

Comparison study of modeling covariance matrix for multivariate longitudinal data

Na Young Kwak^a · Keunbaik Lee^{a,1}

^aDepartment of Statistics, Sungkyunkwan University

(Received March 5, 2020; Revised April 2, 2020; Accepted April 9, 2020)

Abstract

Repeated outcomes from the same subjects are referred to as longitudinal data. Analysis of the data requires different methods unlike cross-sectional data analysis. It is important to model the covariance matrix because the correlation between the repeated outcomes must be considered when estimating the effects of covariates on the mean response. However, the modeling of the covariance matrix is tricky because there are many parameters to be estimated, and the estimated covariance matrix should be positive definite. In this paper, we consider analysis of multivariate longitudinal data via two modeling methodologies for the covariance matrix for multivariate longitudinal data. Both methods describe serial correlations of multivariate longitudinal outcomes using a modified Cholesky decomposition. However, the two methods consider different decompositions to explain the correlation between simultaneous responses. The first method uses enhanced linear covariance models so that the covariance matrix satisfies a positive definiteness condition; in addition, and principal component analysis and maximization-minimization algorithm (MM algorithm) were used to estimate model parameters. The second method considers variance-correlation decomposition and hypersphere decomposition to model covariance matrix. Simulations are used to compare the performance of the two methodologies.

Keywords: correlation matrix, hypersphere decomposition, modified Cholesky decomposition, positive definite, variance-correlation decomposition

1. 서론

한 개체를 여러 번 반복 측정된 자료를 경시적 자료(longitudinal data)라 한다. 같은 개체에서 측정값이 여러 번 측정되므로 경시적 자료에서는 모든 관측치가 서로 독립이라는 통계적 가정이 성립하지 않는다. 이러한 통계적 가정이 성립하지 않은 채로 공변량의 효과를 추정한다면 그 결과는 편향이 일어날 수 있다. 따라서 경시적 자료를 분석하기 위해서는 흔히 사용되는 횡단자료(cross-sectional data) 분석과는 달리, 관측치 간의 상관성을 설명하는 공분산행렬의 모형화에 초점을 맞추어야 한다. 이미 일변량 경시적 자료의 공분산행렬(covariance matrix) 모형화를 제안하는 방법들은 활발히 연구되었다. 그

This project was supported by Basic Science Research Program through the National Research Foundation of Korea (KRF) funded by the Ministry of Education, Science and Technology (NRF-2019R1F1A1058553).

¹Corresponding author: Department of Statistics, Sungkyunkwan University, 25-2, Sungkyunkwan-ro, Jongno-gu, Seoul 03063, Korea. E-mail: keunbaik@skku.edu

리나 다변량 경시적 자료(multivariate longitudinal data)에 대한 공분산행렬 모형화의 경우, 추정해야 할 모수가 많고, 추정된 공분산행렬이 양정치성(positive-definiteness)을 만족해야 하므로 일반변인 경우에 비해 많은 연구가 이루어지지 않았다. 따라서 본 연구에서는 다변량 경시적 자료에 대한 공분산행렬 모형화에 초점을 맞추도록 한다.

다변량 경시적 자료분석에서 공분산행렬의 모형화에서 다음의 3가지 상관관계를 고려해야 한다 (Kim과 Zimmerman, 2012). 1) 같은 속성의 응답변수들이 반복측정될 때에 발생하는 상관관계; 2) 서로 다른 속성들 간의 시간에 따라 발생하는 상관관계; 3) 같은 시간에서 변수들 간의 상관관계. 여기서 1)과 2)의 상관관계는 시간에 의한 상관관계(serial correlations)이며, 3)의 상관관계는 시간과는 무관한 상관관계이다. 이러한 3가지의 상관관계는 공분산행렬에서 설명되며, 이 행렬을 모형화 하기 위한 방법들이 몇몇 논문에서 제시되었다. 이 논문들은 모두 Pourahmadi (1999)의 수정된 콜레스키 분해(modified Cholesky decomposition; MCD)를 이용하여 모형화 하였다. Pourahmadi (1999)는 다변량 경시적 자료의 공분산행렬을 모형화하기 위해 MCD를 이용하였고, 이 때 MCD는 공분산행렬을 직접 분해하는 것이 아니라, 공분산행렬의 역행렬을 분해할 때에 적용하였다. 즉, 공분산행렬의 역행렬을 일반화 자기회귀모수(generalized autoregressive parameteres; GARPs)와 혁신분산(innovation variances; IVs)으로 분해함으로써, 제약은 없으면서도 통계적으로 의미 있는 모수로 분해하였고, 이렇게 만들어진 공분산행렬은 양정치성을 만족하게 된다. 또한, GARPs와 IVs을 선형/로그선형모형 식으로 표현하여 추정함으로써 모수의 수를 줄일 수 있다.

다변량 경시적 자료분석에서 MCD를 이용한 방법으로 Xu와 Mackenzie (2012)와 Kim과 Zimmerman (2012)이 있다. 이 두 논문 모두 공분산행렬을 일반화 자기회귀행렬(generalized autoregressive matrices; GARMs)과 혁신공분산행렬(innovation covariance matrices; ICMs)로 분해하였다. 공분산행렬을 위의 두 행렬로 분해함으로써 공분산행렬을 간단한 모형으로 나타낼 수 있으며 통계적인 해석을 하기도 유용하다. GARMs는 블록 하삼각행렬(block lower triangular matrix) 형태이며, 각 요소들은 다변량 경시적 자료에서 발생하는 시간에 의한 상관관계인 1)과 2)의 상관관계를 설명할 수 있다. ICMs는 블록 대각행렬(block diagonal matrix)으로 나타낼 수 있고, 이 행렬은 다변량 경시적 자료에서 발생하는 마지막 상관관계인 3)의 상관관계를 설명할 수 있다. 이러한 두 행렬로 공분산행렬이 분해되므로 이 두 행렬은 유일하게 존재하며, IVMs가 양정치성을 만족하면 자연스럽게 공분산행렬도 양정치성을 만족한다 (Kim과 Zimmerman, 2012). 따라서 IVMs가 양정치성을 만족하도록 모형화 하는 것이 중요하다. 이 ICMs를 모형화 하는 방법에 따라서 여러 가지 방법들이 제안되었다 (Xu와 Mackenzie, 2012; Kim과 Zimmerman, 2012; Kohli 등, 2016; Lee 등, 2019).

Xu와 Mackenzie (2012)는 혁신공분산행렬을 스펙트럴 분해(spectral decomposition)를 이용하여 고유벡터(eigenvector)로 구성된 직교행렬과 고유값(eigenvalues)을 가지는 대각행렬로 분해한다. 그리고 ICMs의 로그변환은 결국 고유값의 로그변환으로 연결되면 이는 고유값의 양수를 보장받을 수 있다. 그리고 이를 통하여 IVMs의 양정치성을 만족시키는 모형이다. GARMs의 요소와 IVMs에서 만들어진 고유치는 선형회귀 및 로그 선형회귀 모형으로 모형화 할 수 있다. Xu와 Mackenzie (2012)의 공분산행렬 모형화 방법은 자료가 단조적인 중도탈락(monotone dropout)에 의해 불균형 형태일 때에도 개체별로 공분산행렬을 모형화 하는 것이 가능하다는 장점이 있다. 그러나 각 시점의 값 간의 상관관계를 설명하기 위한 공분산 행렬인 IVMs의 모수들을 추정하기 어렵고 해석하기 어렵다는 단점이 있다. Kim과 Zimmerman (2012)는 Xu와 Mackenzie (2012)와 같이 공분산행렬을 MCD를 이용하여 GARMs과 IVMs로 분해하였고, IVMs를 MCD를 다시 이용하여 모형화 하였다. IVMs는 MCD를 이용하여 분해함에 따라 모형화에 대한 제약을 제거하여 간단화 하였으나, 동일 시점에서의 반응값들에 순서를 부여하여 실제 데이터에 적용하기에는 한계가 있다는 단점이 있다. Kohli 등 (2016)은 Kim과 Zimmerman

(2012)과 Xu와 Mackenzie (2012)와 같이 MCD를 이용하고, IVMs를 알려진 양정치행렬과 미지의 양수인 모수의 선형 결합으로 나타냄으로써 IVMs의 양정치성을 만족시켰다. 따라서 알려진 양정치 행렬을 선택하는 것이 핵심이며, 총 세가지 방법으로 양정치행렬을 선택하였다. IVMs를 직접적으로 모형화한 기존 방법들과는 달리, Lee 등 (2019)은 IVMs를 분산-상관 분해(variance-correlation decomposition)를 하여, 혁신표준편차(innovation standard deviations; ISDs)와 상관계수(correlations)로 분해하여 모형화 하였다. 상관계수행렬을 모형화 할 때에는 양정치성과 모든 원소들이 -1 과 1 사이의 값을 가져야하며 대각원소는 반드시 1 이어야 한다는 제약조건이 있다. Lee 등 (2019)은 상관행렬을 모형화하기 위해 초구분해(hypersphere decomposition; HD) 방법을 이용하였다. 분산-상관 분해를 이용하여 IVMs 모형화 했기 때문에 Kohli 등 (2016)과는 달리 해석이 가능하다는 장점이 있다. 또한, Kim과 Zimmerman (2012)과는 달리 같은 시점에서의 반응값에 순서를 부여하지 않고, MCD를 이용하여 공분산 행렬을 모형화 하였다. 우리는 위에서 제시한 다변량 경시적 자료에 대한 공분산행렬의 모형화 방법 중에 Kohli 등 (2016)의 방법과 Lee 등 (2019)의 방법을 중점적으로 살펴보고자 한다. 그 이유는 Kohli 등 (2016) 방법은 IVMs의 모형화에 다양한 방법을 제안하였고, Lee 등 (2019)은 자연스러운 해석이 가능한 모형이라서 선택하였다. 그리고 Lee 등 (2019)의 모형이 Kim과 Zimmerman (2012)의 모형 보다 편향과 효율성에서 더 나은 것으로 확인되었기 때문이다 (Lee 등, 2019). 이 논문에서 우리는 이 방법들을 모의실험을 통해 여러 경우에 대해 비교해 보고자 한다.

