

การแก้ปัญหาหุ้สมพันธ์ในการถดถอยโลจิสติก ด้วยตัวประมาณริตจ์โดยวิธีบูตสแตรป์ขั้นตอนเดียว

The Remedy Multicollinearity Problems with a Ridge Logistic Regression Estimator by One-Step Bootstrapping

นพวรรณ มั่นคง (Noppawan Mankong)^{1*} ดร.มานัดฎ์ คำกอง (Dr.Manad Kamkong) **

พุฒิพงษ์ พุกกะมาน (Putipong mBookkamana)***

บทคัดย่อ

การศึกษาครั้งนี้มีวัตถุประสงค์เพื่อสร้างตัวประมาณริตจ์ในการวิเคราะห์การถดถอยโลจิสติกโดยวิธีบูตสแตรป์ขั้นตอนเดียว เมื่อเกิดหุ้สมพันธ์ระหว่างตัวแปรอิสระ โดยเปรียบเทียบวิธีตัวประมาณริตจ์โดยวิธีบูตสแตรป์ขั้นตอนเดียว วิธีความควรจะเป็นสูงสุด วิธีความควรจะเป็นสูงสุดโดยวิธีบูตสแตรป์ขั้นตอนเดียว และวิธีตัวประมาณริตจ์ เกณฑ์ที่ใช้ในการเปรียบเทียบประสิทธิภาพในการประมาณ คือ ค่าความคลาดเคลื่อนกำลังสองเฉลี่ย (MSE) ซึ่งทำการจำลองข้อมูลและการวิเคราะห์ผลโดยใช้โปรแกรม R กำหนดจำนวนตัวแปรอิสระ 4 ตัวแปร ที่มีระดับความสัมพันธ์ 0.4, 0.5, 0.6 และขนาดตัวอย่าง 60, 120, 240 และ 400 ผลการศึกษาสรุปตั้งนี้ในสถานการณ์จำลองทั้งหมดโดยส่วนใหญ่พบว่า การประมาณค่าสัมประสิทธิ์การถดถอยโลจิสติกด้วยตัวประมาณริตจ์โดยวิธีบูตสแตรป์ขั้นตอนเดียวให้ค่าความคลาดเคลื่อนกำลังสองเฉลี่ยต่ำกว่าวิธีการประมาณอื่นๆ เมื่อขนาดตัวอย่างน้อย และเมื่อขนาดตัวอย่างใหญ่ขึ้นประสิทธิภาพในการประมาณของแต่ละวิธีไม่แตกต่างกัน

ABSTRACT

The purpose of this study is to construct a ridge logistic regression estimator using one-step bootstrapping for coping with multicollinearity. The one-step bootstrap ridge estimator method, the maximum likelihood estimator method, the one-step bootstrap maximum likelihood estimator method and the ridge estimator method are compared in term of the mean square error (MSE). The R program is used for the both cases of simulation and generation data. Four independent variables with various both correlation coefficient levels of 0.4, 0.5 and 0.6 and sample size 60, 120, 240 and 400 are studied. The results are shown that the logistic regression coefficients of the one-step bootstrap

¹ Correspondent author: satung_sut@hotmail.com

* นักศึกษา หลักสูตรวิทยาศาสตรมหาบัณฑิต สาขาวิชาสถิติประยุกต์ คณะวิทยาศาสตร์ มหาวิทยาลัยเชียงใหม่

** ผู้ช่วยศาสตราจารย์ ภาควิชาสถิติ คณะวิทยาศาสตร์ มหาวิทยาลัยเชียงใหม่

*** รองศาสตราจารย์ ภาควิชาสถิติ คณะวิทยาศาสตร์ มหาวิทยาลัยเชียงใหม่

ridge estimator method has minimum mean square error when compared with for small sample sizes. For large sample size, however, the efficiency obtained from all methods are not different.

คำสำคัญ : ตัวประมาณริดจ์ พหุสัมพันธ์ วิธีบูตสเตรปซันตอนเดียว

Key Words: Ridge estimator, Multicollinearity, One-step bootstrapping method

บทนำ

การวิเคราะห์การถดถอยโลจิสติก (Logistic regression analysis) เป็นเทคนิคการวิเคราะห์ความสัมพันธ์ของตัวแปรในรูปของการทำนายโอกาสหรือความน่าจะเป็นของการเกิดขึ้นหรือไม่เกิดเหตุการณ์ที่สนใจ สำหรับในกรณีที่ตัวแปรตามมีลักษณะเป็นตัวแปรเชิงกลุ่ม (Categorical variable) กับตัวแปรอิสระอย่างน้อย 1 ตัวแปรเป็นได้ทั้งตัวแปรเชิงปริมาณหรือตัวแปรเชิงกลุ่ม เพื่อไปทำนายโอกาสของการเกิดเหตุการณ์ที่สนใจของตัวแปรตาม เป็นเทคนิคการวิเคราะห์ข้อมูลทางสถิติที่นิยมใช้กันอย่างแพร่หลายในด้านต่างๆ อาทิ ทางด้านเศรษฐศาสตร์ ด้านการเงิน ด้านการตลาด งานวิจัยทางการแพทย์ เป็นต้น เนื่องจากมีการสร้างตัวแบบทางคณิตศาสตร์เพื่อใช้ในการทำนายตัวแปรที่ต้องการศึกษาได้อย่างมีประสิทธิภาพและช่วยในการตัดสินใจได้มากขึ้น ในงานวิจัยส่วนใหญ่นิยมใช้การวิเคราะห์การถดถอยโลจิสติกแบบสองกลุ่ม (Binary logistic regression analysis) คือ ตัวแปรตามจะมีค่าเพียง 2 ค่า คือ 0 กับ 1 โดยที่ 0 แทนเหตุการณ์ที่ไม่สนใจ และ 1 แทนเหตุการณ์ที่สนใจ เช่น ตัวแปรตาม เป็น 0 ถ้าผู้ป่วยไม่เป็นโรคเบาหวาน หรือเป็น 1 ถ้าผู้ป่วยเป็นโรคเบาหวาน ตัวแปรตามเป็น 0 ถ้าลูกค้าไม่ต้องการซื้อสินค้า หรือเป็น 1 ถ้าลูกค้าต้องการซื้อสินค้า เป็นต้น ในการวิเคราะห์การถดถอยโลจิสติกนั้นมีเงื่อนไขเกี่ยวกับตัวแปรอิสระต้องไม่มีความสัมพันธ์กันสูง เนื่องจากอาจก่อให้เกิดปัญหาพหุสัมพันธ์ระหว่างตัวแปรอิสระ (Multicollinearity) เนื่องจากมีตัวแปรอิสระที่ใช้ในการวิเคราะห์มากกว่าหนึ่งตัวแปร

