

A Comparison of Influence Diagnostics in Linear Mixed Models¹⁾

Jang-Taek Lee²⁾

Abstract

Standard estimation methods for linear mixed models are sensitive to influential observations. However, tools and concepts for linear mixed model diagnostics are rudimentary until now and research is heavily demanded in linear mixed models.

In this paper, we consider two diagnostics to evaluate the effects of individual observations in the estimation of fixed effects for linear mixed models. Those are Cook's distance and COVRATIO. Results of our limited simulation study suggest that the Cook's distance is not good statistical quantity in linear mixed models. Also calibration point for COVRATIO seems to be quite conservative.

Keywords : Cook's distance, COVRATIO, Linear mixed models.

1. 서론

회귀모형이나 고정효과모형에 있어서 여러 가지 진단방법들에 대한 눈부신 연구들은 랜덤효과와 혼합모형에까지 확장되어 진행되어지고 있다. 1980년대 후반부터는 본격적으로 회귀모형에서 혼합모형으로 일반화되어서 진단의 문제가 언급되기 시작했는데, 혼합모형의 진단 중에서 제일 먼저 관심을 가졌던 분야는 분산성분의 점추정 이었으며, 실제로 분산성분추정치의 값이 음수가 되는 경우가 발생하는 중요한 원인이 영향력 관찰점이라고 생각되어지기 때문이었다. Hocking(1983)은 이러한 연구의 선구자로서 분산성분추정치가 음수가 되는 원인을 연구하였으며, 수상한 데이터를 규명하기 위한 진단도구들을 제안하였다. 그 후 Hocking과 Bremer(1987)는 균형혼합모형에서 AVE 추정량과 진단방법을 제안하였으나 불균형혼합모형에 있어서 AVE 추정량의 우수성은 입증되지 못했으며, 셀이 비어있는 경우에도 사용할 수가 없다. Fellner(1986)는 혼합모형에 있어서 분산성분을 추정하는데 이상점의 영향력을 제한최우추정량(REMLE)을 이용하여 구하였다. 또한 Beckman, Nachtsheim과 Cook(1987)은 혼합모형에 사용하는 일반적인 가정에 대한 타당성을 평가하는 우도원리에 입각한 새로운 방법을 제안하였다.

1990년대에는 보다 주목할 만한 논문들이 발표되었는데, Christensen, Pearson과 Johnson(1992)

1) This work was supported by Dankook University Research Fund in 2003

2) Professor, Division of Information and Computer Science, Dankook University, San 8, HanNam-Dong, YongSan-Gu, Seoul, Korea
E-mail : jtlee@dankook.ac.kr

은 REMLE를 이용한 혼합모형의 진단 문제를 다루었다. 그들은 특정 관측치를 제거한 후 나머지 관측치를 이용하여 영향력을 측정하는 소거법을 위한 편리한 계산방법 및 고정효과와 분산성분에 대한 진단도구들을 제안하였다. 그리고 Hurtado(1993)는 최우추정량(MLE)을 이용한 혼합모형의 진단에 사용되는 통계량을 제안하였는데, 분산성분의 비율이 기지인 경우에는 제안된 통계량을 이용하여 정확하게 한 개의 관측치가 영향력 관찰점인가를 알 수 있으며, 미지인 경우에는 몇 가지 가정 아래에서 근사적으로 적용 가능하다.

본 연구에서는 가장 보편화된 회귀모형에서의 영향력 측도들이 혼합모형에서도 적용 가능한지를 구체적으로 모의실험을 통하여 살펴보고자 한다. 고려된 진단통계량은 회귀모형에서 기준치(calibration point)가 분명하게 제안되어 있는 쿡의 통계량(Cook, 1977)과 Belsley, Kuh와 Welsch(1980)가 제안한 COVRATIO를 사용하였다.

본 논문의 구성은 2절에서는 선형혼합모형과 고려된 영향력 측도에 대하여 간략하게 서술하고, 3절에서는 본 연구에서 사용된 분산성분추정량의 종류, 모의실험에 사용된 모형의 종류 및 조건, 모의실험의 과정과 추정된 쿡 통계량과 COVRATIO에 대한 결과를 다양한 영향력 관찰점의 발생 위치 및 개수별로 각각 나누어서 서술한다. 끝으로 4절에서는 본 연구의 결론이 주어진다.

