือความร้อน โดยให้อุณหภูมิหรือความร้อนแยกจนวน แล้วบันทึกเวลา หรือจำนวนชั่วโมงที่จนวนนั้นจะเสื่อมสภาพ ในระหว่างการทดลองจนวนอันที่มีการเสื่อมสภาพ จะเป็นข้อมูลที่ไม่ถูกตัดทิ้ง (Uncensored Data) เมื่อลิสต์การทดลอง จนวนที่ยังคงอยู่ในสภาพใช้งานได้จะเป็นจนวนที่ไม่ทราบอายุการใช้งานที่แท้จริง เพราะฉะนั้นตัวเลขอายุการใช้งานที่บันทึกไว้ เมื่อลิสต์การทดลองจะเป็นข้อมูลที่ถูกตัดทิ้งทางขวา

ตัวอย่างข้อมูลที่ถูกตัดทิ้งทางขวาในการทดลองทางด้านการแพทย์ เช่น ในโครงการทดลองเปลี่ยนหัวใจคน ใช้¹ Stanford Heart Transplantation Program เมื่อคัดเลือกการทดลองว่าจำนวนวันที่มีชีวิตอยู่รอดของคนใช้หลังการผ่าตัดเปลี่ยนหัวใจ ช้าอยู่กับอายุของคนใช้ที่เข้ารับการรักษาหรือไม่ ในการทดลองนี้ตัวแปรตาม คือ จำนวนวันที่มีชีวิตอยู่รอดของคนใช้หลังการผ่าตัด และตัวแปรอิสระ คืออายุของคนใช้ที่เข้ารับการรักษา โดยเริ่มทดลองเมื่อ 1 ตุลาคม ค.ศ. 1967 และลิสต์การทดลอง ในวันที่ 1 เมษายน ค.ศ. 1974 ในระหว่างการทดลองมีคนใช้ 69 คน เข้ารับการเปลี่ยนหัวใจ และสังเกตว่าคนใช้จะมีชีวิตอยู่รอดหลังจากเปลี่ยนหัวใจ แล้ว โดยนับเป็นจำนวนวันที่มีชีวิตอยู่รอด ดังนั้นคนใช้ที่เสียชีวิตในระยะเวลาของการทดลองจะเป็นข้อมูลที่ไม่ถูกตัดทิ้ง แต่คนใช้ที่มีชีวิตอยู่รอด เมื่อลิสต์การทดลองจะไม่ทราบอายุที่มีชีวิตอยู่รอดหลังการเปลี่ยนหัวใจที่แน่นอนได้ จะทราบก็แต่อายุที่อยู่รอด เมื่อลิสต์การทดลองเท่านั้น ข้อมูลล่วงนี้จะเป็นข้อมูลที่ถูกตัดทิ้งทางขวา

¹ Rupert G. Miller. "Least Squares Regression with Censored Data." Biometrika (1976), 63, 3, 456-458

ทางด้านการประกันภัย เช่น การประกันภัยรถยนต์ ข้อมูลค่าเสียหาย (Claim) ให้เป็นตัวแปรตาม สำหรับตัวแปรอิสระ เช่น อายุคนขับ, ขนาดซีซีของรถ, ประเภทรถ และขนาดบรรทุก เป็นต้น ปกติข้อมูลค่าเสียหายทางด้านการประกันภัยจะจัดระดับค่าความเสียหายที่มีหน่วยเป็นบาท และจะกำหนดจำนวนเงินค่าเสียหายที่บริษัทรับประกันสูงสุดไม่เกิน C_k ถ้า T เป็นจำนวนเงินค่าเสียหายจริงมีค่ามากกว่าจำนวนเงินสูงสุดที่ต้องรับผิดชอบตามกรมธรรม์ คือ $T > C_k$ บริษัทรับประกันจะจ่ายเงินค่าเสียหายจำนวน C_k เท่านั้น และบันทึกการจ่ายเงินค่าเสียหายจำนวน C_k ในกรณีนี้บริษัทรับประกันจะไม่ทราบจำนวนเงินค่าเสียหายล่วงหน้าที่เกินกว่าจำนวนเงินที่ต้องรับผิดชอบตามกรมธรรม์ที่แน่นอนได้ ข้อมูลลักษณะ เช่นนี้ เป็นข้อมูลที่ถูกตัดตั้งทางขวา

