

# 내포오차성분을 가정한 패널회귀모형에서 추정량의 효율에 관한 비교

송석헌<sup>1)</sup> 전명식<sup>2)</sup> 정병철<sup>3)</sup>

## 요약

본 논문에서는 내포오차성분을 가지는 패널회귀모형에서 회귀계수에 대하여 다양한 추정량들을 유도하고, 추정량들의 효율성을 모의실험을 통하여 평균제곱오차의 기준에서 비교하였다. 모의실험 결과, 제안된 FGLS추정량들은 GLS추정량과 효율성에서 서로 큰 차이를 보이지 않았으며, 계산상 더욱 복잡한 ML, REML추정량 및 MIVQUE와 거의 비슷한 효율성을 보여주었다.

주요 용어: 패널회귀모형, FGLS, ML, REML, MIVQUE.

## 1. 서론

내포오차성분모형(nested error components model)은 실험계획에서 내포된 실험에 대하여 주로 사용하는 모형이다. 이와 같은 모형에서의 주된 관심은 각 그룹의 효과와 내포된 그룹의 효과를 나타내는 분산성분들의 추정에 있다. Sahai (1976)는 내포오차성분모형에서 분산성분에 대한 여러 추정량들의 평균제곱오차를 비교하였으며, Graybill과 Wang (1979), Wang과 Graybill (1981)과 Burdick과 Graybill (1985)등은 여러가지 분산성분의 비(ratio)에 대한 신뢰구간에 대하여 연구하였다. 그러나, 경영학이나, 경제학, 교육학 등 여러분야에서 내포오차성분모형을 회귀모형의 관점에서 적용해야 하는 경우가 흔히 발생한다 (Fuller와 Battese (1973), Pantula와 Pollock (1985), Baltagi (1995)).

실험계획의 관점과는 달리 회귀분석의 관점에서는 주요 관심이 분산성분의 추정보다는 설명변수와 반응변수의 관련성을 나타내는 회귀계수의 추정에 있게 된다. 그러나, 패널회귀모형에서 분산성분이 알려지지 않은 경우라면 회귀계수에 대한 일반화 최소제곱(generalized least squares, GLS)추정량은 구할 수 없게 된다. 이 경우에는 일반적으로 일차적으로 분산성분을 추정하고, 추정된 분산성분을 이용하여 회귀계수를 추정하는 '추정가능' 일반화 최소제곱(feasible generalized least squares, FGLS)추정법이 이용된다.

이에 본 연구에서는 분산성분의 추정을 위하여 분산분석표에서 얻어지는 여러형태의 ANOVA추정량들을 유도하여 이를 이용한 여러형태의 FGLS추정량들을 제안하려 한다. 또

1) (136-701) 서울시 성북구 안암동 5가 1, 고려대학교 통계학과 조교수

E-mail: ssong@kucncx.korea.ac.kr

2) (136-701) 서울시 성북구 안암동 5가 1, 고려대학교 통계학과 교수

E-mail: jhun@kucncx.korea.ac.kr

3) (136-701) 서울시 성북구 안암동 5가 1, 고려대학교 통계연구소 선임연구원

E-mail: bcjung@kustat.korea.ac.kr

한 최대우도(maximum likelihood, ML), 제한적 최대우도(restricted maximum likelihood, REML) 추정법 및 최소분산 이차불편추정법(minimum variance quadratic unbiased estimation, MIVQUE)을 이용하여 회귀계수에 대한 추정량들을 제시하고, 이러한 다양한 추정량의 효율성을 모의실험을 통하여 비교하고자 한다. 더불어 분산성분의 효율적인 추정이 회귀계수의 추정에 미치는 효과도 알아보하고자 한다.

이 논문의 구성은 다음과 같다. 먼저 2장에서는 내포오차성분을 따르는 패널회귀모형을 다루고, 3장에서는 이러한 모형에서의 회귀계수에 대한 여러형태의 추정량들을 유도한다. 4장에서는 여러 추정량들의 효율성을 평균제곱오차의 기준에서 비교하기 위하여 모의 실험을 실시한다. 마지막으로 5장에서는 결론을 정리한다.