논문의 구성은 다음과 같다. 2장에서 MCD를 자세하게 설명한 후, MCD를 기반으로 한 Kohli 등 (2016)의 방법과 Lee 등 (2019)의 방법을 제시한다. 3장에서는 이 두 가지 방법을 이용하여 모의실험을 진행하고 그 결과를 논의한다. 4장에서는 3장에서의 결과를 바탕으로 결론을 내린다.

2. 공분산 행렬의 모형화

2.1. 단변량 경시적 자료에서의 MCD를 이용한 모형화

Pourahmadi (1999)는 공분산행렬을 모형화 하기 위해서 다음의 모형을 제안하였다. $Y_i = (Y_{i1}, \dots, Y_{in})^T$ 는 i 번째 대상자의 응답변수 벡터이다 ($i = 1, \dots, N$). 여기서 Y_{it} 는 i 번째 대상자의 t 번째 응답변수이다 ($t = 1, \dots, n$). \hat{Y}_{it} 를 $\hat{Y}_{it} = \mu_{it} + \sum_{j=1}^{t-1} \phi_{it,j}(Y_{ij} - \mu_{ij})$ 라고 정의한다. 그리고 계산의 편의를 위해 $\mu_{ij} = E(Y_{ij}) = 0$ 으로 가정한다. 따라서 $\hat{Y}_{it} = \sum_{j=1}^{t-1} \phi_{it,j}Y_{ij}$ 이고, 이것은 마치 회귀식에서의 추정식과 유사하다. 이 때의 예측오차는 ε_{it} 이며 $\sigma_{it}^2 = \text{var}(\varepsilon_{it})$ 이다. 즉,

$$\varepsilon_{it} = Y_{it} - \hat{Y}_{it} = Y_{it} - \sum_{j=1}^{t-1} \phi_{it,j}Y_{ij}, \quad t = 1, \dots, n \quad (2.1)$$

이며, 연속된 예측오차끼리는 서로 독립이다. 대각원소가 1 이고 하삼각원소가 $\hat{Y}_{it} = \sum_{j=1}^{t-1} \phi_{it,j}Y_{ij}$ 의 회귀계수에 대한 음수 값인 행렬을 T_i , 대각원소가 $\sigma_{it}^2 = \text{var}(Y_{it} - \hat{Y}_{it})$ 인 행렬을 D_i 라고 정의한다. 즉, $\varepsilon_i = T_i Y_i$ 으로 표현할 수 있고, 여기에 분산을 취하면 다음과 같다.

$$\text{cov}(\varepsilon_i) = T_i \text{cov}(Y_i) T_i^T = T_i \Sigma T_i^T = D_i \quad (2.2)$$

라는 식이 도출되는 것을 확인할 수 있다. 식 (2.2)에서 행렬 T_i 의 요소를 일반화 자기회귀모수, 행렬 D_i 를 혁신분산이라 지칭한다. T_i 가 하삼각 행렬이고 D_i 가 대각행렬이므로 Σ 는 대칭행렬이 된다. T_i 와 D_i 는 통계적으로 해석이 용이하며, Σ 를 추정하기 위해서는 $(1/2)n(n-1)$ 개의 $\phi_{it,j}$ 와 $\log(\sigma_{it}^2)$ 만 추정하면 된다. Σ_i 를 추정하기 위해 T_i 와 D_i 의 원소만 추정하면 되며, 모수의 수를 줄이기 위한 선형 및 로

그 선형모형이 제안되었다. T_i 의 원소인 $\phi_{it,j}$ 와 D_i 는 주대각 요소인 σ_{it} 의 모형화는 다음과 같다.

$$\phi_{it,j} = z_{it,j}^T \gamma, \quad (2.3)$$

$$\log(\sigma_{it}^2) = h_{it}^T \lambda. \quad (2.4)$$

이때 $z_{it,j}$ 와 h_{it} 는 개체-특정적 공변량이며, γ 와 λ 는 알려지지 않은 모수이다.

2.2. 다변량 경시적 자료에서의 MCD를 이용한 모형화

Pourahmadi (1999)가 제안한 MCD는 단변량 경시적 자료에 한정된 방법이기 때문에 다변량 경시적 자료에서 MCD를 사용하기 위해서는 수정된 블록 콜레스키 분해(modified Cholesky block decomposition; MCB) 방법을 사용해야 한다. 따라서 MCD에서는 하나의 값이었던 T_i 와 D_i 의 원소들이 MCB에서는 행렬로 표현되어야 한다.

2.2.1. 수정된 블록 콜레스키 분해 K 개의 특성이 있는 다변량 경시적 자료에서 Y_{itk} 는 i 번째 대상자의 t 번째 반복수의 k 번째 속성을 나타낸다($i = 1, \dots, N; t = 1, \dots, n; k = 1, \dots, K$). 이 경우 다음과 같은 선형모형으로 표현할 수 있다.

$$\begin{aligned} Y_{i1k} &= x_{i1}^T \beta_k + \epsilon_{i1k}, \\ Y_{itk} &= x_{it}^T \beta_k + \sum_{j=1}^{t-1} \sum_{g=1}^K \phi_{itk,kg} (Y_{ijg} - x_{ij}^T \beta_g) + \epsilon_{itk}, \end{aligned} \quad (2.5)$$

여기서 x_{it} 는 $p \times 1$ 공변량 벡터, β_k 는 미지의 회귀계수 벡터이다. 식 (2.5)를 행렬 형태로 나타내면

$$T_i(Y_i - X_i \beta) = \epsilon_i \quad (2.6)$$

이다. 여기서

$$T_i = \begin{pmatrix} I & 0 & \cdots & 0 \\ -\Phi_{i21} & I & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ -\Phi_{in1} & -\Phi_{in2} & \cdots & I \end{pmatrix}, \quad Y_i = \begin{pmatrix} Y_{i1} \\ Y_{i2} \\ \vdots \\ Y_{in} \end{pmatrix}, \quad X_i = \begin{pmatrix} X_{i1} \\ X_{i2} \\ \vdots \\ X_{in} \end{pmatrix}, \quad \epsilon_i = \begin{pmatrix} \epsilon_{i1} \\ \epsilon_{i2} \\ \vdots \\ \epsilon_{in} \end{pmatrix},$$

여기서 Φ_{itj} 를 $K \times K$ 의 행렬로 그 요소는 $\phi_{itj,lm}$ ($l, l = 1, \dots, K$)이며, $Y_{it} = (Y_{it1}, \dots, Y_{itK})^T$, $X_{it} = I_K \otimes x_{it}^T$ 는 $K \times p$ 행렬이며, 그리고 $\epsilon_{it} = (\epsilon_{it1}, \dots, \epsilon_{itn})^T$ 이다. \otimes 는 크로네크곱을 의미한다. 그리고 Φ_{itj} 를 자기회귀모수행렬(generalized autoregressive parameter matrix; GARPM)이라고 지칭한다.

단변량 경시적 자료에서의 MCD와 마찬가지로 예측오차벡터인 ϵ_i 는 다음을 만족한다고 가정한다.

$$\epsilon_i = \begin{pmatrix} \epsilon_{i1}^T \\ \vdots \\ \epsilon_{in}^T \end{pmatrix}^T \sim N(0, D_i),$$

여기서 $D_i = \text{diag}\{D_{i1}, \dots, D_{in}\}$ 를 혁신분산행렬(innovation covariance matrix; ICM)이라 지칭하고, $D_{it} = \text{var}(\epsilon_{it})$ 이다.

MCD와 마찬가지로 MCB에서의 T_i 도 회귀계수 행렬과 같이 제약을 받지 않고, 서론에서 제시한 시

간에 의한 응답변수들의 상관관계(상관관계 1)과 2))를 설명이 가능하다. GARP의 원소인 $\phi_{itj,lm}$ 은 시간이나 개체특징적 공변량 벡터인 w_{itj} 을 이용하여 추정할 수 있다.

$$\phi_{itj,lm} = w_{itj}^T \alpha_{lm}, \quad (2.7)$$

이때 α_{lm} 은 미지의 모수로 이루어져 있는 $a \times 1$ 벡터이며, w_{itj} 은 개체특징적 공변량 벡터이다. 그러나 D_i 는 K 개의 결과에서의 횡단적인 종속을 의미하며, 양정치성을 만족해야 하는 제약이 있다. D_i 는 통계적으로 해석이 가능한 동시에 최대한 제약받지 않도록 추정하는 것이 중요하다. 이러한 조건들에 부합하는 D_i 를 추정하기 위하여 Kohli 등 (2016)과 Lee 등 (2019)은 각자 다른 방법론을 제시하였다. D_i 의 모형화는 여러 결과값들의 순서를 필요로 하지 않으며, 추정된 Σ 의 양정치성을 만족시킨다.