เนื่องจากข้อมูลล่วงหน้าถูกตัดตั้ง มีความลำบากที่จะต้องนำมาใช้วิเคราะห์เพื่อวางแผน หรือปรับปรุงแผน และเพื่อการตัดสินใจ จึงควรทำการประมาณค่าที่ถูกตัดตั้ง โดยใช้วิธีการที่มีความเหมาะสมจะได้ค่าประมาณที่มีค่าใกล้เคียงกับค่าจริง และมีจำนวนข้อมูลเพียงพอที่จะนำไปศึกษาตามวัตถุประสงค์ เป็นต้นว่า เพื่อสร้างตัวแบบที่เหมาะสม หรือเพื่อศึกษาถึงการแจกแจงของตัวแปรตาม ตัวอย่างทางด้านการประกันภัย เมื่อทำการประมาณจำนวนเงินค่าเสียหาย T ที่มากกว่า C_k จะมีข้อมูลค่าความเสียหายล่วงหน้าที่ได้จากการประมาณของค่าที่ถูกตัดตั้ง ข้อมูลเหล่านี้รวมทั้งข้อมูลล่วงหน้าที่ไม่ถูกตัดตั้ง สามารถนำข้อมูลทั้งหมดมาศึกษารูปแบบการแจกแจงของค่าเสียหาย (Loss Distribution) และอัตราค่าความเสียหาย (Loss Ratio) ซึ่งข้อมูลเหล่านี้จะให้ประโยชน์ในการสร้างตารางอัตราค่าเบี้ยประกันภัย เป็นต้น

ในการวิจัยครั้งนี้จะประมาณค่าพารามิเตอร์ β_i , $i=0,1$ จากสมการถดถอยเชิงเส้นอย่างง่าย โดยใช้วิเคราะห์กับข้อมูลที่ค่าลังเกตของตัวแปรตามมีค่าถูกตัดตั้ง ซึ่งการนำข้อมูลที่มีค่าถูกตัดตั้งมาใช้ประมาณค่าพารามิเตอร์นี้ อาจจะกระทำได้ในลักษณะต่อไปนี้

กรณีแรก จะถือว่าค่าที่ถูกตัดตั้ง เมื่อ้อนเป็นค่าที่ไม่ถูกตัดตั้ง แล้วจึงทำการวิเคราะห์จากข้อมูลทั้งหมด

กรณีที่สอง ไม่สนใจข้อมูลล่วงหน้าที่เป็นค่าถูกตัดตั้ง นั่นคือ จะทำการวิเคราะห์ข้อมูลเฉพาะกับข้อมูลที่ไม่ถูกตัดตั้ง เท่านั้น ดังนั้น จำนวนข้อมูลที่นำมาวิเคราะห์จึงมีขนาดตัวอย่างน้อยกว่ากรณีแรก

จากทั้งสองกรณีดังกล่าว เมื่อประมาณค่าพารามิเตอร์ ด้วยวิธีกำลังสองต่ำสุดจะทำให้ได้ตัวประมาณเที่ยงเฉียง และโดยเฉลี่ยการประมาณค่าจะต่ำกว่าความเป็นจริง หรือจะทำให้ได้ช่วงความเชื่อมั่นแคบกว่าความเป็นจริง² ดังนั้น ในการประมาณค่าพารามิเตอร์ กรณีที่ข้อมูลของตัวแปรตามมีค่าที่ถูกตัดทั้งนั้น ควรที่จะศึกษาวิธีการอื่น ๆ ที่มีประสิทธิภาพมากกว่าวิธีกำลังสองต่ำสุด นั่นคือ สนใจที่จะศึกษาวิธีการประมาณด้วยภาวะน่าจะเป็นสูงสุด และวิธีการของบัคเลร์และเจมส์ และใช้วิเคราะห์กับข้อมูลกรณีแรก คือ จะถือว่าค่าลังเกตที่ถูกตัดทิ้ง เสมือนค่าลังเกตที่ไม่ถูกตัดทิ้ง

1.2 วัตถุประสงค์ของการวิจัย

เพื่อเปรียบเทียบความถูกต้องของ การประมาณค่าตัวแปรตาม ในสมการถดถอยเชิงเส้นอย่างง่าย เมื่อค่าลังเกตของตัวแปรตามเป็นค่าถูกตัดทางขวา โดยประมาณค่าพารามิเตอร์ ด้วยวิธี