2. 선형혼합모형과 영향력 측도

이 절에서는 선형혼합모형과 본 연구에서 사용한 영향력 측도를 간략하게 설명한다.

2.1 선형혼합모형

선형혼합모형은 다음과 같이 기술된다.

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e} \quad (2.1)$$

여기서 \mathbf{y} 는 알려진 $n \times 1$ 관측치 벡터, \mathbf{X} 는 고정효과에 관련된 $n \times p$ 계획행렬, \mathbf{Z} 는 랜덤효과에 관련된 $n \times q$ 계획행렬, $\boldsymbol{\beta}$ 는 고정효과로 불리는 $p \times 1$ 열벡터, \mathbf{u} 는 랜덤효과로 불리는 $q \times 1$ 열벡터, 그리고 \mathbf{e} 는 $n \times 1$ 오차벡터이다. 또한 \mathbf{u} 에 포함된 랜덤효과들에 대응되는 분산성분들을 표시하기 위하여 \mathbf{u} 를 c 개의 부분벡터 $\mathbf{u}' = (\mathbf{u}_1' | \mathbf{u}_2' | \cdots | \mathbf{u}_c')$ 과 같이 분할하고 \mathbf{u}_i' 에 대응되는 계획행렬을 Z_i 로 두면 식(2.1)은 다음과 같은 식(2.2)로 표현할 수 있다.

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + Z_1\mathbf{u}_1 + Z_2\mathbf{u}_2 + \cdots + Z_c\mathbf{u}_c + \mathbf{e} \quad (2.2)$$

식(2.2)에서 Z_i 의 차수는 $n \times m_i$ 행렬, \mathbf{u}_i 는 $m_i \times 1$ 열벡터이며, q 와 m_i 의 관계식은 $q = \sum_{i=1}^c m_i$ 와 같다. 아울러 다음과 같은 사실이 성립한다고 가정한다.

$$\mathbf{u}_i \sim N(\mathbf{0}, \sigma_i^2 I_{m_i}), \quad \mathbf{e} \sim N(\mathbf{0}, \sigma^2 I_n), \quad \text{Cov}(\mathbf{u}_i, \mathbf{e}') = 0, \quad \text{Cov}(\mathbf{u}_i, \mathbf{u}_j) = 0, \quad \forall i \neq j. \quad (2.3)$$

여기서 I_{m_i} 와 I_n 은 각각 차수가 $m_i \times m_i$, $n \times n$ 인 항등행렬이다. 또한 γ_i 를 $i = 1, 2, \dots, c$ 에 대하여 $\gamma_i = \sigma_i^2 / \sigma^2$ 으로 정의하면 관측치 벡터 \mathbf{y} 의 분산행렬 V 는 $V = \sigma^2 H$, $H = I_n + \sum_{i=1}^c \gamma_i Z_i Z_i'$ 으로 나타낼 수 있다. 따라서 식(2.2)는 다음과 같은 선형모형으로 기술된다.

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}, \quad \mathbf{e} \sim N(\mathbf{0}, \sigma^2 H) \quad (2.4)$$

한편 H 는 정칙행렬이며 대칭이므로 $H = TT'$ 의 관계를 만족하는 정칙행렬 T 가 존재하며, 역행렬 T^{-1} 을 식(2.4)의 좌변과 우변에 각각 곱하면, $e^* = T^{-1}e$ 의 분산행렬이 $\sigma^2 I_n$ 가 되어 β 를 추정하는데 BLUE인 최소제곱추정량을 사용할 수 있다. 그러므로 변환된 모형에 대한 여러 가지 회귀진단도구들을 사용하면, 혼합모형 y 에 대한 진단도구들을 구할 수 있다.

2.2 영향력 측도

관측치의 영향력을 판단하기 위해서는 영향력 측도가 필요하다. 지금까지 제안된 회귀모형에 대한 영향력 측도는 매우 많은 종류가 있으나 크게 나누어 소거법, 무한소교란법, 국소 영향력, 대치법과 같은 네 가지로 분류된다. 이 중에서 소거법이 가장 많이 사용되며 많은 측도들이 이 방법에 의해 유도되었다. 4가지 방법에 대한 상세한 설명과 관계식은 강근석과 김충락(1999) 등에 나와 있다. 본 연구에서는 소거법을 이용한 측도 중에서 쿡 통계량과 COVRATIO 두 가지만을 고려하였다. 왜냐하면 두 가지 통계량 이외에도 앤드류-프레기본 통계량, DFBETAS, DFFITS와 같은 회귀진단통계량이 많이 있으나 기준치가 명확하게 제안되어 있지 않은 이유로 본 연구에서는 제외하였다.