2.2.2. 향상된 혁신공분산행렬을 이용한 모형화 Kohli 등 (2016)은 D_i 를 향상된 선형 공분산 모형 (enhanced linear covariance models; LCM)을 이용하여 추정하는 것이다. 이를 식으로 표현하면 다음과 같다.

$$D_{it}^{(q)} = \tau_{t1} M_{t1} + \dots + \tau_{tq} M_{tq}, \quad (2.8)$$

이때의 M_{tl} ($l = 1, \dots, q$)는 알려진 양정치 행렬이며 τ_{tl} 은 미지의 음수가 아닌 모수를 의미한다. 또한 τ_{tl} 은 $D_{it}^{(q)}$ 의 변동을 설명하는 모수로 해석되기도 한다. 이러한 제약은 추정된 Σ 의 양정치성을 보장한다. 따라서 LCM 모형을 사용하는 데 있어 가장 중요한 사안은 적절한 M_{it} 를 선택하는 것이다. Kohli 등 (2016)은 적절한 양정치행렬 M_{it} 를 선택하기 위해 세 가지 방법을 제안하였다. 첫 번째 방법은 모수적인 방법으로써, 알려진 공분산 기반 행렬 묶음에서 M_{it} 를 선택하는 것이다. 두 번째 방법은 비모수적인 방법으로써, 주성분분석(principal component analysis; PCA)을 이용하여 혁신행렬의 고유벡터를 이용해 M_{it} 를 생성하는 것이다. 마지막 방법은 단변량 regressogram 개념을 확장하여 M_{it} 의 양정치성을 만족시키는 것이다.

우선 첫 번째 모수적인 방법은 compound symmetry (CS)나 AR(1)과 같은 모수적 상관관계 구조에서 M_{it} 를 선택하는 것이다. Qu 등 (2000)과 Zhou와 Qu (2012)는 CS와 AR(1)모형이 경시적 자료의 상관관계나 의존관계를 잘 설명할 수 있다는 것을 증명하였다. 따라서 M_{it} 를 AR(1), CS, AR-CS 상관관계 모형에서 선택하는 것을 고려한다. 모수적 방법은 M_{it} 를 쉽게 선택할 수 있다는 장점이 있지만, 상관관계에 대한 사전 정보가 없으면 잘못된 추론을 할 가능성이 있다는 단점이 존재한다. 두 번째 방법은 D 의 표본을 스펙트럴 분해하는 비모수적 방법이다. D 의 표본의 작은 고유값에 대응하는 고유벡터를 이용하여 M_{it} 를 추정할 수 있다. 이 방법은 M_{it} 를 잘못 구축할 위험은 피할 수 있으나, 불균형한 데이터의 경우 표본 공분산행렬이 존재하지 않을 수도 있다는 단점이 있다. 세 번째 방법은 그래프를 이용하여 D_i^{-1} 나 $\log(D_i)$ 의 원소들을 선택하는 방법이다.

2.2.3. 초구분해를 이용한 모형화 Lee 등 (2019)은 Kohli 등 (2016)과는 달리 이분산성까지 고려하여 공분산행렬을 모형화 하고자 하였다. 이분산성과 양정치성을 만족시키기 위해 D_{it} 를 직접 모형화 하는 것이 아니라 분산-상관 분해를 이용하여 혁신분산과 상관행렬로 분해하여 모형화 하였다. 이렇게 혁신분산과 상관행렬의 모수를 추정함으로써 모형화 결과가 해석 불가능한 Kohli 등 (2016)의 경우와는 달리, Lee 등 (2019)의 모형화의 결과는 해석이 가능하다는 장점이 있다. 또한, 이 방법론은 반응변수 간의 순서가 없어도 D_{it} 의 모형화가 가능하므로 실제 데이터로의 적용이 더 자연스럽다. 상관행렬을 모형화 하는 방법론으로 Lee 등 (2019)은 초구분해 방법을 제안하였다. 초구분해가 상관행렬의 양정치성을 만족시키므로, 자연스럽게 D_i 의 양정치성까지 만족시킬 수 있다. Lee 등 (2019)의 따름정리에 의해

D_i 가 양정치성을 만족하면 Σ_i 도 자동적으로 양정치성을 만족한다. 따라서 D_i 의 양정치성을 만족시키기 위해 D_{it} 를 다음과 같은 식으로 표현한다.

$$D_{it} = C_{it}R_iC_{it}, \quad (2.9)$$

이때 C_{it} 는 혁신표준편차를 대각원소로 가지는 대각행렬이며, R_i 는 각 시점에서의 반응변수들 간의 상관행렬이다. 따라서 C_{it} 와 R_i 를 각각 모형화하면 D_{it} 를 추정할 수 있다. 우선 C_{it} 의 대각원소인 혁신표준편차를 추정하기 위한 식은 다음과 같다.

$$\log(\sigma_{itk}) = h_{it}^T \lambda_k, \quad (2.10)$$

여기서 λ 는 알려지지 않은 모수이며 h_{it} 는 시간이나 개체 특정적 공변량을 의미한다. $\log(\sigma_{itk})$ 가 혁신표준편차를 로그 변환한 형태이므로 $\log(\sigma_{itk})$ 은 제약이 없다. 따라서 D_{it} 의 양정치성은 R_i 의 양정치성 여부에 영향을 받는다.

상관계수행렬 R_i 을 모형화 할 때 고려해야 할 사항으로 총 세 가지가 있다. 첫 번째 고려사항은 R_i 가 양정치성을 만족해야 하는 것이다. 두 번째는 R_i 의 주대각 원소를 제외한 모든 원소는 -1 과 1 사이의 값을 가져야 하며, 마지막 고려사항은 R_i 의 대각원소는 반드시 1이어야 한다는 것이다. 이 세 가지 고려사항을 모두 만족시키는 R_i 을 모형화하기 위하여 Lee 등 (2019)은 초구 분해 방법을 사용하였다. 초구 분해 방법을 이용하면 R_i 는 다음과 같이 분해된다.

$$R_i = F_i F_i^T, \quad (2.11)$$

여기서

$$F_i = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ f_{i21} & f_{i22} & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ f_{iK1} & f_{iK2} & f_{iK3} & \cdots & f_{iKK} \end{pmatrix}, \quad (2.12)$$

각각의 F_i 의 원소인 f_{ilm} 은 다음과 같은 식으로 구할 수 있다.

$$f_{ilm} = \begin{cases} \cos(\omega_{ilm}), & \text{for } m = 1, l = 2, \dots, K; \\ \cos(\omega_{ilm}) \prod_{r=1}^{m-1} \sin(\omega_{ilr}), & \text{for } 2 \leq m < l \leq K; \\ \prod_{r=1}^{m-1} \sin(\omega_{ilr}), & \text{for } l = m; m = 2, \dots, K; \end{cases}$$

따라서 제약이 없는 $\omega_{ilm} \in (0, \pi)$ 를 추정하기 위해 알려지지 않은 모수인 v 와 공변량 벡터인 g_{ilm} 을 이용하여 다음과 같은 식으로 표현한다.

$$\log\left(\frac{\omega_{ilm}}{\pi - \omega_{ilm}}\right) = g_{ilm}^T v. \quad (2.13)$$

모형 (2.13)과 같은 식으로 ω_{ilm} 를 재모수화함에 따라 ω_{ilm} 는 $(0, \pi)$ 의 범위를 만족한다. ω_{ilm} 를 추정함으로써 R_i 를 추정할 수 있고, 그 결과로 D_{it} 를 추정할 수 있다.

3. 모의실험

이 절에서 모의실험을 통하여 제 2절에서 언급한 두 가지 방법의 성능을 비교하고자 한다. 모의실험마다 500개의 자료집합을 생성하였다. 각각의 자료집합에는 응답변수를 3종류의 속성($K = 3$)으로 하고, 표본의 크기(N)와 반복수(n)는 경우에 따라서 다르게 하였다. 구체적인 수는 각각의 소절에서 제시하였다.