- 1.2.1 วิธีกำลังสองต่ำสุด (Ordinary Least Squares Method)
- 1.2.2 วิธีการประมาณด้วยภาวะน่าจะเป็นสูงสุด (Maximum Likelihood Estimation Method)
- 1.2.3 วิธีการของบัคเลร์และเจมส์ (Buckley and James Method)

1.3 สมมติฐานทางการวิจัย

การประมาณค่าตัวแปรตาม ในสมการถดถอยเชิงเส้นอย่างง่าย เมื่อตัวแปรตามมีค่าที่ถูกตัดทิ้งทางขวา ในแต่ละวิธีจะให้ค่าประมาณที่แตกต่างกัน วิธีการประมาณด้วยภาวะน่าจะเป็นสูงสุด จะให้ค่าประมาณตัวแปรตามที่ใกล้เคียงกับค่าจริงมากกว่าวิธีกำลังสองต่ำสุด และวิธีการของบัคเลร์และเจมส์

² Josef Schmee and Gerald J. Hahn. "A Simple Method for Regression Analysis with Censored Data." Technometrics 21(4); (1979): 417-418.

1.4 ข้อตกลงเบื้องต้น

1.4.1 ศึกษารูปแบบสมการถดถอยอย่างง่าย มีรูปแบบดังนี้

$$T_i = \beta_0 + \beta_1 X_i + \varepsilon_i, \quad i = 1, 2, \dots, n+m$$

T_i เป็นตัวแปรตาม

X_i เป็นตัวแปรอิสระ

β_i เป็นพารามิเตอร์ที่ไม่ทราบค่า, $i = 0, 1$

ε_i เป็นค่าความคลาดเคลื่อนสุ่ม

1.4.2 การแจกแจงของค่าที่ถูกตัด และค่าที่ไม่ถูกตัด เป็นอิสระต่อกัน

1.4.3 ตัวแปรตาม (Dependent Variable) เท่านั้นที่เป็นค่าถูกตัดทั้ง

1.4.4 ศึกษาการแข่งpong ประเภทค่าที่ถูกตัดทั้ง เป็นการตัดทางขวา

1.5 ขอบเขตของการวิจัย

1.5.1 ใน การวิจัยครั้งนี้จะทำการเปรียบเทียบการประมาณค่าตัวแปรตาม ในสมการถดถอยเชิงเส้นอย่างง่าย เมื่อค่าลังกอกของตัวแปรตามเป็นค่าถูกตัดทางขวา โดยประมาณค่าพารามิเตอร์ด้วยวิธี

1.5.1.1 วิธีกำลังสองต่ำสุด

1.5.1.2 วิธีการประมาณด้วยภาวะน่าจะเป็นสูงสุด

1.5.1.3 วิธีการของบัคเลร์และเจมส์

1.5.2 ศึกษาเมื่อกรณีของค่าที่ไม่ถูกตัดทั้ง T_i มีรูปแบบเป็น

$$T_i = \beta_0 + \beta_1 X_i + \varepsilon_i, \quad i = 1, \dots, n$$

กำหนดค่าของพารามิเตอร์ ค่าของ X_i และ ε_i เป็นค่าใด ๆ ซึ่งในการวิจัยครั้งนี้กำหนดพารามิเตอร์ $\beta_0 = 2$, $\beta_1 = 1$, X_i มีการแจกแจงแบบปกติ มีค่าเฉลี่ยเป็น 20 มีค่าความแปร

ปัրวนเป็น $60 \text{ } \varepsilon_1$ มีการแจกแจงแบบปกติ มีค่าเฉลี่ยเป็น 0 มีค่าความแปรปรวนเป็น 16 ทุกเงื่อนไขการวิจัย³

1.5.3 ค่าที่ถูกตัดทิ้ง c_1, c_2, \dots เป็นอิสระกัน และมีการแจกแจงเป็น

1.5.3.1 การแจกแจงแบบสม่ำเสมอ (Uniform Distribution)