อาจทำให้ $|X'WX|$ มีค่าเข้าใกล้ศูนย์ และทำให้ค่าเฉพาะ (Eigen value) บางค่าของเมทริกซ์ $X'WX$ มีค่าต่ำมาก ซึ่งจะส่งผลให้ค่าความคลาดเคลื่อนกำลังสองเฉลี่ยมีค่าสูง และค่าประมาณสัมประสิทธิ์การถดถอยที่ได้มีความผิดพลาดไปจากค่าจริง ขาดความแม่นยำในการพยากรณ์ของตัวแบบ ซึ่งในทางปฏิบัติเป็นไปได้ยากที่ตัวแปรอิสระจะไม่มีความสัมพันธ์กัน ดังนั้นจึงจำเป็นต้องมีการแก้ไขปัญหาดังกล่าว [1, 2]

โดยวิธีการแก้ปัญหาคือการเกิดพหุสัมพันธ์ระหว่างตัวแปรอิสระมีหลายวิธี อาทิ เมื่อตัวแปรอิสระมีความสัมพันธ์กันสูง ผู้วิจัยสามารถตัดตัวแปรอิสระบางตัวออกจากตัวแบบได้ แต่บางครั้งอาจเป็นเรื่องยากสำหรับนักวิจัยที่จะตัด ตัวแปรอิสระตัวใดตัวหนึ่งออกจากตัวแบบ เนื่องจากตัวแปรอิสระทุกตัวอาจมีผลต่อการเปลี่ยนแปลงของตัวแปรตาม อีกทั้งความสัมพันธ์ระหว่างตัวแปรอิสระยังไม่ชัดเจน วิธีการวิเคราะห์ปัจจัย (Factor analysis) เป็นเทคนิคในการจัดกลุ่มหรือรวมตัวแปรที่มีความสัมพันธ์ไว้ในกลุ่มหรือ Factor เดียวกัน ซึ่งตัวแปรที่อยู่ในกลุ่มเดียวกันจะมีความสัมพันธ์กันมาก ส่วนตัวแปรที่อยู่ในคนละกลุ่มจะมีความสัมพันธ์กันค่อนข้างน้อยมาก หรือไม่มีความสัมพันธ์กัน วิธีการคัดเลือกตัวแปรอิสระสำหรับตัวแบบที่นิยมใช้กันอยู่แพร่หลายในปัจจุบัน ได้แก่ วิธีการคัดเลือกตัวแปรแบบไปข้างหน้า (Forward selection) วิธีการกำจัดตัวแปรแบบถดถอยหลัง (Backward elimination) และวิธีการถดถอยแบบขั้นตอน (Stepwise regression) และในปี ค.ศ.1979 Efron [3] ได้เสนอวิธี

การบูตสแตรป (Bootstrap method) ซึ่งเป็นวิธีการสุ่มตัวอย่างซ้ำเพื่อสร้างตัวอย่างชุดใหม่จากตัวอย่างชุดที่มีเพียงชุดเดียว โดยมีการสุ่มตัวอย่างแบบใส่คืน (Resampling with replacement) ซึ่งตัวอย่างชุดใหม่จะมีขนาดเท่ากับตัวอย่างชุดเดิมที่มีอยู่ ซึ่งวิธีการบูตสแตรปได้ถูกนำมาใช้กันอย่างกว้างขวาง Moulton และ Zeger [4] ได้ประยุกต์ใช้วิธีการบูตสแตรปทั้งในตัวแบบเชิงเส้น (Linear model) และตัวแบบเชิงเส้นวงนัยทั่วไป (Generalized Linear Models : GLM) ซึ่งมีหลักการประมาณค่าสัมประสิทธิ์การถดถอยบูตสแตรปโดยใช้ขั้นตอนแรกหรือการย้อนซ้ำครั้งแรกของกระบวนการกำลังสองน้อยสุดถ่วงน้ำหนักแบบย้อนซ้ำ (Iteratively Weighted Least Squares : IWLS) ในตัวแบบเชิงเส้นมาประยุกต์ใช้กับตัวแบบเชิงเส้นวงนัยทั่วไป จึงเรียกวิธีนี้ว่า วิธีการบูตสแตรปขั้นตอนเดียว นอกจากนี้ ยังมีวิธีการถดถอยริดจ์ (Ridge regression method) ที่ถูกเสนอโดย Hoerl และ Kennard [5] มาใช้ในการแก้ปัญหาความสัมพันธ์ระหว่างตัวแปรอิสระในการประมาณค่าสัมประสิทธิ์ถดถอยเชิงพหุ โดยไม่ต้องตัดตัวแปรอิสระออกจากตัวแบบ โดยมีหลักการในการลดค่าความคลาดเคลื่อนกำลังสองเฉลี่ยของการประมาณค่าสัมประสิทธิ์การถดถอยให้น้อยลง ดังนั้นจึงทำให้ตัวประมาณค่าสัมประสิทธิ์โดยวิธีการถดถอยริดจ์ให้ค่าความคลาดเคลื่อนกำลังสองเฉลี่ยต่ำกว่าวิธีกำลังสองน้อยที่สุด โดย Cessie และ Houwelling [6] ได้ประยุกต์ใช้วิธีการถดถอยริดจ์ ในการวิเคราะห์การถดถอยโลจิสติกด้วยวิธีนิวตันราฟสัน เป็นต้น ดังนั้นผู้วิจัยจึงสนใจศึกษาการ สร้างตัวประมาณริดจ์ในการวิเคราะห์การถดถอยโลจิสติกโดยใช้วิธีการบูตสแตรปขั้นตอนเดียวในการประมาณค่าสัมประสิทธิ์การถดถอยของตัวแบบโลจิสติกในกรณีที่เกิดปัญหาความสัมพันธ์

วัตถุประสงค์ของงานวิจัย

1. เพื่อสร้างตัวประมาณริดจ์ในการวิเคราะห์การถดถอยโลจิสติกโดยวิธีบูตสแตรปขั้นตอนเดียวในการประมาณค่าสัมประสิทธิ์การถดถอยโลจิสติกเมื่อเกิดพหุสัมพันธ์ระหว่างตัวแปรอิสระ
2. เพื่อเปรียบเทียบวิธีการประมาณค่าสัมประสิทธิ์การถดถอยโลจิสติก เมื่อเกิดพหุสัมพันธ์ระหว่างตัวแปรอิสระ ด้วยวิธีดังนี้
 - 2.1 วิธีตัวประมาณริดจ์โดยวิธีบูตสแตรปขั้นตอนเดียว (One-step bootstrapping ridge estimator method)
 - 2.2 วิธีความควรจะเป็นสูงสุด (Maximum likelihood estimator method)
 - 2.3 วิธีความควรจะเป็นสูงสุดโดยวิธีบูตสแตรปขั้นตอนเดียว (One-step bootstrapping maximum likelihood estimator method)
 - 2.4 วิธีตัวประมาณริดจ์ (Ridge estimator method)