(1) 쿡 통계량

소거법에 의한 측도 중에서는 쿡 통계량(Cook's distance)이 가장 많이 사용되고 있으며, 회귀모형인 경우에는 SAS, SPSS, MINITAB 등의 여러 통계패키지에서 간단한 명령어에 의해 쉽게 구해진다. 혼합모형에 있어서 i 번째 고정효과의 영향력을 평가하기 위하여 고정효과에 대한 쿡 통계량은

$$D_\beta = (\mathbf{b} - \mathbf{b}_{[i]})'(\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})(\mathbf{b} - \mathbf{b}_{[i]})/p \quad (2.5)$$

로 정의되며, 이 경우 \mathbf{b} 는 β 의 일반화최소제곱추정량 $\mathbf{b} = (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{y}$ 이며, $\mathbf{b}_{[i]}$ 는 i 번째 관측치를 제외한 $n-1$ 개의 관측치로 구한 β 의 일반화최소제곱추정량이다. 분산행렬 \mathbf{V} 가 알려져 있는 경우에는 D_β 를 자유도가 p 인 χ^2 -분포를 이용하여 영향력 판찰점 여부를 판단할 수 있으며, 추정된 경우에는 근사적으로 사용할 수 있다(Christensen et al, 1992).

(2) COVRATIO

i 번째 관측치의 $\text{Cov}(\mathbf{b})$ 의 추정치에 대한 영향력을 고려하기 위해 $\text{Cov}(\mathbf{b}) = \sigma^2 (\mathbf{X}' \mathbf{H}^{-1} \mathbf{X})^{-1}$ 와 $\text{Cov}(\mathbf{b}_{[i]}) = \sigma^2_{[i]} (\mathbf{X}_{[i]}' \mathbf{H}^{-1} \mathbf{X}_{[i]})^{-1}$ 를 비교하는 것이 바람직하다. 여기서 $\mathbf{X}_{[i]}$ 는 i 번째 관측치를 제거시킨 $(n-1) \times p$ 행렬이다. 그런데 이들은 모두 행렬이므로 스칼라로 바꾸어 주는 작업이 필요하며, Belsley, Kuh와 Welsch(1980)는 회귀모형에 대한 행렬식의 비인 COVRATIO를 제안하였는데, 혼합모형에 대해서는 σ^2 의 추정치로 s^2 , i 번째 관측치를 제외하고 구한 $\sigma^2_{[i]}$ 의 추정치로 $s^2_{[i]}$ 을 사용하여 다음과 같이 정의할 수 있다.

$$\text{COVRATIO}_i = \frac{\det[s^2_{[i]} (\mathbf{X}_{[i]}' \mathbf{H}^{-1} \mathbf{X}_{[i]})^{-1}]}{\det[s^2 (\mathbf{X}' \mathbf{H}^{-1} \mathbf{X})^{-1}]} \quad (2.6)$$

식(2.6)에서 $\det[\cdot]$ 는 행렬식을 의미하며, 회귀모형인 경우에 COVRATIO_i 의 기준치로 Belsley, Kuh와 Welsch(1980)는 $|\text{COVRATIO}_i - 1| \geq 3p/n$ 을 제시한 바 있다.

3. 모의실험

이 절에서는 쿠 통계량과 COVRATIO_i 의 기준치들이 분산성분이 미지인 선형혼합모형에도 적용될 수 있는지를 모의실험을 통하여 알아본다.

3.1 사용된 분산성분추정량

분산성분들이 미지인 경우에는 데이터로부터 먼저 분산성분을 추정하고, 영향력 측도에 이를 대입하여 구하여만 한다. 혼합모형의 분산성분을 추정하기 위한 방법들은 첫째로 불편추정량의 원리에 입각한 추정량, 둘째는 우도원리에 입각한 추정량, 셋째는 베이지안 방법에 입각한 추정량들로 크게 분류된다. 이 중에서 본 연구에서는 대부분의 통계패키지가 제공하는 가장 보편화된 네 가지 추정량인 분산분석추정량(ANOVA 추정량), 최소분산이차불편추정량(MIVQUE), 최우추정량(MLE)과 제한최우추정량(REMLE)을 사용하였다. 여러 가지 분산성분추정량에 관한 보다 상세한 내용은 Rao와 Kleffe(1988) 또는 Searle(1987)의 책에서 찾아볼 수 있다.