$U(a, b) :$

$$f(c) = \frac{1}{b-a}, a < c < b$$

1.5.3.2 การแจกแจงแบบปกติตัดปลายทางซ้าย (Left-Truncated Normal Distribution) กำหนดค่า μ, σ^2 เพื่อให้เกิดค่าที่ถูกตัดทิ้งตามขอบเขตที่กำหนด

$$f(c^*) = \frac{f(c)}{1-F(c_o)}, c_o = a$$

$$f(c) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(c-\mu)^2}{2\sigma^2}\right), c > a$$

1.5.3.3 การแจกแจงแบบไวบูลล์ตัดปลายทางซ้าย (Left-Truncated Weibull Distribution) กำหนดค่า α, β เพื่อให้เกิดค่าที่ถูกตัดทิ้งตามขอบเขตที่กำหนด

$$f(c^*) = \frac{f(c)}{1-F(c_o)}, c_o = a$$

$$f(c) = \alpha\beta c^{\beta-1} \exp(-\alpha c^\beta), c > a$$

1.5.3.4 เมื่อค่าถูกตัดทิ้งเป็นฟังก์ชันเชิงเส้น ในรูปแบบเดียวกันกับ T_1 นั้นคือ $C_i = \beta_0 + \beta_1 X_i + \varepsilon_i, i = 1, \dots, m$, m คือ จำนวนค่าลังกาที่ถูกตัดทิ้ง $\beta_0 = 2, \beta_1 = 1, X_i$ มีการแจกแจงแบบปกติ มีค่าเฉลี่ยเป็น 20 มีค่าความแปรปรวนเป็น 60 และ ε_i มีการแจกแจงแบบปกติ มีค่าเฉลี่ยเป็น 0 มีค่าความแปรปรวนเป็น 25

³ ผู้วิจัยได้ทดลองด้วยค่าพารามิเตอร์ ค่าเฉลี่ย และความแปรปรวน หลาย ๆ ค่า พบว่าผลสรุปไม่แตกต่างกัน

1.5.4 ศึกษาเมื่อการเพลิดล่วงของข้อมูลที่ถูกตัดทิ้งเป็น 10%, 20%, 30% และ 40% ของขนาดตัวอย่าง

1.5.5 ศึกษาเมื่อกรณีขนาดตัวอย่างเป็น 5 ระดับคือ 10, 15, 30, 50 และ 70

การวิจัยครั้งนี้จำลองข้อมูลให้มีสถานการณ์ตามที่ต้องการศึกษาโดยใช้เทคนิค การจำลองแบบมอนติ คาร์โล (Monte Carlo Simulation Technique) จากเครื่องคอมพิวเตอร์ AMDAHL 5860 เขียนโปรแกรมด้วยภาษาฟอร์TRAN (FORTRAN 77) ทำการจำลองข้อมูลชั้้ ๆ กัน จำนวน 1,000 ครั้ง ในแต่ละสถานการณ์

1.6 เกณฑ์การตัดสินใจ

เกณฑ์การตัดสินใจว่าการประมาณค่าพารามิเตอร์ด้วยวิธีใดใช้ได้กว่าจะนิยามโดย การเปรียบเทียบค่าความคลาดเคลื่อนระหว่างค่าประมาณของตัวแปรตาม กับค่าจริง ในรูปของ ค่ารากที่สองของค่าเฉลี่ยของความคลาดเคลื่อนกำลังสอง, RMSE (The Square Root of Mean Squares Error) วิธีการได้ให้ค่า RMSE ต่ำกว่าจะเป็นวิธีการประมาณที่ดีกว่า

1.7 ประโยชน์ที่คาดว่าจะได้รับ

1.7.1 เพื่อเป็นแนวทางให้กับวิจัยมีผลสรุปและหลักฐานในการเลือกวิธีการประมาณ ประมาณค่าพารามิเตอร์ในสมการทดถอยเชิงเส้นอย่างง่าย เมื่อค่าสัมเกตของตัวแปรตามมีค่าที่ ถูกตัดทางขวา จะทำให้ได้ค่าประมาณของตัวแปรตามที่ใกล้เคียงกับค่าจริง

1.7.2 เพื่อเป็นแนวทางในการศึกษา และเปรียบเทียบวิธีการประมาณค่าพารามิเตอร์ในสมการทดถอยเชิงเส้นอย่างง่าย เมื่อมีการตัดค่าทางขวาของตัวแปรตาม ในสถานการณ์ อื่น ๆ ต่อไป เช่น เมื่อเป็นสมการทดถอยเชิงเส้นพหุคุณ