วิธีการวิจัย

การวิเคราะห์การถดถอยโลจิสติกแบบสองกลุ่ม

ตัวแปรตามเป็นตัวแปรเชิงกลุ่มซึ่งมีค่าข้อมูลเพียง 2 ค่า คือ 0 กับ 1 กับตัวแปรอิสระที่สามารถเป็นได้ทั้งตัวแปรเชิงกลุ่มและตัวแปรเชิงปริมาณ โดยที่ $y_i = 0$ จะแทนเหตุการณ์ที่ไม่สนใจ ด้วยความน่าจะเป็นที่จะไม่เกิดเหตุการณ์ที่สนใจ คือ $1 - p_i = P(Y_i = 0 | x_{ij})$ และ $y_i = 1$ จะแทนเหตุการณ์ที่สนใจ ด้วยความน่าจะเป็นของการเกิดเหตุการณ์ที่สนใจ คือ $p_i = P(Y_i = 1 | x_{ij})$ ดังนั้นเมื่อนำข้อมูลมาทำการวิเคราะห์ซึ่งอยู่ในรูป (x_{ij}, y_i) จะได้

$$Y_i = \begin{cases} 1 & \text{ด้วยความน่าจะเป็น } p_i \\ 0 & \text{ด้วยความน่าจะเป็น } 1 - p_i \end{cases} \quad (1)$$

โดยให้ p_i แทน ความน่าจะเป็นของการเกิดเหตุการณ์ที่สนใจ

$1-p_i$ แทน ความน่าจะเป็นของการไม่เกิดเหตุการณ์ที่สนใจ

และ $Y_i \sim \text{Bernoulli}(p_i)$ เมื่อ $i = 1, 2, \dots, n$ โดยที่ n คือ ขนาดตัวอย่าง ดังนั้น

$$P(Y_i = y_i) = \begin{cases} p_i^{y_i} (1 - p_i)^{1-y_i} & ; y_i = 0, 1 \\ 0 & ; y_i = \text{ค่าอื่น} \end{cases} \quad (2)$$

โดยที่ y_i แทน ค่าของตัวแปรตามหน่วยที่ i ; $i = 1, 2, \dots, n$

x_{ij} แทน ค่าที่ i ; $i = 1, 2, \dots, n$ ของตัวแปรอิสระที่ j ; $j = 1, 2, \dots, k$

i แทน ขนาดตัวอย่าง $\rightarrow n$

j แทน ตัวแปรอิสระ $\rightarrow k$

Hosmer และ Lemeshow [7] จะสามารถเขียนสมการแสดงความน่าจะเป็นของการเกิดเหตุการณ์ที่สนใจ (Probability of event) ซึ่งเป็นสมการการถดถอยโลจิสติก (3) เนื่องจากตัวแบบการถดถอยโลจิสติกไม่ได้อยู่ในรูปสมการเชิงเส้นตรงดังนั้นสามารถแปลงสมการ (3) ให้อยู่ในรูปเชิงเส้นได้ดังสมการ (4) ดังนี้

$$p_i = P(Y_i = 1 | \mathbf{x}_i) = \frac{e^{\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_j x_{ij} + \varepsilon_i}}{1 + e^{\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_j x_{ij} + \varepsilon_i}} = \frac{e^{\beta' \mathbf{x}_i + \varepsilon_i}}{1 + e^{\beta' \mathbf{x}_i + \varepsilon_i}} \quad (3)$$

$$\text{logit}(p_i) = \ln\left(\frac{p_i}{1 - p_i}\right) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_j x_{ij} + \varepsilon_i \quad (4)$$

โดยที่ β_j คือ ค่าประมาณสัมประสิทธิ์การถดถอยโลจิสติกของตัวแปรอิสระที่ j ; $j = 1, 2, \dots, k$

ε_i คือ ค่าความคลาดเคลื่อนสุ่ม ; $i = 1, 2, \dots, n$

วิธีการถดถอยริดจ์ในการถดถอยโลจิสติก

การวิเคราะห์การถดถอยโลจิสติกเมื่อเกิดปัญหาพหุสัมพันธ์ระหว่างตัวแปรอิสระ ทำการพิจารณาค่าเฉพาะของเมทริกซ์ $\mathbf{X}'\mathbf{W}\mathbf{X}$ พบว่าค่าเฉพาะจะมีค่าน้อยมาก ส่งผลให้ค่าความคลาดเคลื่อนกำลังสองเฉลี่ยสูงขึ้นและทำให้ค่าประมาณสัมประสิทธิ์การถดถอยที่ได้ผิดพลาดไปจากค่าจริงดังที่กล่าวมาข้างต้นฉะนั้นจึงสามารถแก้ปัญหาโดยการบวกค่าคงที่ค่าหนึ่งกับสมาชิกทุกตัวบนเส้นทแยงมุมของเมทริกซ์ $\mathbf{X}'\mathbf{W}\mathbf{X}$ เพื่อให้ค่าเฉพาะมีค่าสูงขึ้น จึงเรียกวิธีการนี้ว่า การถดถอยริดจ์โลจิสติก (Logistic ridge regression)

เมื่อใส่ \ln ลอการิทึมธรรมชาติของฟังก์ชันความควรจะเป็นสูงสุด (Maximum likelihood function) ของ β คือ

$$\ln l(\beta) = \sum_{i=1}^n [y_i (\beta'x_i) - \log(1 + e^{\beta'x_i})] \quad (5)$$

ซึ่ง Duffy และ Santner [8] ได้ทำการพิจารณาค่าสูงสุดของลอการิทึมธรรมชาติของฟังก์ชันความควรจะเป็นเมื่อมีการปรับด้วยขนาดของ β และค่าพารามิเตอร์ c เป็นตัวควบคุมการลดลงของพารามิเตอร์ β ดังนี้

$$\ln l(\beta^*) = \sum_{i=1}^n [y_i (\beta'x_i) - \log(1 + e^{\beta'x_i})] - \frac{1}{2} c \|\beta\|^2 \quad (6)$$

โดย $\|\beta\| = \left(\sum_{j=0}^k \beta_j^2 \right)^{\frac{1}{2}}$ คือ ขนาดของเวกเตอร์พารามิเตอร์ β

กำหนดให้

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix}_{n \times (k+1)}, \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}_{(k+1) \times 1}$$

ดังนั้นค่าสูงสุดของสมการ (6) จะถูกกำหนดโดย $\hat{\beta}^*$ และค่าพารามิเตอร์ c เป็นค่าที่เกิดจากการคำนวณสามารถประมาณ $\hat{\beta}^*$ ด้วยวิธีความควรจะเป็นสูงสุดโดยแก้สมการด้วยกระบวนการวิธีนิวตันราฟสันวิธีหาอนุพันธ์เทียบพารามิเตอร์ $\beta_0, \beta_1, \dots, \beta_k$ ของ $\ln l(\beta^*)$ เมื่อกำหนดให้ β_j แทน $\beta_0, \beta_1, \beta_2, \dots, \beta_k$ ดังนี้

$$\begin{aligned} \frac{\partial \ln l(\beta^*)}{\partial \beta_j} &= \sum_{i=1}^n \left[y_i x_{ij} - \left(\frac{e^{\beta'x_i}}{1 + e^{\beta'x_i}} \right) x_{ij} \right] - c \beta_j \\ &= \mathbf{X}'(\mathbf{y} - \mathbf{p}) - c \mathbf{I} \beta = 0 \end{aligned} \quad (7)$$