3.2 모의실험계획

모의실험에서는 가장 간단한 혼합모형인 일원변량모형을 선택하였다. 그리고 고려된 모의실험계획은 <표 1>과 같다. <표 1>에서 실험계획패턴의 팔호표기는 (n_1, n_2, \dots, n_p) 를 의미한다. 금내상 관계수 $\rho = \sigma_a^2 / (\sigma_a^2 + \sigma_e^2)$ 는 0.1부터 0.9까지 0.1 간격으로 선택되어졌으며 <표 1>에서 설명된 여러 가지 경우에 대하여 각각 100,000개의 자료가 생성되어졌다.

<표 1> 모의실험에 사용된 모형의 종류와 조건

모형	고려된 조건	사용된 조건의 값
$y_{ij} = \mu + a_i + e_{ij}$, $i = 1, 2, \dots, p; j = 1, 2, \dots, n_i$.	총 관측치의 개수	$n = 15, 30$
	p 의 값	$p = 3, 6$
	확률분포	$a_i \sim N(0, \sigma_a^2), e_i \sim N(0, \sigma_e^2)$
	사용된 모수의 값	$\mu = 0; \sigma^2 = 1$
	금내상 관계수	$\rho = 0.1, 0.2, \dots, 0.9$
	독립성 가정	a_i, e_{ij} : 서로 독립
	실험계획패턴	$P1 = (5, 5, 5), P2 = (2, 5, 8), P3 = (5, 5, 5, 5, 5), P4 = (2, 2, 5, 5, 8, 8)$

3.3 모의실험의 과정

일반적으로 영향력 관찰점이 이상점이 아닐 수도 있고, 이상점이 영향력 관찰점이 아닐 수도 있지만 실질적인 대부분의 경우에는 이상점이 영향력 관찰점이 되기 때문에 고정효과에 대한 영향력 관찰점의 후보로는 생성된 첫 번째 자료의 값에 5와 10을 각각 더한 값으로 정의를 하였다. 실

영향력 관찰점이 2개 존재하는 경우인 C11인 경우는 같은 셀에 영향력 관찰점이 2개 존재하는 C07보다는 양호하나 여전히 영향력 관찰점을 잘 지적하지 못한다.

2. 불균형디자인 P2인 경우, 영향력 관찰점의 개수가 1개인 경우에 C02와 C05를 C01과 비교하면 불균형정도보다는 영향력 관찰점이 존재하는 셀의 위치에 더욱 민감함을 살펴볼 수 있다. 쿠의 통계량은 셀의 도수가 적은 곳에 영향력 관찰점이 위치하는 경우가 영향력 관찰점을 좀 더 잘 지적하나 COVRATIO는 그 반대이다. 영향력 관찰점이 2개인 C08과 C12를 비교하면 C12의 COVRATIO만이 영향력 관찰점을 지적하고 나머지 경우는 지적을 못한다.
3. 균형디자인 P3인 경우, C03과 C01, C09와 C07, C13과 C11를 비교하면 셀의 개수가 늘어나면 쿠의 통계량은 영향력 관찰점을 지적할 가능성은 더 떨어지나 COVRATIO는 그 반대이다. 또한 불균형디자인 P4인 경우도 셀의 개수가 늘어나면 C04과 C06을 C02와 C05와 각각 비교하여 살펴볼 때 균형디자인과 마찬가지 결론이 성립한다.
4. 영향적 관찰점을 조사하는 경우에 분산성분추정량의 선택은 ANOVA추정량, MIVQUE와 REML은 비슷한 효율성을 가지나 MLE는 세 가지 추정량보다 쿠의 통계량 및 COVRATIO에 대하여 영향력관찰점을 잘 지적하기 때문에 네 가지 추정량 중에서 영향력 측도에 사용될 가장 바람직한 추정량이라고 판단되어진다.