Table 3.1. $N = 100, n = 5$

β	Lee		Kohli(AR)		Kohli(CS)		Kohli(AR-CS)	
	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)
β_{10} (0.3)	0.312 (0.865 _{0.870})	3.962 88.000	0.282 (0.288 _{0.288})	5.756 96.000	0.292 (0.309 _{0.276})	2.396 96.600	0.289 (0.315 _{0.279})	3.630 98.400
β_{11} (-0.1)	-0.110 (0.703 _{0.672})	8.309 85.000	-0.075 (0.337 _{0.311})	24.490 96.800	-0.088 (0.360 _{0.320})	11.740 97.000	-0.082 (0.368 _{0.326})	17.401 97.800
β_{12} (0.2)	0.262 (0.829 _{0.702})	15.309 89.200	0.190 (0.130 _{0.121})	4.885 96.000	0.197 (0.139 _{0.123})	1.167 96.800	0.192 (0.142 _{0.125})	3.861 97.200
β_{13} (0.3)	0.281 (0.730 _{0.698})	8.934 86.800	0.301 (0.014 _{0.013})	0.353 95.200	0.300 (0.015 _{0.014})	0.075 96.600	0.300 (0.016 _{0.014})	0.259 97.400
β_{20} (0.2)	0.223 (0.594 _{0.509})	13.050 89.200	0.219 (0.307 _{0.281})	9.900 96.200	0.204 (0.325 _{0.295})	20.157 97.800	0.202 (0.333 _{0.279})	10.179 98.000
β_{21} (-0.1)	-0.129 (0.532 _{0.483})	30.493 91.000	-0.125 (0.366 _{0.339})	25.745 96.600	-0.109 (0.143 _{0.134})	9.902 96.800	-0.107 (0.147 _{0.129})	7.526 96.400
β_{22} (0.2)	0.217 (0.692 _{0.629})	12.942 82.600	0.210 (0.137 _{0.125})	5.244 96.600	0.204 (0.143 _{0.134})	2.162 96.800	0.199 (0.147 _{0.129})	0.171 96.400
β_{23} (0.3)	0.339 (0.822 _{0.821})	24.583 89.200	0.298 (0.015 _{0.014})	0.412 96.600	0.299 (0.016 _{0.015})	0.169 96.800	0.300 (0.016 _{0.014})	0.083 96.600
β_{30} (0.2)	0.237 (0.998 _{0.984})	15.834 90.000	0.190 (0.315 _{0.300})	4.960 95.400	0.221 (0.342 _{0.300})	10.844 98.000	0.200 (0.348 _{0.280})	0.136 98.800
β_{31} (-0.2)	-0.228 (0.702 _{0.698})	13.200 85.000	-0.184 (0.370 _{0.345})	7.561 95.800	-0.213 (0.404 _{0.353})	6.786 97.600	-0.206 (0.412 _{0.338})	3.359 98.800
β_{32} (0.2)	0.253 (0.824 _{0.794})	24.196 84.200	0.193 (0.133 _{0.125})	3.425 95.400	0.204 (0.145 _{0.127})	2.023 97.800	0.203 (0.148 _{0.121})	1.725 98.400
β_{33} (0.4)	0.434 (0.883 _{0.829})	6.292 89.200	0.400 (0.014 _{0.013})	0.203 95.600	0.399 (0.015 _{0.014})	0.132 97.600	0.399 (0.016 _{0.013})	0.076 98.200
F.norm	4.838		4.922		5.373		5.502	

3.1. Kohli 등 (2016)에서 생성한 자료집합

우선 Kohli 등 (2016)이 제안한 방식을 이용하여 평균과 AR(1)구조의 공분산행렬을 아래에 제시된 방법으로 다변량 정규분포에서 난수를 발생시켰다.

$$Y_{itk} = x_{it}^T \beta_k + \sum_{j=1}^{t-1} \sum_{g=1}^K \phi_{itj,kg} (Y_{ijg} - x_{ij}^T \beta_g) + e_{it}, \tag{3.1}$$

여기서 공변량은 $x_{it}^T = (1, \text{Group}_i, \text{Time}_{it}, \text{Time}_{it} \times \text{Group}_i)$ 이며, $\text{Time}_{it} = (t - 1)/10$ 이며 Group_i 는 0 또는 1의 값을 가진다. 이 경우 회귀계수의 차원은 $p = 4$ 이며 그 값은 아래와 같다.

$$\beta_1 = (0.3, -0.1, 0.2, 0.3)^T, \quad \beta_2 = (0.2, -0.1, 0.2, 0.3)^T, \quad \beta_3 = (0.2, -0.2, 0.2, 0.4)^T.$$

일반화자기회귀모수의 회귀식 (2.7)에 있는 모수인 α_{lm} 과 향상된 혁신공분산행렬 (2.8)의 모수인 τ_{it} 의 값은 아래와 같다.

$$(\alpha_{11}, \alpha_{12}, \alpha_{13}, \alpha_{21}, \alpha_{22}, \alpha_{23}, \alpha_{31}, \alpha_{32}, \alpha_{33}) = (0.3, 0.4, 0.1, 0.1, 0.3, 0.1, 0.1, 0.3, 0.2),$$

$$(\tau_1, \tau_2, \tau_3) = (1.0, 0.5, 0.7).$$

Table 3.2. $N = 100, n = 10$

β	Lee		Kohli(AR)		Kohli(CS)		Kohli(AR-CS)	
	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)
β_{10}	0.318 (0.794 _{0.759})	6.309 88.000	0.291 (0.157 _{0.137})	2.712 97.600	0.292 (1.862 _{0.142})	2.460 98.800	0.298 (0.195 _{0.144})	0.429 99.600
β_{11}	-0.115 (-0.1) (0.643 _{0.534})	14.392 90.200	-0.093 (0.133 _{0.117})	6.112 97.000	-0.095 (0.157 _{0.120})	4.860 98.800	-0.099 (0.164 _{0.124})	0.213 99.200
β_{12}	0.263 (0.2) (0.584 _{0.539})	8.394 85.000	0.198 (0.042 _{0.037})	0.990 97.000	0.198 (0.050 _{0.038})	0.677 98.800	0.199 (0.053 _{0.039})	0.041 99.000
β_{13}	0.285 (0.3) (0.498 _{0.438})	4.742 88.400	0.300 (0.004 _{0.003})	0.063 96.800	0.300 (0.004 _{0.003})	0.004 98.800	0.300 (0.005 _{0.003})	0.013 99.000
β_{20}	0.220 (0.2) (0.434 _{0.423})	6.306 89.600	0.205 (0.164 _{0.150})	2.561 96.400	0.201 (0.200 _{0.153})	0.795 98.400	0.201 (0.211 _{0.160})	0.834 99.200
β_{21}	-0.113 (-0.1) (0.403 _{0.307})	9.329 90.200	-0.102 (0.138 _{0.128})	2.298 97.200	-0.098 (0.168 _{0.130})	1.085 98.400	-0.098 (0.176 _{0.133})	1.940 99.200
β_{22}	0.220 (0.2) (0.594 _{0.528})	10.391 83.600	0.201 (0.048 _{0.045})	0.639 96.000	0.200 (0.057 _{0.045})	0.047 98.400	0.199 (0.061 _{0.046})	0.177 99.200
β_{23}	0.363 (0.3) (0.693 _{0.684})	7.807 88.200	0.299 (0.005 _{0.004})	0.064 95.800	0.299 (0.006 _{0.005})	0.027 96.400	0.299 (0.006 _{0.004})	0.009 99.400
β_{30}	0.284 (0.2) (0.959 _{0.892})	14.236 90.000	0.200 (0.124 _{0.114})	0.135 95.200	0.201 (0.143 _{0.118})	0.652 97.000	0.197 (0.150 _{0.107})	1.295 99.200
β_{31}	-0.218 (-0.2) (0.632 _{0.611})	6.493 89.200	-0.199 (0.108 _{0.100})	0.093 95.600	-0.201 (0.125 _{0.101})	0.552 98.600	-0.200 (0.131 _{0.096})	0.297 99.900
β_{32}	0.264 (0.2) (0.813 _{0.743})	27.394 84.600	0.199 (0.033 _{0.030})	0.402 94.800	0.199 (0.038 _{0.030})	0.151 98.400	0.200 (0.040 _{0.029})	0.010 99.600
β_{33}	0.433 (0.4) (0.843 _{0.892})	4.302 92.000	0.400 (0.003 _{0.002})	0.030 95.600	0.400 (0.003 _{0.002})	0.018 98.600	0.400 (0.003 _{0.002})	0.013 99.600
F.norm	4.493		4.530		5.248		5.395	

표본의 크기는 100, 300, 500으로 하였고, 각각 500개의 자료집합을 만들었다.