และกำหนดให้ β_m แทน $\beta_0, \beta_1, \beta_2, \dots, \beta_k$ หาอนุพันธ์ย่อยอันดับสองเทียบ β_m จะได้

$$\begin{aligned} \frac{\partial^2 \ln l(\beta^*)}{\partial \beta_j \partial \beta_m} &= -\sum_{i=1}^n \left[x_{ij} \left(\frac{(1 + e^{\beta'x_i})(e^{\beta'x_i})x_{im} - (e^{\beta'x_i})(e^{\beta'x_i})x_{im}}{(1 + e^{\beta'x_i})^2} \right) \right] - cI \\ &= -(X'WX + cI) \end{aligned} \tag{8}$$

เมื่อกำหนดให้ $W = \text{diag}(p_i(1 - p_i))$ เป็นเมทริกซ์ทแยงมุมขนาด $n \times n$

ดังนั้นจะได้ตัวประมาณใหม่ด้วยวิธีการถดถอยริดจ์ หรือเรียกว่า ตัวประมาณริดจ์ คือ

$$\hat{\beta}^* = \hat{\beta}_{r+1}^{\text{ridge}} = \hat{\beta}_r + (X'WX + cI)^{-1} [X'(y - p) - cI\hat{\beta}_r]$$

เมื่อ $Z = X\hat{\beta}_r + W^{-1}(y - p)$ จะได้

$$\hat{\beta}_{r+1}^{\text{ridge}} = (X'WX + cI)^{-1} X'WZ \tag{9}$$

เมื่อ c คือ ค่าพารามิเตอร์ริดจ์ สำหรับ $r = 0, 1, 2, \dots$ ซึ่ง $\hat{\beta}_0$ เป็นเวกเตอร์ของตัวประมาณเริ่มต้นที่กำหนดให้เวกเตอร์ของตัวประมาณเริ่มต้นจากเวกเตอร์พารามิเตอร์ของวิธีความควรจะเป็นสูงสุด และถ้าผลต่างระหว่างเวกเตอร์ของพารามิเตอร์ที่ได้ในรอบติดกันมีค่าต่างกันน้อยมากจะถือว่าไม่แตกต่างกัน โดยจะกำหนดเกณฑ์ว่า $|\hat{\beta}_{r+1}^{\text{ridge}} - \hat{\beta}_r^{\text{ridge}}| < 0.000001 \times \mathbf{1}$ เมื่อ $\mathbf{1}$ เป็นเวกเตอร์ ดังนั้น $\hat{\beta}_{r+1}^{\text{ridge}}$ จะเป็นค่าที่ยอมรับได้

โดยที่ $\hat{\beta}_{r+1}^{\text{ridge}}$ คือ เวกเตอร์ของตัวประมาณสัมประสิทธิ์การถดถอยที่คำนวณได้จากรอบที่ $r + 1$

$\hat{\beta}_r^{\text{ridge}}$ คือ เวกเตอร์ของตัวประมาณสัมประสิทธิ์การถดถอยที่คำนวณได้จากรอบที่ r

ตัวประมาณริดจ์โดยวิธีบูตสแตรปขึ้นตอนเดียว

Moulton และ Zeger [4] ได้ประยุกต์ใช้วิธีการบูตสแตรปในการวิเคราะห์การถดถอยโลจิสติกโดยขึ้นอยู่กับการสุ่มซ้ำค่าความคลาดเคลื่อนสำหรับตัวแบบการถดถอยโลจิสติก เรียกว่า One-step residual resampling จากตัวประมาณสัมประสิทธิ์การถดถอยคือ $\hat{\beta}^{\text{ridge}} = (X'WX + cI)^{-1} X'WZ$ ทำการคำนวณค่าความคลาดเคลื่อนเพียร์สันมาตรฐาน (Standardized pearson residuals) ดังนี้

$$r_i = \frac{y_i - \hat{p}_i}{\sqrt{w_i(1 - h_i)}} \tag{10}$$

โดยที่ $w_i = \hat{p}_i(1 - \hat{p}_i)$ เมื่อ $\hat{p}_i = \frac{1}{1 + \exp^{-\hat{\beta}'x_i}}$ และ h_i คือ ค่าของตัวที่ i ที่อยู่บนเส้นทแยงมุมของ

เมทริกซ์ H คือ เมทริกซ์ที่เกิดจากค่าของตัวแปรอิสระเท่านั้น เรียกว่า Hat matrix และทำการปรับค่าเฉลี่ยของค่าความ

คลาดเคลื่อนเพียร์สันมาตรฐาน ได้ดังนี้ $\hat{\varepsilon}_i = r_i - \bar{r}_i$ เมื่อ $\bar{r}_i = \frac{1}{n} \sum_{i=1}^n r_i$

$$H = W^{1/2} X(X'WX)^{-1} X'W^{1/2} \tag{11}$$

เมื่อได้ค่าความคลาดเคลื่อนเพียร์สันมาตรฐานจึงสามารถเข้าขั้นตอนวิธีการบูตสแตรปได้ ดังนี้

เมื่อได้ค่าความคลาดเคลื่อนเพียร์สันมาตรฐานจึงสามารถเข้าขั้นตอนวิธีการบูตสแตรปได้ ดังนี้

ขั้นตอนที่ 1 จาก $\hat{\varepsilon}_1, \hat{\varepsilon}_2, \dots, \hat{\varepsilon}_n$ สร้างชุดตัวอย่างบูตสแตรปขนาด n แบบใส่คืน จะได้ค่าความคลาดเคลื่อน $\hat{\varepsilon}_1^*, \hat{\varepsilon}_2^*, \dots, \hat{\varepsilon}_n^*$ เมื่อ $\hat{\varepsilon}_i^*$ เป็นตัวอย่างที่สุ่มได้ตัวที่ i จาก $\hat{\varepsilon}_1, \hat{\varepsilon}_2, \dots, \hat{\varepsilon}_n$

ขั้นตอนที่ 2 ทำการคำนวณ $\hat{\beta}^{*ridge}$ ซึ่งเป็นตัวประมาณสัมประสิทธิ์การถดถอยโลจิสติกของข้อมูลชุดใหม่ จาก

$$\begin{aligned}\hat{\beta}^{*ridge} &= (\mathbf{X}'\mathbf{W}\mathbf{X} + c\mathbf{I})^{-1} \mathbf{X}'\mathbf{W}\mathbf{Z}^* \\ &= (\mathbf{X}'\mathbf{W}\mathbf{X} + c\mathbf{I})^{-1} \mathbf{X}'\mathbf{W}(\mathbf{X}\hat{\beta} + \mathbf{W}^{-1/2}\hat{\varepsilon}^*)\end{aligned}\quad (12)$$