10을 더하여 만든 영향력 관찰점의 결과인 <표5>를 <표3>과 <표4>와 비교하면 쿠 통계량은 상대적으로 모든 실험계획패턴에서 영향력 관찰점을 잘 지적한다. 따라서 쿠 통계량은 영향력 관찰점의 크기에 영향을 많이 받는다고 할 수 있다. 반면에 COVRATIO는 빈도수가 적은 셀에 영향력 관찰점이 존재하는 C02와 C04, 또 같은 셀에 영향력 관찰점이 2개가 존재하는 C07에서 C10에는 영향력 관찰점의 크기의 변화가 큰 영향을 미치지 않음을 확인할 수 있었다.

4. 결론

본 연구에서는 영향력 측도로 회귀모형에서 기준치가 명확하게 제안되어 있는 쿠 통계량과 COVRATIO를 사용하여 선형혼합모형에도 적용가능한지를 모의실험을 통하여 살펴보았다. 제한된 모의실험의 결과, 두 통계량의 값은 영향력 관찰점의 개수 및 발생 셀의 위치에 매우 큰 영향을 받는 것으로 판명되었으며, 회귀모형의 기준치를 사용한 COVRATIO는 전반적으로 너무 보수적으로 판정하며 또한 영향력 관찰점이 많고, 같은 셀에 많이 발생하는 경우에는 적당하지 않은 것으로 간주되어진다. 그리고 χ^2 -분포를 이용한 쿠 통계량의 기준치는 혼합모형에서 영향력 관찰점이 매우 특이한 극단값인 경우를 제외한 대부분의 경우에 적당하지 않은 것으로 판명되었다. 또한 본 연구는 선형혼합모형의 고정효과에 대한 영향력 관찰점의 지적에 대한 논의였지만 Hurtado (1993)의 연구에 의하면 고정효과에 대한 영향력 관찰점은 랜덤효과의 영향력 관찰점이며, 오차항의 분산성분추정에 같이 영향을 미친다고 알려져 있으며, 따라서 랜덤효과와 분산성분에의 추정에도 비슷한 결론에 도달될 것으로 간주되어진다.

본 연구에서는 논의되지 않았지만, 혼합모형에 있어서 두 통계량의 타당한 기준치를 정하는 문제는 향후 논의되어야 할 중요한 주제가 될 것이다. 또한 본 연구에서 사용한 랜덤효과의 공분산 구조가 특정구조가 아닌 일반적인 구조를 가질 때 쿠 통계량과 COVRATIO에 대한 보다 폭넓은

접근도 반드시 도전하여야 할 중요과제라고 생각되어진다.

참 고 문 헌

- [1] 강근석, 김충락 (1999). 「회귀분석」, 교우사.
- [2] Beckman, R. J., Nachtsheim, C. J., and Cook, R. D. (1987). Diagnostics for Mixed-Model Analysis of Variance, *Technometrics*, 29, 413-426.
- [3] Belsley, D. A., Kuh, E., and Welsch, R. E. (1980). *Regression Diagnostics*, New York: John Wiley.
- [4] Christensen, R., Pearson, L. M. and Johnson, W. (1992). Case-Deletion Diagnostics for Mixed Models, *Technometrics*, 34, 38-44.
- [5] Cook, R. D. (1977). Detection of Influential Observations in Linear Regression, *Technometric*, 19, 15-18.
- [6] Cook, R. D. and Weisberg, S. (1982). *Residuals and Influence in Regression*, New York: Chapman and Hall.
- [7] Fellner, W. H. (1986). Robust Estimation of Variance Components, *Technometrics*, 28, 51-60.
- [8] Hocking, R. R. (1983). A Diagnostic Tool for Mixed Models with Application to Negative Estimates of Variance Components, *SUGI 8*, 711-716.
- [9] Hocking, R. R. and Bremer, R. H. (1987). Estimation of Variance Components in Mixed Factorial Models Including Model-Based Diagnostics, *SUGI 12*, 1162-1167.
- [10] Hurtado, G. (1993). *Detection of Influential Observations in Linear Mixed Models*, Ph.D. dissertation, North Carolina State University.
- [11] Rao, C. R. and Kleffe, J. (1988). *Estimation of Variance Components and Applications*, Amsterdam: North-Holland.
- [12] Searle, S. R. (1987). *Linear Models for Unbalanced Data*, John Wiley & Sons, New York.

[2002년 12월 접수, 2003년 3월 채택]