이러한 500개의 자료집합을 제 2절에서 제시한 Kohli 등 (2016)의 모형과 Lee 등 (2019)이 제안한 모형을 적합시켰다. Kohli 등 (2016)의 모형 적합 시에는 공분산행렬을 모수적인 방법인 AR(1), CS, 그리고 AR-CS구조로 가정한 모형을 각각 적합시켰다. 그 적합한 결과들은 각 모수의 추정치 평균(mean), 상대편향의 절댓값(absolute relative bias; ARB), 추정량의 표준오차의 평균(standard error; SE), 추정치들의 표준편차(standard deviation; SD), 그리고 포함확률(coverage probability; CP)으로 요약되었다. 또한 Kohli 등 (2016)이 제안한 모형과 Lee 등 (2019)이 제안한 모형 간의 비교를 위해 프로베니우스 노름(Frobenius norm)을 이용하여 추정된 공분산행렬을 비교하고자 한다.

$$\text{F.norm}(\hat{\Sigma}) = \text{tr} \left\{ \left(\hat{\Sigma} \Sigma^{-1} - I \right)^T \left(\hat{\Sigma} \Sigma^{-1} - I \right) \right\},$$

Tables 3.1과 3.2는 표본의 크기가 100이며, 각각 반복수가 5, 10일 때에 공분산행렬을 앞서 제시한 4가지 모형을 가정하고 적합한 결과를 나타낸 것이다. 그리고 Figure 3.1의 (a)는 모든 평균모수들의 추정 정도를 한 번에 파악하기 위해 ARB값의 합인 $\sum_{i=1}^3 \sum_{j=0}^3 (|\hat{\beta}_{ij} - \beta_{ij}| / \beta_{ij}) \times 100$ 를 계산하여 그 결과를 그림으로 나타내었다. 반복수가 증가하면 편향이 줄어드는 것을 알 수 있다. 그리고 Kohli 등 (2016)의 모든 방법들이 Lee 등 (2019)의 방법보다는 편향이 적음을 알 수 있다. 이는 공분산행렬의 구조를 난

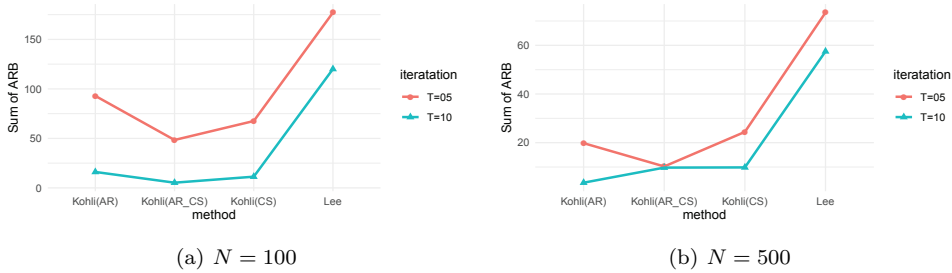


Figure 3.1. Figure 3.1: Sum of ARBs of mean parameters in Lee *et al.* (2020)'s and Kohli *et al.* (2016)'s models using the data generated from Kohli *et al.* (2016)'s model.

Table 3.3. $N = 500, n = 5$

β	Lee		Kohli(AR)		Kohli(CS)		Kohli(AR-CS)	
	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)
β_{10}	0.307 (0.382 _{0.311})	2.394 91.400	0.294 (0.129 _{0.125})	1.813 94.800	0.305 (0.138 _{0.116})	1.891 97.800	0.299 (0.140 _{0.109})	0.090 98.800
β_{11}	-0.114 (-0.1) (0.399 _{0.328})	13.391 89.200	-0.093 (0.151 _{0.146})	6.668 95.200	-0.110 (0.161 _{0.134})	10.319 97.800	-0.102 (0.164 _{0.132})	2.381 98.400
β_{12}	0.223 (0.2) (0.528 _{0.493})	5.730 90.000	0.197 (0.058 _{0.056})	1.098 94.600	0.204 (0.062 _{0.053})	2.198 97.800	0.200 (0.063 _{0.050})	0.403 99.000
β_{13}	0.292 (0.3) (0.211 _{0.209})	3.088 92.600	0.300 (0.006 _{0.006})	0.065 95.200	0.299 (0.007 _{0.006})	0.179 97.200	0.299 (0.007 _{0.005})	0.050 98.800
β_{20}	0.210 (0.2) (0.293 _{0.234})	4.399 92.600	0.195 (0.138 _{0.122})	2.288 96.400	0.201 (0.145 _{0.131})	0.531 97.600	0.202 (0.148 _{0.125})	1.100 98.000
β_{21}	-0.109 (-0.1) (0.278 _{0.300})	7.203 90.000	-0.094 (0.164 _{0.147})	5.521 95.800	-0.098 (0.172 _{0.159})	1.842 97.200	-0.101 (0.176 _{0.151})	1.389 97.800
β_{22}	0.212 (0.2) (0.494 _{0.440})	7.392 89.200	0.197 (0.061 _{0.055})	1.489 96.400	0.198 (0.064 _{0.059})	0.677 97.200	0.199 (0.065 _{0.056})	0.142 97.600
β_{23}	0.311 (0.3) (0.413 _{0.401})	8.133 91.600	0.300 (0.006 _{0.006})	0.138 96.800	0.300 (0.007 _{0.006})	0.056 97.400	0.300 (0.007 _{0.006})	0.022 96.400
β_{30}	0.213 (0.2) (0.293 _{0.246})	7.742 88.200	0.199 (0.141 _{0.128})	0.257 96.000	0.202 (0.153 _{0.132})	1.300 97.200	0.194 (0.155 _{0.130})	2.940 98.000
β_{31}	-0.209 (-0.2) (0.533 _{0.510})	4.300 91.400	-0.199 (0.166 _{0.152})	0.433 96.200	-0.207 (0.181 _{0.155})	3.711 96.600	-0.197 (0.184 _{0.153})	1.227 98.400
β_{32}	0.214 (0.2) (0.429 _{0.412})	6.483 92.000	0.199 (0.059 _{0.055})	0.433 96.800	0.203 (0.064 _{0.056})	0.632 96.200	0.199 (0.066 _{0.055})	0.446 97.600
β_{33}	0.417 (0.4) (0.131 _{0.117})	3.205 94.600	0.400 (0.006 _{0.006})	0.003 96.000	0.399 (0.007 _{0.006})	0.088 96.200	0.400 (0.007 _{0.006})	0.037 97.300
F.norm	7.193		8.745		11.224		11.914	

수를 발생한 모형과 다르게 가정하였기에 발생한 현상이다. 회귀계수에 대한 정확성을 모형별로 비교하면, 반복수가 증가함에 따라 모든 모형의 SE와 SD는 감소한다. 그러나 포함확률의 변화는 반복수가 증가하는 것과 큰 상관이 없는 것을 알 수 있다.

F.norm 기준으로 살펴보았을 때, 반복수 관계없이 Lee 등 (2019) 모형의 F.norm 값이 가장 작다. 이는 Kohli 등 (2016)의 방법으로 난수를 발생시켜도 Lee 등 (2019) 모형이 공분산 행렬을 더 잘 추정한다고 볼 수 있다. 또한, 4개의 모형 모두 반복수가 커짐에 따라 F.norm 값이 작아졌다.

Table 3.4. $N = 500, n = 10$

β	Lee		Kohli(AR)		Kohli(CS)		Kohli(AR-CS)	
	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)
β_{10} (0.3)	0.314 (0.193 _{0.188})	5.842 92.000	0.299 (0.007 _{0.006})	0.249 95.800	0.302 (0.083 _{0.065})	0.941 99.200	0.297 (0.086 _{0.064})	0.813 99.800
β_{11} (-0.1)	-0.104 (0.392 _{0.228})	3.180 92.600	-0.098 (0.060 _{0.050})	1.169 95.400	-0.101 (0.070 _{0.055})	1.814 99.200	-0.098 (0.073 _{0.055})	1.998 99.600
β_{12} (0.2)	0.232 (0.321 _{0.319})	4.843 88.800	0.199 (0.019 _{0.017})	0.131 94.800	0.200 (0.022 _{0.017})	0.401 99.000	0.199 (0.023 _{0.017})	0.324 99.800
β_{13} (0.3)	0.293 (0.132 _{0.131})	2.385 92.400	0.300 (0.001 _{0.001})	0.007 95.000	0.299 (0.002 _{0.001})	0.029 98.600	0.300 (0.002 _{0.001})	0.029 99.400
β_{20} (0.2)	0.208 (0.243 _{0.203})	3.492 92.600	0.199 (0.075 _{0.069})	0.230 95.800	0.195 (0.089 _{0.070})	2.150 98.600	0.203 (0.093 _{0.070})	1.674 98.800
β_{21} (-0.1)	-0.109 (0.293 _{0.249})	6.923 95.000	-0.100 (0.063 _{0.059})	0.402 95.800	-0.097 (0.075 _{0.060})	2.733 97.800	-0.102 (0.078 _{0.059})	2.920 99.200
β_{22} (0.2)	0.206 (0.243 _{0.221})	6.423 91.600	0.199 (0.022 _{0.020})	0.097 96.800	0.198 (0.025 _{0.020})	0.643 98.600	0.200 (0.027 _{0.020})	0.475 99.600
β_{23} (0.3)	0.321 (0.219 _{0.234})	4.312 92.000	0.300 (0.002 _{0.002})	0.017 95.600	0.300 (0.002 _{0.002})	0.055 98.600	0.299 (0.002 _{0.002})	0.030 99.200
β_{30} (0.2)	0.232 (0.143 _{0.148})	8.433 94.600	0.198 (0.054 _{0.043})	0.941 96.800	0.198 (0.064 _{0.048})	0.588 99.200	0.198 (0.066 _{0.049})	0.540 99.200
β_{31} (-0.2)	-0.209 (0.290 _{0.242})	3.543 93.200	-0.199 (0.049 _{0.042})	0.187 96.600	-0.199 (0.056 _{0.043})	0.082 99.200	-0.198 (0.058 _{0.044})	0.671 99.200
β_{32} (0.2)	0.239 (0.092 _{0.083})	5.694 95.400	0.200 (0.015 _{0.013})	0.127 96.400	0.200 (0.017 _{0.013})	0.223 99.000	0.199 (0.018 _{0.014})	0.270 99.400
β_{33} (0.4)	0.418 (0.149 _{0.132})	2.438 95.000	0.399 (0.001 _{0.001})	0.001 96.200	0.399 (0.001 _{0.001})	0.171 98.200	0.400 (0.001 _{0.001})	0.013 99.000
F.norm	3.819		4.527		4.842		4.976	