ขั้นตอนที่ 3 ทำซ้ำในขั้นตอนที่ 1-2 เท่ากับจำนวนที่บูตสแตรป โดยกำหนดให้จำนวนที่บูตสแตรปเป็น B ครั้ง จะได้ $\hat{\beta}^{*ridge(1)}, \hat{\beta}^{*ridge(2)}, \dots, \hat{\beta}^{*ridge(B)}$

ขั้นตอนที่ 4 นำตัวประมาณสัมประสิทธิ์การถดถอย $\hat{\beta}^{*ridge(1)}, \hat{\beta}^{*ridge(2)}, \dots, \hat{\beta}^{*ridge(B)}$ มาหาค่าเฉลี่ย จะได้ตัวประมาณสัมประสิทธิ์การถดถอยโลจิสติกด้วยวิธีบูตสแตรป จาก $\bar{\beta}^{*ridge} = \frac{1}{B} \sum_{b=1}^B \hat{\beta}_b^{*ridge}$ หรือเรียกว่า ตัวประมาณริจด์โดยวิธีบูตสแตรปขั้นตอนเดียว

วิธีดำเนินการวิจัย

ในงานวิจัยครั้งนี้ทำการศึกษาวิธีการประมาณค่าสัมประสิทธิ์การถดถอยโลจิสติก (β) แบบสองกลุ่ม เมื่อเกิดพหุสัมพันธ์ระหว่างตัวแปรอิสระ โดยสร้างสถานการณ์ที่ต้องการศึกษาภายใต้เงื่อนไขดังนี้คือ จำนวนตัวแปรอิสระ (k) ที่ใช้ในการศึกษามีทั้งหมด 4 ตัวแปร ขนาดตัวอย่าง (n) ที่ใช้ในการศึกษา คือ 60, 120, 240 และ 400 โดยกำหนด (β) เริ่มต้น ตามแนวคิดของ Newhouse และ Oman [9] โดยให้ $\beta_0 = 0$ และจะเลือกใช้เวกเตอร์เฉพาะ (Eigen vector) ที่สอดคล้องกับค่าเฉพาะ (Eigen value) ที่มีค่าน้อยที่สุดของเมทริกซ์ $\mathbf{X}'\mathbf{X}$ ที่ให้ $\beta'\beta = 1$ และสร้างตัวแปรอิสระให้มีรูปแบบความสัมพันธ์ตามที่ต้องการ ดังต่อไปนี้

$$x_{ij} = \left(\frac{(1-\rho)}{\rho} \right)^{\frac{1}{2}} z_{ij} + z_{i(k+1)} ; i = 1, 2, \dots, n ; j = 1, 2, \dots, k ; j \neq k \quad (13)$$

เมื่อ z_{ij} เป็นตัวเลขสุ่มที่มีการแจกแจงปรกติมาตรฐาน ที่มีค่าเฉลี่ยเท่ากับศูนย์ ค่าความแปรปรวนเท่า กับหนึ่ง และเป็นอิสระต่อกัน ค่าสัมประสิทธิ์สหสัมพันธ์ระหว่างตัวแปรอิสระ (ρ) ที่ศึกษามีระดับความสัมพันธ์ คือ $\rho = 0.40, \rho = 0.50, \rho = 0.60$

การสร้างค่าพารามิเตอร์ริจด์ (c) จะสร้างตามแนวคิดของ Kibria และคณะ [10] ที่ได้เสนอให้ใช้ค่า c ทั้ง 2 วิธี เมื่อเกิดปัญหาพหุสัมพันธ์ระหว่างตัวแปรอิสระในระดับต่ำและระดับสูง ดังนี้

$$c_1 = \prod_{j=1}^{k+1} \left(\frac{1}{q_j} \right)^{\frac{1}{k+1}} \tag{14}$$

$$c_2 = \text{median} \left(\frac{1}{q_j} \right) \tag{15}$$

เมื่อ $m_j = \sqrt{\frac{\hat{\sigma}^2}{\hat{\alpha}_j^2}}$, $q_j = \frac{\lambda_{\max}}{(n-k)\hat{\sigma}^2 + \lambda_{\max}\hat{\alpha}_j^2}$, $\hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - \hat{p}_i)^2}{n-k-1}$, $\hat{\alpha}_j^2$ คือ สมาชิกตัวที่ j ของ

เมทริกซ์ $\gamma\mathbf{B}$ เมื่อกำหนดให้ γ เป็นเมทริกซ์เชิงตั้งฉากขนาด $(k+1) \times (k+1)$ ซึ่งแต่ละแนวตั้งฉากของเมทริกซ์คือเวกเตอร์เฉพาะที่มีค่าสอดคล้องกับค่าเฉพาะของเมทริกซ์ $\mathbf{X}'\mathbf{W}\mathbf{X}$ และ $\mathbf{X}'\mathbf{W}\mathbf{X} = \gamma'\Lambda\gamma$ เมื่อ $\Lambda = \text{diag}(\lambda_j)$ และ λ_j คือ ค่าเฉพาะตัวที่ j ของเมทริกซ์ $\mathbf{X}'\mathbf{W}\mathbf{X}$ สร้างตัวแปรตาม (y_i) ซึ่งมีการแจกแจงแบร์นูลลีที่มีรูปแบบความสัมพันธ์ดังสมการ (3) และแปลงให้อยู่ในรูปเชิงเส้นได้ดังสมการ (4)

งานวิจัยครั้งนี้ใช้การจำลองแบบเทคนิคมอนติคาร์โล (Monte carlo simulation) กระทำซ้ำ จำนวน 1000 ครั้ง และมีการกระทำซ้ำของวิธีการบูตสแตรป 500 ครั้ง ในแต่ละสถานการณ์ของการจำลอง โดยใช้โปรแกรม R 3.1.1 [11] แล้ว ทำการประมาณค่าสัมประสิทธิ์การถดถอยโลจิสติกทั้ง 4 วิธี คือ 1) วิธีตัวประมาณริคต์โดยวิธีบูตสแตรปขั้นตอนเดียว 2) วิธีความควรจะเป็นสูงสุด 3) วิธีความควรจะเป็นสูงสุดโดยวิธีบูตสแตรปขั้นตอนเดียว 4) วิธีตัวประมาณริคต์จากนั้นคำนวณเกณฑ์ที่ใช้ในการเปรียบเทียบความแม่นยำของการประมาณค่าสัมประสิทธิ์การถดถอยโลจิสติกของแต่ละวิธีคือ ค่าความคลาดเคลื่อนกำลังสองเฉลี่ย (Mean Squares Error: MSE) สามารถหาได้จาก

$$MSE(\hat{\beta}_r) = \frac{1}{1000} \sum_{r=1}^{1000} (\hat{\beta}_r - \beta)'(\hat{\beta}_r - \beta)' \tag{16}$$