Tables 3.3과 3.4은 표본의 크기가 500일 때 4가지 모형의 결과를 나타낸 것이다. 그리고 Figure 3.1의 (b)는 모든 평균모수들의 추정정도를 한 번에 파악하기 위해 ARB값의 합을 나타내었다. F.norm 기준으로 살펴보았을 때, 반복수 관계없이 Lee 등 (2019) 모형의 F.norm 값이 가장 작다. 이는 표본의 크기가 100일 때뿐만 아니라 500일 때에도 Lee 등 (2019) 모형이 공분산 행렬을 더 잘 추정한다고 볼 수 있다. 또한, 4개의 모형 모두 반복수가 커짐에 따라 F.norm 값이 작아졌다. 따라서 반복수가 커짐에 따라 모든 모형이 공분산 행렬을 잘 추정한다고 해석할 수 있다. 회귀계수에 대한 정보 기준으로 모형을 비교하면, 반복수가 증가함에 따라 모든 모형의 SE와 SD는 감소한다. 그러나 포함확률의 변화는 반복수가 증가하는 것과 큰 상관이 없는 것을 알 수 있다.

Kohli 등 (2016)이 제안한 AR, CS, AR-CS 모형을 비교하면, 반복수가 적을 때에는 세 모형의 SE, SD 값의 차이가 크다. AR 모형의 SE와 SD가 가장 작을 때가 많으며, AR-CS 모형의 SE와 SD 값이 가장 큰 경우가 많다. 이러한 사실 때문에 AR-CS 모형에서의 포함확률이 가장 크게 나타난다. 그러나 반복수가 커짐에 따라 각 모형의 SE와 SD 차이가 작아지는 것을 알 수 있다.

3.2. Lee 등 (2019)에서 생성한 자료집합

두 번째로 Lee 등 (2019)이 제안한 모형에서 난수를 발생시켰다. 응답변수의 수, 대상자의 수 및 반

Table 3.5. $N = 100, n = 5$

β	Lee		Kohli(AR)		Kohli(CS)		Kohli(AR-CS)	
	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)
β_{10} (0.3)	0.301 (1.155 _{1.152})	0.350 95.400	0.228 (0.124 _{0.145})	23.808 85.800	0.156 (0.480 _{0.391})	10.188 94.200	0.163 (0.482 _{0.383})	9.9780 94.200
β_{11} (-0.1)	-0.100 (0.222 _{0.217})	0.170 94.800	-0.025 (0.730 _{0.492})	47.429 82.000	-0.014 (0.781 _{0.625})	60.813 97.000	-0.016 (0.718 _{0.650})	56.410 96.600
β_{12} (0.2)	0.206 (0.296 _{0.287})	3.086 95.800	0.170 (0.985 _{0.829})	14.591 84.600	0.210 (0.184 _{0.181})	15.318 78.600	0.225 (1.842 _{1.743})	25.186 74.600
β_{13} (0.3)	0.305 (0.423 _{0.421})	1.866 95.200	0.189 (0.821 _{0.814})	36.769 86.200	0.153 (0.152 _{0.137})	41.082 89.800	0.149 (0.152 _{0.142})	45.461 89.000
β_{20} (0.2)	0.202 (0.129 _{0.124})	1.225 97.000	0.128 (0.107 _{0.163})	35.887 87.000	0.061 (0.587 _{0.608})	46.910 91.200	0.125 (0.585 _{0.634})	26.260 92.000
β_{21} (-0.1)	-0.102 (0.184 _{0.183})	2.919 95.200	-0.011 (0.671 _{0.610})	58.841 82.000	-0.026 (0.527 _{0.408})	57.399 98.800	-0.095 (0.527 _{0.419})	13.440 98.800
β_{22} (0.2)	0.205 (0.245 _{0.243})	2.824 94.200	0.176 (0.090 _{0.128})	11.565 83.600	0.181 (0.170 _{0.162})	8.907 89.200	0.113 (1.711 _{1.655})	26.727 89.200
β_{23} (0.3)	0.310 (0.350 _{0.357})	3.347 94.400	0.193 (0.088 _{0.118})	35.441 71.000	0.247 (0.130 _{0.132})	45.332 62.400	0.257 (1.296 _{1.290})	37.365 86.000
β_{30} (0.2)	0.203 (0.137 _{0.140})	1.820 95.400	0.081 (0.123 _{0.128})	59.089 82.800	0.176 (0.499 _{0.471})	10.384 94.200	0.159 (0.501 _{0.482})	14.035 94.000
β_{31} (-0.2)	-0.203 (0.197 _{0.200})	1.866 94.800	-0.053 (0.712 _{0.468})	45.382 92.400	-0.015 (0.455 _{0.351})	77.547 96.400	-0.039 (0.455 _{0.369})	58.033 96.800
β_{32} (0.2)	0.202 (0.261 _{0.262})	1.324 94.000	0.199 (0.099 _{0.107})	0.282 92.400	0.521 (0.209 _{0.195})	87.758 95.600	0.685 (0.209 _{0.198})	99.796 94.600
β_{33} (0.4)	0.409 (0.373 _{0.377})	2.395 94.200	0.207 (0.089 _{0.097})	48.095 82.200	0.863 (0.173 _{0.167})	69.842 95.600	1.067 (1.725 _{1.739})	86.998 93.000
F.norm	4.796		4.802		4.811		4.918	

복수는 소절 3.1에서 제시된 것과 동일하며, 회귀계수 값도 동일하다. 공분산행렬의 구조는 동일한 AR(1)이며 이 공분산행렬을 분해하여 만들어진 모수들인 일반화자기회귀모수, 혁신분산모수 그리고 계구모수의 선형모형인 식 (2.8), (2.10), 그리고 (2.13)에서 만들어지는 모수들의 값은 아래와 같다.

$$\begin{aligned}
 (\alpha_{11}, \alpha_{12}, \alpha_{13}, \alpha_{21}, \alpha_{22}, \alpha_{23}, \alpha_{31}, \alpha_{32}, \alpha_{33}) &= (0.3, 0.4, 0.1, 0.1, 0.3, 0.1, 0.1, 0.3, 0.2), \\
 (\lambda_1, \lambda_2, \lambda_3) &= (0.2, 0.2, 0.2), \\
 (\nu_1, \nu_2, \nu_3) &= (-0.5, -0.4, -0.3).
 \end{aligned}$$

표본의 크기(N)가 100, 500으로 하고 500개의 자료집합을 생성하였다. 이렇게 만들어진 자료집합을 소절 3.1에서 적합시킨 4개의 모형이 다시 적합시켰다. 그 결과로 각 모수의 추정치 평균(mean), 상대편향의 절댓값(ARB), 추정량의 표준오차의 평균(SE), 추정치들의 표준편차(SD), 그리고 포함확률(CP)을 계산하였다. 또한 Kohli 등 (2016)이 제안한 모형과 Lee 등 (2019)이 제안한 모형 간의 비교를 위해 프로베니우스 노름을 다시 사용하여 비교하였다.

Tables 3.5과 3.6은 표본의 크기가 100일 때 4가지 모형의 결과를 표로 나타낸 것이다. 또한, 모든 평균모수들의 추정정도를 한 번에 파악하기 위해 ARB값의 합인 $\sum_{i=1}^3 \sum_{j=0}^3 (|\hat{\beta}_{ij} - \beta_{ij}| / \beta_{ij}) \times 100$ 를 계

Table 3.6. $N = 100, n = 10$

β	Lee		Kohli(AR)		Kohli(CS)		Kohli(AR-CS)	
	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)
β_{10}	0.304 (0.165 _{0.170})	1.424 95.000	0.122 (0.483 _{0.438})	32.536 93.400	0.231 (0.129 _{0.161})	22.950 85.200	0.233 (0.126 _{0.154})	22.048 87.200
β_{11}	-0.102 (0.235 _{0.240})	2.483 95.200	-0.033 (0.719 _{0.698})	36.643 95.000	-0.037 (0.071 _{0.048})	33.234 81.200	-0.025 (0.722 _{0.502})	40.025 78.200
β_{12}	0.231 (0.659 _{0.667})	0.667 94.800	0.126 (0.808 _{0.819})	23.063 72.400	0.164 (0.102 _{0.142})	17.978 83.400	0.156 (0.100 _{0.141})	21.600 79.800
β_{13}	0.285 (0.942 _{0.938})	4.729 95.400	0.553 (0.511 _{0.426})	47.774 87.000	0.184 (0.082 _{0.123})	38.416 65.200	0.184 (0.082 _{0.119})	38.371 64.800
β_{20}	0.206 (0.147 _{0.149})	3.446 95.000	0.190 (0.583 _{0.716})	20.467 89.400	0.124 (0.108 _{0.182})	37.633 72.400	0.124 (0.108 _{0.174})	37.778 72.600
β_{21}	-0.109 (0.209 _{0.213})	9.537 95.000	-0.126 (0.526 _{0.438})	37.305 98.000	-0.055 (0.674 _{0.601})	20.055 71.100	-0.077 (0.673 _{0.613})	17.922 70.000
β_{22}	0.216 (0.577 _{0.578})	8.388 94.600	1.195 (1.689 _{1.624})	67.659 88.000	0.172 (0.090 _{0.140})	13.861 78.800	0.172 (0.090 _{0.134})	13.779 81.200
β_{23}	0.323 (0.824 _{0.821})	7.807 95.600	0.257 (0.977 _{1.005})	57.099 56.600	0.187 (0.092 _{0.129})	37.614 66.600	0.180 (0.089 _{0.128})	39.938 64.000
β_{30}	0.201 (0.154 _{0.163})	0.752 93.800	0.396 (0.502 _{0.543})	20.269 91.800	0.072 (0.127 _{0.141})	63.551 80.600	0.086 (0.125 _{0.138})	56.767 81.000
β_{31}	-0.206 (0.220 _{0.231})	3.306 93.600	-0.093 (0.453 _{0.362})	39.534 98.000	0.104 (0.699 _{0.459})	50.521 84.000	-0.048 (0.705 _{0.477})	79.975 91.100
β_{32}	0.248 (0.610 _{0.622})	24.196 93.800	0.337 (1.058 _{1.083})	68.571 92.400	0.199 (0.103 _{0.121})	0.267 89.400	0.189 (0.101 _{0.119})	5.238 90.000
β_{33}	0.401 (0.872 _{0.878})	0.424 95.200	0.686 (1.005 _{1.052})	47.170 90.200	0.206 (0.092 _{0.110})	48.311 75.400	0.198 (0.090 _{0.108})	50.457 79.600
F.norm	4.121		5.343		6.522		7.131	

산하였고 그 결과를 Figure 3.2의 (a)에 나타내었다. 회귀계수는 우리가 예상하는 바와 같이 Lee 등 (2019)의 모형으로 적합한 경우에 편향이 적었다. 하지만 Kohli 등 (2016)의 모형을 적합시에는 편향이 많이 있음을 알 수 있다. 이는 3.1소절에서 Lee 등 (2019) 모형적합에서 나타나는 편향보다 훨씬 더 큼을 알 수 있다. 이를 통하여 Lee 등 (2019) 모형이 Kohli 등 (2016)의 모형적합 보다 좀 더 강건함(robust)을 알 수 있다.

F.norm 기준으로 살펴보았을 때, 반복수와 관계없이 Lee 등 (2019) 모형의 F.norm 값이 가장 작다. 이는 Kohli 등 (2016)의 방법들보다 Lee 등 (2019) 모형이 공분산 행렬을 더 잘 추정한다고 볼 수 있다. 반복수가 증가함에 따라 Lee 등 (2019) 모형의 F.norm 값은 작아졌지만, Kohli 등 (2016)의 세 모형의 F.norm 값은 커졌다. 회귀계수에 대한 정보 기준으로 모형을 비교하면, 반복수와는 상관없이 Lee 등 (2019) 모형의 포함확률이 가장 큰 것을 알 수 있다. 상대편향의 절댓값 측면에서도 대부분 Lee 등 (2019) 모형에서의 값이 가장 작으므로, Lee 등 (2019) 모형이 회귀계수 또한 잘 추정한다고 해석할 수 있다.

Tables 3.7과 3.8은 표본의 크기가 500일 때 4가지 모형의 결과를 표로 나타낸 것이다. 그리고 모든 평균모수들의 추정정도를 한 번에 파악하기 위해 ARB값의 합인 $\sum_{i=1}^3 \sum_{j=0}^3 (|\hat{\beta}_{ij} - \beta_{ij}| / \beta_{ij}) \times 100$ 를 계산하였고 그 결과를 Figure 3.2의 (b)에 나타내었다. F.norm 기준으로 살펴보았을 때, 반복수와 관계

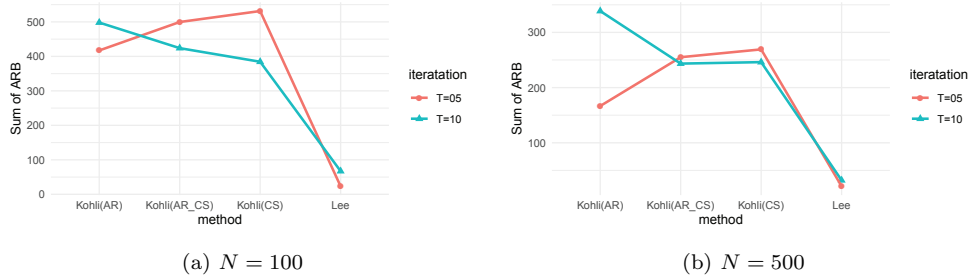


Figure 3.2. Sum of ARBs of mean parameters in Lee *et al.* (2020)'s and Kohli *et al.* (2016)'s models using the data generated from Lee *et al.* (2020)'s model.

Table 3.7. $N = 500, n = 5$

β	Lee		Kohli(AR)		Kohli(CS)		Kohli(AR-CS)	
	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)
β_{10}	0.296 (0.074 _{0.073})	1.114 95.800	0.185 (0.212 _{0.182})	21.165 92.000	0.163 (0.209 _{0.223})	28.308 92.400	0.172 (0.183 _{0.139})	24.333 91.000
β_{11}	-0.098 (-0.1) (0.106 _{0.108})	1.357 94.400	-0.060 (0.317 _{0.287})	33.910 95.400	-0.075 (0.359 _{0.373})	20.845 91.200	-0.081 (0.366 _{0.364})	17.379 91.200
β_{12}	0.198 (0.2) (0.299 _{0.287})	0.532 96.000	0.128 (0.452 _{0.426})	24.082 87.200	0.224 (0.093 _{0.100})	14.396 90.000	0.208 (0.283 _{0.244})	3.385 92.600
β_{13}	0.298 (0.3) (0.423 _{0.411})	0.448 97.000	0.408 (0.673 _{0.600})	23.850 89.000	0.160 (0.158 _{0.139})	39.848 87.600	0.156 (0.228 _{0.283})	40.097 90.000
β_{20}	0.198 (0.2) (0.066 _{0.063})	0.789 96.200	0.194 (0.254 _{0.290})	4.015 92.400	0.173 (0.232 _{0.298})	15.288 88.400	0.181 (0.265 _{0.233})	12.633 88.200
β_{21}	-0.099 (-0.1) (0.945 _{0.954})	0.049 95.000	-0.104 (0.232 _{0.187})	3.594 97.000	-0.088 (0.247 _{0.235})	17.391 91.400	-0.090 (0.173 _{0.169})	20.112 92.400
β_{22}	0.193 (0.2) (0.261 _{0.259})	3.350 94.800	0.134 (0.752 _{0.693})	10.011 86.200	0.179 (0.399 _{0.323})	9.395 87.000	0.125 (0.098 _{0.093})	14.838 88.400
β_{23}	0.303 (0.3) (0.370 _{0.380})	1.313 94.400	0.273 (0.567 _{0.562})	9.494 93.400	0.267 (0.632 _{0.655})	20.934 88.200	0.261 (0.533 _{0.518})	23.776 89.800
β_{30}	0.196 (0.2) (0.070 _{0.067})	1.892 97.200	0.194 (0.220 _{0.215})	3.022 92.400	0.187 (0.238 _{0.277})	7.394 90.400	0.188 (0.163 _{0.132})	8.332 92.200
β_{31}	-0.192 (-0.2) (0.099 _{0.098})	3.771 96.400	-0.179 (0.201 _{0.165})	7.001 88.200	-0.025 (0.211 _{0.236})	46.384 89.000	-0.083 (0.234 _{0.284})	35.365 90.800
β_{32}	0.201 (0.2) (0.276 _{0.272})	0.888 95.000	0.247 (0.406 _{0.388})	8.382 93.600	0.372 (0.118 _{0.106})	32.686 87.200	0.402 (0.100 _{0.094})	36.443 89.200
β_{33}	0.377 (0.4) (0.392 _{0.380})	5.687 95.600	0.663 (0.751 _{0.779})	18.109 88.600	0.636 (0.073 _{0.078})	16.325 89.400	0.688 (0.097 _{0.093})	18.355 90.000
F.norm	4.139		4.852		6.389		8.378	

없이 Lee 등 (2019) 모형의 F.norm 값이 가장 작다. 이는 표본의 크기가 100일 때뿐만 아니라 500일 때에도 Lee 등 (2019) 모형이 공분산 행렬을 더 잘 추정한다고 볼 수 있다. 또한, 4개의 모형 모두 반복수가 커짐에 따라 F.norm 값이 작아졌다. 따라서 반복수가 커짐에 따라 모든 모형이 공분산 행렬을 잘 추정한다고 해석할 수 있다. 회귀계수에 대한 정보 기준으로 모형을 비교하면, 반복수가 증가함에 따라 모든 모형의 SE와 SD는 감소한다. 그러나 포함확률의 변화는 반복수가 증가하는 것과 큰 상관이 없는 것