โดยที่ r คือ การทำซ้ำในการจำลอง $r = 1, 2, \dots, 1000$, $\hat{\beta}_r$ คือเวกเตอร์ของตัวประมาณสัมประสิทธิ์การถดถอยโลจิสติกที่ได้ในแต่ละวิธีของการทำซ้ำ r ครั้ง ในการจำลองของแต่ละสถานการณ์ และ β คือ เวกเตอร์ของค่าสัมประสิทธิ์การถดถอยโลจิสติกที่กำหนดไว้เป็นค่าคงที่ใดๆ ตามแนวคิดของ Newhouse และ Oman [9]

ผลการวิจัย

จากตารางที่ 1 กรณีที่ระดับความสัมพันธ์เท่ากับ 0.40 ที่ขนาดตัวอย่างเท่ากับ 60, 120 และ 240 พบว่า ตัวประมาณริตจโดยวิธีบูตสแตรปชั้นตอนเดียวจะให้ค่า MSE ต่ำกว่าวิธีอื่น โดยเฉพาะค่าพารามิเตอร์ริตจ c_1 จะให้ค่า MSE ต่ำที่สุด เมื่อขนาดตัวอย่างเท่ากับ 400 พบว่า ค่า MSE มีค่าใกล้เคียงกัน โดยที่ตัวประมาณริตจให้ค่า MSE ต่ำที่สุด เมื่อ c_1

กรณีที่ระดับความสัมพันธ์เท่ากับ 0.50 ที่ขนาดตัวอย่างเท่ากับ 60, 120 และ 240 พบว่า ตัวประมาณริตจ โดยวิธีบูตสแตรปชั้นตอนเดียวจะให้ค่า MSE ต่ำกว่าวิธีอื่น โดยเฉพาะค่าพารามิเตอร์ริตจ c_1 จะให้ค่า MSE ต่ำที่สุด เมื่อขนาดตัวอย่างเท่ากับ 400 พบว่า ค่า MSE มีค่าใกล้เคียงกัน โดยที่ตัวประมาณริตจให้ค่า MSE ต่ำที่สุด

กรณีที่ระดับความสัมพันธ์เท่ากับ 0.60 ที่ขนาดตัวอย่างเท่ากับ 60, 120 และ 240 พบว่า ตัวประมาณริตจโดยวิธีบูตสแตรปชั้นตอนเดียวจะให้ค่า MSE ต่ำกว่าวิธีอื่น โดยเฉพาะค่าพารามิเตอร์ริตจ c_1 จะให้ค่า MSE ต่ำที่สุด เมื่อขนาดตัวอย่างเท่ากับ 400 พบว่า ค่า MSE มีค่าใกล้เคียงกัน โดยที่ตัวประมาณริตจให้ค่า MSE ต่ำที่สุด เมื่อ c_2

จากภาพที่ 1-3 พบว่าขนาดตัวอย่างเพิ่มขึ้น ทำให้ค่า MSE ลดลง เมื่อตัวอย่างขนาดใหญ่ การประมาณค่าแต่ละ วิธีจะให้ค่า MSE ใกล้เคียงกันมาก เมื่อระดับความสัมพันธ์เพิ่มขึ้น ทำให้ค่า MSE ที่ได้จะมีค่าเพิ่มขึ้น และจากการจำลองสถานการณ์ส่วนใหญ่การประมาณค่าสัมประสิทธิ์การถดถอยโลจิสติกด้วยตัวประมาณริตจโดยวิธีบูตสแตรปชั้นตอนเดียวจะให้ค่า MSE ต่ำที่สุด โดยส่วนใหญ่ประมาณค่าพารามิเตอร์ริตจ c ด้วยวิธี c_1 จะให้ค่า MSE ต่ำที่สุด และวิธีรองลงมา และยังพบว่าตัวประมาณริตจ จะให้ค่า MSE ต่ำกว่าวิธีความควรจะเป็นสูงสุดโดยวิธีบูตสแตรปชั้นตอนเดียว และวิธีความควรจะเป็นสูงสุด

สรุปผลการวิจัย

การจำลองข้อมูลและวิเคราะห์ผล กรณีที่ตัวแปรอิสระมีทั้งหมด 4 ตัวแปร ขนาดตัวอย่างที่ใช้ในการศึกษา คือ 60, 120, 240 และ 400 และระดับความสัมพันธ์ของตัวแปรอิสระที่ระดับ 0.40, 0.50 และ 0.60 สามารถสรุปผลได้ดังนี้

ภายใต้ระดับความสัมพันธ์เดียวกันเมื่อพิจารณาขนาดตัวอย่าง $n \leq 120$ พบว่า วิธีตัวประมาณริตจโดยวิธีบูตสแตรปชั้นตอนเดียวทั้งวิธีประมาณค่าพารามิเตอร์ริตจ c_1 และ c_2 จะให้ค่า MSE ที่ต่ำกว่าวิธีการประมาณอื่น ๆ แต่เมื่อขนาดตัวอย่างเพิ่มขึ้น $n > 240$ วิธีการประมาณค่าสัมประสิทธิ์การถดถอยโลจิสติกทั้ง 4 วิธีให้ค่า MSE ที่ต่ำ และไม่แตกต่างกัน และเมื่อระดับความสัมพันธ์เพิ่มขึ้นภายใต้ขนาดตัวอย่างที่เท่ากัน พบว่าค่า MSE มีแนวโน้มเพิ่มสูงขึ้นด้วย

ดังนั้นจากการจำลองสถานการณ์ทั้งหมด โดยส่วนใหญ่การประมาณค่าสัมประสิทธิ์การถดถอยโลจิสติกด้วยตัวประมาณริตจโดยวิธีบูตสแตรปชั้นตอนเดียวจะให้ค่า MSE ต่ำกว่าวิธีการประมาณอื่น ๆ เมื่อขนาดตัวอย่างน้อย โดยเฉพาะวิธีประมาณค่าพารามิเตอร์ริตจ c_1 แต่เมื่อขนาดตัวอย่างใหญ่ขึ้นประสิทธิภาพในการประมาณไม่แตกต่างกัน

ข้อเสนอแนะ

จากการศึกษาครั้งนี้เมื่อพิจารณาระดับความสัมพันธ์ของตัวแปรอิสระในทุกระดับความสัมพันธ์ พบว่าเมื่อขนาดตัวอย่างใหญ่ประสิทธิภาพในการประมาณไม่แตกต่างกัน ดังนั้นเมื่อเกิดพหุสัมพันธ์ระหว่างตัวแปรอิสระ วิธีการประมาณค่าสัมประสิทธิ์การถดถอยโลจิสติกทุกวิธีสามารถนำมาใช้ในการประมาณค่าพารามิเตอร์ได้

เมื่อพิจารณาระดับความสัมพันธ์ของตัวแปรอิสระในทุกระดับความสัมพันธ์กรณีที่ขนาดตัวอย่างเล็กพบว่า วิธีตัวประมาณริตจโดยวิธีบูตสแตรปชั้นตอนเดียวจะมีประสิทธิภาพในการประมาณที่ดีกว่าวิธีอื่น ๆ

ดังนั้นจึงเสนอ แนะนำให้ใช้วิธีตัวประมาณริตจโดยวิธีบูตสเตรปชั่นตอนเดียวในการประมาณค่าสัมประสิทธิ์การถดถอยโลจิสติกกรณี ที่เกิดพหุสัมพันธ์ระหว่างตัวแปรอิสระเมื่อขนาดตัวอย่างเล็ก

กิตติกรรมประกาศ

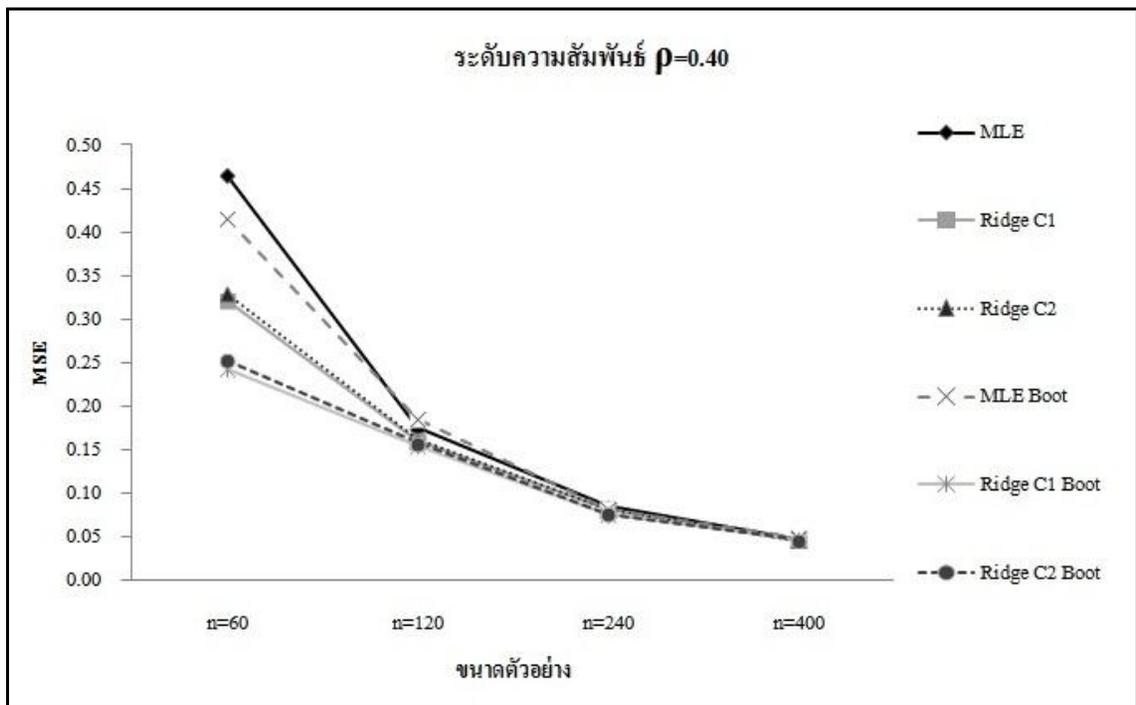
งานวิจัยฉบับนี้สำเร็จลุล่วงด้วยดี เนื่องจากได้รับความช่วยเหลือจากอาจารย์ที่ปรึกษาวิทยานิพนธ์ ผู้วิจัยขอกราบขอบพระคุณท่านอาจารย์เป็นอย่างสูง ที่คอยให้คำปรึกษา แนะนำเกี่ยวกับงานวิจัย คอยตรวจสอบ แก้ไขงานวิจัยฉบับนี้ และขอขอบพระคุณบัณฑิตวิทยาลัย มหาวิทยาลัยเชียงใหม่ที่ให้การสนับสนุน อีกทั้งต้องขอขอบพระคุณผู้ทรงคุณวุฒิที่ช่วยตรวจสอบ แก้ไขงานวิจัยให้สมบูรณ์ยิ่งขึ้น สุดท้ายนี้ผู้วิจัยขอกราบขอบพระคุณ บิดา มารดา ญาติพี่น้อง และเพื่อน ๆ ที่คอยให้กำลังใจเสมอมา

เอกสารอ้างอิง

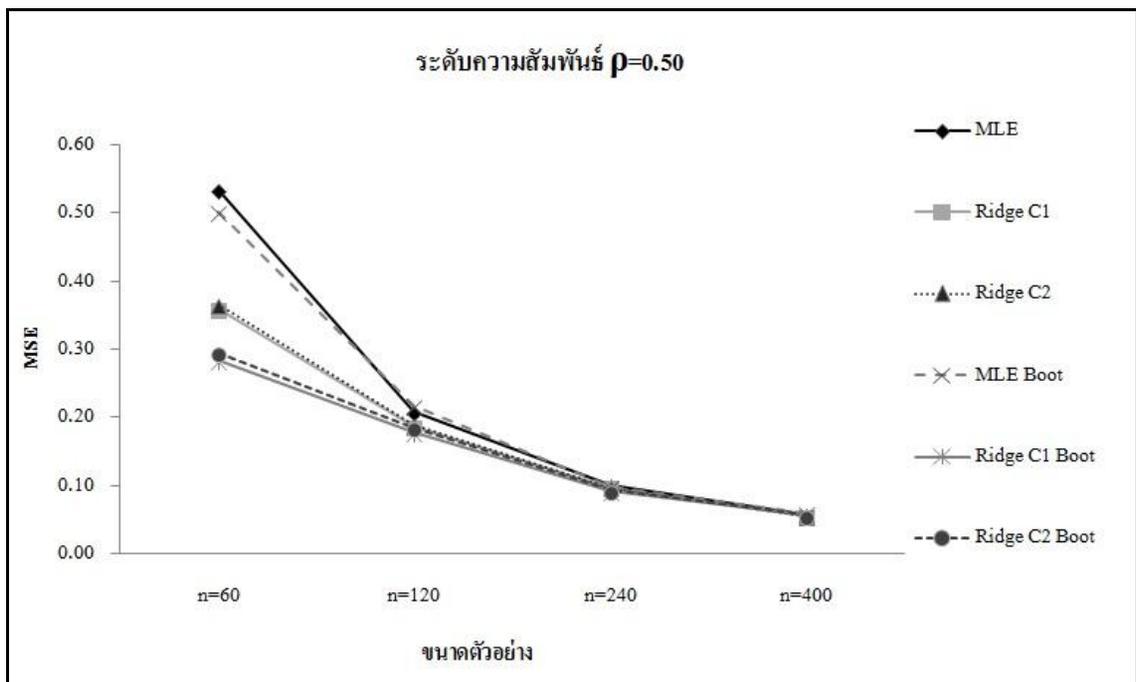
1. Palasri W, Durangwaana S. Estimation of multiple regression coefficients with multicollinearity by Ridge regression bootstrapping method. Proceedings of the 11th graduate research conference; 2010 KKU, Khon Kaen, Thailand.
2. Boonpat S, Vanichbuncha K. Optimal Ridge parameter for solving multicollinearity problem in binary logistic regression. Proceedings of national operations research conference; 2011 Sep 8-9; Bangkok, Thailand.
3. Efron B. Bootstrap Methods Another look at the Jackknife. *Annals of Statistics*. 1979; 7: 1-26.
4. Moulton LH, Zeger SL. Bootstrapping generalized linear models. *Computational Statistics and Data Analysis*. 1991; 11: 53-63.
5. Hoerl AE, Kennard RW. Ridge regression: biased estimation for non orthogonal Problems. *Technometrics*. 1970; 12: 69-82.
6. Cessie LS, Houwelingen VJC. Ridge estimators in logistic regression. *applied statistics - Journal of the Royal Statistical Society Series*. 1992; 41: 191-201.
7. Hosmer DW, Lemeshow S. *Applied logistic regression*. 2nd ed. New York. 2000.
8. Duffy DE, Santner TJ. On the small sample properties of norm-restricted maximum likelihood estimators for logistic regression model. *Communications in Statistics - Theory and Methods*. 1989; 18: 959-980.
9. Newhouse JP, Oman SD. An Evaluation of Ridge estimators. The Rand Corporation. 1971; P-716-PR: 1-28. Hosmer DW, Lemeshow S. *Applied logistic regression*. 2nd ed. New York. 2000.
10. Kibria BMG, Mansson K, Shukur G. Performance of some logistic ridge regression estimators.
11. The comprehensive R archive network (CRAN) [Internet] 2008 [update 2014 Jul 10; cited 2014 Jul 23]. Available from: <http://www.cran.r-project.org/index.html>.