Table 3.8. $N = 500, n = 10$

β	Lee		Kohli(AR)		Kohli(CS)		Kohli(AR-CS)	
	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)	Mean (SE _{SD})	ARB CP(%)
β_{10}	0.297 (0.3)	0.852 95.000	0.234 (0.058)	21.858 77.000	0.236 (0.073)	19.134 82.000	0.242 (0.078)	16.332 84.200
β_{11}	-0.091 (-0.1)	8.216 94.600	-0.083 (0.341)	10.239 92.000	-0.073 (0.032)	17.389 89.200	-0.071 (0.227)	18.373 88.200
β_{12}	0.202 (0.2)	1.330 95.000	0.160 (0.045)	19.895 76.600	0.163 (0.051)	17.733 78.000	0.158 (0.068)	21.229 81.400
β_{13}	0.290 (0.3)	3.251 96.600	0.184 (0.037)	38.534 84.200	0.186 (0.044)	37.110 81.200	0.187 (0.053)	36.385 83.800
β_{20}	0.195 (0.2)	2.430 96.200	0.126 (0.050)	36.892 86.800	0.153 (0.056)	18.331 82.600	0.151 (0.052)	19.989 87.200
β_{21}	-0.094 (-0.1)	5.104 94.200	-0.124 (0.317)	25.391 86.400	-0.077 (0.277)	10.378 87.600	-0.081 (0.192)	9.341 88.400
β_{22}	0.205 (0.2)	2.677 95.000	0.171 (0.041)	14.112 78.000	0.168 (0.046)	16.776 81.400	0.175 (0.065)	10.378 85.000
β_{23}	0.294 (0.3)	1.7000 96.000	0.183 (0.040)	38.783 76.400	0.198 (0.048)	20.311 81.400	0.197 (0.044)	21.534 83.400
β_{30}	0.196 (0.2)	1.840 95.000	0.079 (0.058)	60.452 79.400	0.093 (0.064)	49.433 82.600	0.097 (0.076)	47.988 86.000
β_{31}	-0.195 (-0.2)	2.171 94.200	-0.132 (0.333)	20.061 92.600	-0.132 (0.208)	19.998 89.600	-0.123 (0.174)	23.423 90.000
β_{32}	0.201 (0.2)	0.975 94.600	0.194 (0.046)	2.947 91.000	0.195 (0.046)	1.733 90.200	0.192 (0.063)	3.112 91.400
β_{33}	0.393 (0.4)	1.535 94.200	0.202 (0.041)	49.496 88.600	0.332 (0.057)	17.777 82.600	0.347 (0.046)	15.379 91.600
F.norm	4.912		9.273		11.634		12.337	

을 알 수 있다.

Kohli 등 (2016)이 제안한 AR, CS, AR-CS 모델을 비교하면, 반복수가 적을 때에는 세 모형의 SE, SD 값의 차이가 크다. AR 모형의 SE와 SD가 가장 작을 때가 많으며, AR-CS 모형의 SE와 SD 값이 가장 큰 경우가 많다. 이러한 사실 때문에 AR-CS 모형에서의 포함확률이 가장 크게 나타난다. 그러나 반복수가 적어짐에 따라 각 모형의 SE와 SD 차이가 작아지는 것을 알 수 있다.

4. 결론

본 논문에서 다변량 경시적 자료분석을 위해 공분산행렬을 모형화 하는 여러 방법을 고찰하였다. Pourahmadi (1999)는 단변량 경시적 자료의 공분산 행렬을 모형화하기 위해 MCD를 이용하였다. MCD에 의하여 만들어진 모수들은 선형/로그선형모형으로 모형화하여 모수의 수를 줄이고 양정치성을 만족시킬 수 있다. 이러한 MCD의 장점 때문에, Kohli 등 (2016)과 Lee 등 (2019)은 MCD를 이용하여 다변량 경시적 자료의 공분산 행렬을 모형화하는 방법론을 제안하였다. 두 방법론 모두 MCD를 이용하였다는 공통점이 있지만 혁신공분산행렬을 모형화하는 방법이 다르다는 차이점이 있다. 우선 Kohli 등 (2016)은 혁신공분산행렬을 직접적으로 추정하지 않고, 이를 미지의 음수가 아닌 모수와 알려진 양정치

행렬의 선형결합 식으로 변환함으로써 추정하였다. 따라서 알려진 양정치 행렬을 어떤 구조에서 모수적 모형화 하느냐에 따라 AR / CS / AR-CS 유형으로 나뉜다. Lee 등 (2019) 또한 혁신공분산행렬을 혁신분산을 대각원소로 가지는 대각행렬과 각 시점에서의 반응변수에 대한 상관행렬의 식으로 분해함으로써 해석 가능한 추정을 하였다.

Kohli 등 (2016)과 Lee 등 (2019)의 성능을 비교하기 위해 모의실험을 진행하였다. 회귀계수와 공분산행렬을 얼마나 잘 추정했는지 확인하기 위해 평균 회귀계수들의 편향, 포함확률과 프로베니우스 노름 등을 계산하였다. 표본의 크기와 반복수와 관계없이 Lee 등 (2019)의 모형이 프로베니우스 노름 값이 가장 작았다. 즉, 주어진 상황 내에서는 Lee 등 (2019)이 공분산 행렬을 가장 정확하게 추정한다는 것을 알 수 있다. 그리고 Lee 등 (2019)은 회귀계수에 추정에서 Kohli 등 (2016)보다 강건함을 보이고 있음을 알 수 있었다.

References

- Kim, C. and Zimmerman, D. L. (2012). Unconstrained models for the covariance structure of multivariate longitudinal data, *Journal of Multivariate Analysis*, **107**, 104–118.
- Kohli, P., Garcia, T. P., and Pourahmadi, M. (2016). Modeling the Cholesky factors of covariance matrices of multivariate longitudinal data, *Journal of Multivariate Analysis*, **145**, 87–100.
- Lee, K., Cho, H., Kwak, M. S. and Jang, E. J. (2019). Estimation of covariance matrix of multivariate longitudinal data using modified Choleksy and hypersphere decompositions, *Biometrics*, **76**, 75–86.
- Pourahmadi, M. (1999). Joint mean-covariance models with applications to longitudinal data: Unconstrained parameterisation, *Biometrika*, **86**, 677–690.
- Qu, A., Lindsay, B. G., and Li, B. (2000). Improving generalised estimating equations using quadratic inference functions, *Biometrika*, **87**, 823–836.
- Xu, J. and Mackenzie, G. (2012). Modelling covariance structure in bivariate marginal models for longitudinal data, *Biometrika*, **99**, 649–662.
- Zhou, J. and Qu, A. (2012). Informative estimation and selection of correlation structure for longitudinal data, *Journal of the American Statistical Association*, **107**, 701–710.

다변량 경시적 자료 분석을 위한 공분산 행렬의 모형화 비교 연구

곽나영^a · 이근백^{a,1}

^a성균관대학교 통계학과

(2020년 3월 5일 접수, 2020년 4월 2일 수정, 2020년 4월 9일 채택)

요약

같은 개체로부터 반복 측정된 자료를 경시적 자료(longitudinal data)라고 한다. 이러한 자료를 분석하려면 흔히 사용되는 횡단 자료 분석과는 다른 분석 방법이 필요하다. 즉, 경시적 자료에서 공변량의 효과를 추정할 때에는 반복 측정된 결과 간의 상관성을 고려해야 하며, 따라서 공분산행렬을 모형화 하는 것이 매우 중요하다. 그러나 추정해야 할 모수가 많고, 추정된 공분산행렬이 양정치성을 만족해야 하므로 공분산 행렬의 모형화는 쉽지 않다. 특히 다변량 경시적 자료분석을 위한 공분산행렬의 모형화는 더욱더 심층적인 방법론을 사용해야 한다. 본 논문은 다변량 경시적 자료분석을 위한 공분산행렬을 모형화하기 위해 두 가지 방법론을 고찰한다. 두 방법 모두 수정된 콜레스키 분해(modified Cholesky decomposition)를 이용하여 시간에 따른 응답변수들의 상관관계를 설명하고 있다. 하지만 같은 시간에서 관측된 응답변수들간의 상관관계를 설명하는 방법이 다르다. 첫 번째 방법론에서는 향상된 선형 공분산 모형(enhanced linear covariance models)을 사용하여 공분산행렬이 양정치성을 만족하도록 한다. 두 번째 방법론에서는 분산-공분산 분해(variance-correlation decomposition)와 초구분해(hypersphere decomposition)을 이용하여 공분산 행렬을 모형화 한다. 이 두 방법론의 성능을 비교하고자 모의실험을 진행한다.

주요용어: 다변량 경시적 자료, 공분산 행렬 모형화, 수정된 콜레스키 분해, 양정치성, 초구분해

이 논문은 정부(교육부)의 재원으로 한국연구재단의 지원을 받아 수행된 기초연구사업임 (과제번호: NRF-2019R1F1A1058553).

¹교신저자: (03063) 서울특별시 종로구 성균관로 25-2, 성균관대학교 통계학과. E-mail: keunbaik@skku.edu