ตารางที่ 1 แสดงค่าความคลาดเคลื่อนกำลังสองเฉลี่ย (MSE)

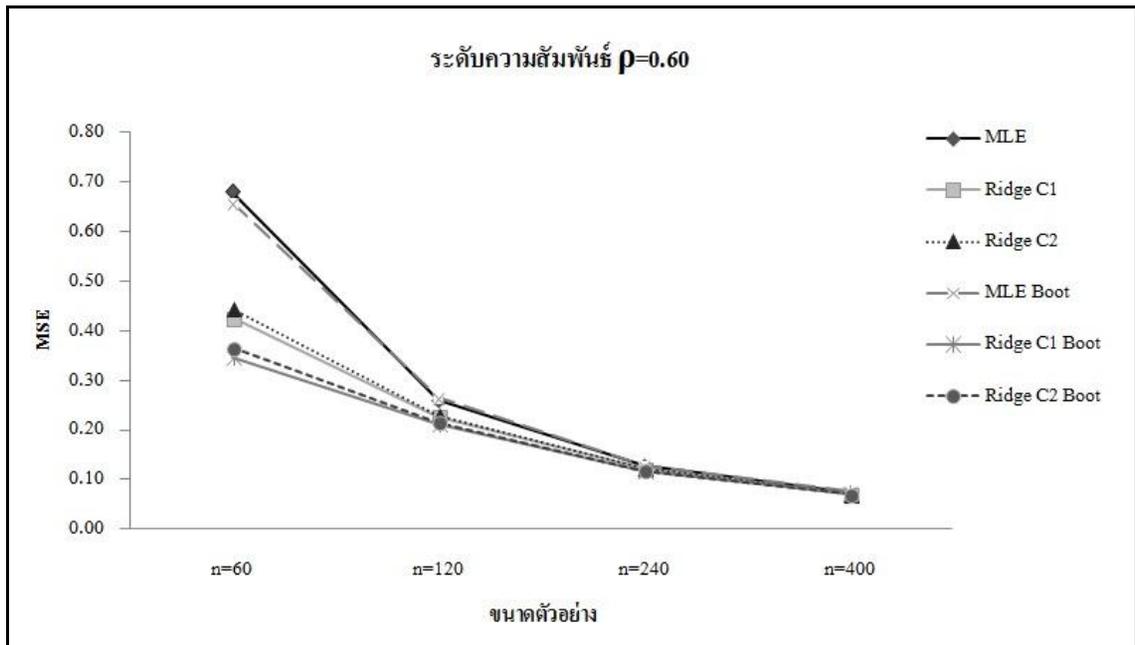
ค่าสัมประสิทธิ์ สหสัมพันธ์ (ρ)	ขนาด ตัวอย่าง (n)	MLE	Ridge C1	Ridge C2	MLE Boot	Ridge C ₁ Boot	Ridge C ₂ Boot
$\rho = 0.40$	60	0.4661	0.3209	0.3286	0.4156	0.2421	0.2523
	120	0.1762	0.1599	0.1613	0.1850	0.1543	0.1578
	240	0.0844	0.0810	0.0810	0.0820	0.0766	0.0767
	400	0.0468	0.0457	0.0458	0.0489	0.0467	0.0468
$\rho = 0.50$	60	0.5334	0.3585	0.3656	0.5025	0.2846	0.2950
	120	0.2081	0.1867	0.1890	0.2174	0.1793	0.1843
	240	0.1014	0.0968	0.0975	0.0994	0.0923	0.0934
	400	0.0572	0.0557	0.0557	0.0598	0.0569	0.0569
$\rho = 0.60$	60	0.6792	0.4253	0.4447	0.6575	0.3473	0.3651
	120	0.2602	0.2273	0.2298	0.2650	0.2109	0.2148
	240	0.1290	0.1221	0.1232	0.1275	0.1169	0.1184
	400	0.0729	0.0706	0.0702	0.0761	0.0715	0.0708



ภาพที่ 1 กราฟแสดงการเปรียบเทียบค่า MSE กรณีที่ระดับความสัมพันธ์เท่ากับ 0.40



ภาพที่ 2 กราฟแสดงการเปรียบเทียบค่า MSE กรณีที่ระดับความสัมพันธ์เท่ากับ 0.50



ภาพที่ 3 กราฟแสดงการเปรียบเทียบค่า MSE กรณีที่ระดับความสัมพันธ์เท่ากับ 0.60

โดยที่	MLE	หมายถึง วิธีความควรจะเป็นสูงสุด
	Ridge C_1	หมายถึง วิธีตัวประมาณริตจ์ โดยใช้ค่า c_1
	Ridge C_2	หมายถึง วิธีตัวประมาณริตจ์ โดยใช้ค่า c_2
	MLE Boot	หมายถึง วิธีความควรจะเป็นสูงสุดโดยวิธีบูตสแตรปซันตอนเดียว
	Ridge C_1 Boot	หมายถึง วิธีตัวประมาณริตจ์โดยวิธีบูตสแตรปซันตอนเดียว โดยใช้ค่า c_1
	Ridge C_2 Boot	หมายถึง วิธีตัวประมาณริตจ์โดยวิธีบูตสแตรปซันตอนเดียว โดยใช้ c_2