## 2. 모형

다음과 같은 패널회귀모형을 고려해보자.

$$y_{ijt} = x'_{ijt}\beta + u_{ijt}, \quad i = 1, 2, \dots, M, \quad j = 1, 2, \dots, N, \quad t = 1, 2, \dots, T \quad (2.1)$$

여기서  $y_{ijt}$ 는  $i$ 번째 그룹내에 내포된  $j$ 번째 그룹의 시점  $t$ 에서의 관측치를 나타내는 반응변수이고,  $x_{ijt}$ 는  $(k+1)$ 개의 변수로 이루어진 설명변수벡터이다. 모형 (2.1)에서 오차항  $u_{ijt}$ 는 다음과 같은 내포오차성분모형(nested error components model)을 따른다고 가정하자(Fuller와 Battese (1973), Baltagi (1995)).

$$u_{ijt} = \mu_i + \nu_{ij} + \varepsilon_{ijt}, \quad (2.2)$$

여기서  $\mu_i$ 는 관측될 수 없는 그룹효과(group specific effect)를 나타내는 확률변수이고  $\nu_{ij}$ 는 내포된 그룹효과(nested subgroup effect)를 나타내는 확률변수이며,  $\varepsilon_{ijt}$ 는 나머지 오차항을 나타낸다.  $\mu_i$ ,  $\nu_{ij}$ 와  $\varepsilon_{ijt}$ 는 서로 독립이며 각각  $\mu_i \stackrel{i.i.d.}{\sim} N(0, \sigma_\mu^2)$ ,  $\nu_{ij} \stackrel{i.i.d.}{\sim} N(0, \sigma_\nu^2)$ 와  $\varepsilon_{ijt} \stackrel{i.i.d.}{\sim} N(0, \sigma_\varepsilon^2)$ 라고 가정하자. 만일 모형 (2.1)에서 설명변수가 상수항만 존재하는 경우라면, 이는 실험계획모형에서 통상적으로 사용하는 내포계획모형이 된다(Sahai (1976), Graybill과 Wang (1979), Burdick과 Graybill (1985)). 모형 (2.1)을 행렬을 이용하여 표현하면 다음과 같다.

$$y = X\beta + u \quad (2.3)$$

$y$ 는  $MNT \times 1$  반응변수벡터,  $X$ 는  $MNT \times (k+1)$  설명변수행렬이고  $\beta$ 는  $(k+1) \times 1$ 인 회귀계수벡터를 나타내며, 더불어 식 (2.2)의 오차항은 다음과 같이 표현된다.

$$u = Z_\mu\mu + Z_\nu\nu + \varepsilon \quad (2.4)$$

$Z_\mu = I_M \otimes i_N \otimes i_T$ ,  $Z_\nu = I_M \otimes I_N \otimes i_T$ ,  $I_N$ 과  $I_T$ 는 각각  $N \times N$ 과  $T \times T$ 인 단위행렬이며,  $i_N$ 과  $i_T$ 는 각각 모든 원소가 1인  $N \times 1$ 과  $T \times 1$ 인 벡터를 나타내며,  $\mu = (\mu_1, \dots, \mu_N)'$ ,  $\nu = (\nu_{11}, \dots, \nu_{MN})'$ ,  $\varepsilon = (\varepsilon_{111}, \dots, \varepsilon_{MNT})'$ 이며  $\otimes$ 는 크로네커 곱(Kronecker product)을 나타낸다. 식 (2.4)의 오차항의 분산-공분산행렬은 다음과 같이 주어진다.

$$\Omega = (I_M \otimes J_N \otimes J_T)\sigma_\mu^2 + (I_M \otimes I_N \otimes J_T)\sigma_\nu^2 + (I_M \otimes I_N \otimes I_T)\sigma_\varepsilon^2 \quad (2.5)$$

$J_N$ 과  $J_T$ 는 각각 모든 원소가 1인  $N \times N$ 과  $T \times T$ 인 행렬이다. 식 (2.5)에서  $\bar{J}_N = J_N/N$ ,  $\bar{J}_T = J_T/T$ 라 하고  $E_N = I_N - \bar{J}_N$ ,  $E_T = I_T - \bar{J}_T$ 라 할 때,  $I_N$ 을  $E_N + \bar{J}_N$ 로  $I_T$ 를  $E_T + \bar{J}_T$ 로 치환하면 식 (2.5)의  $\Omega$ 는 다음과 같이 표현할 수 있다.

$$\Omega = \sigma_\varepsilon^2 Q_1 + \sigma_2^2 Q_2 + \sigma_3^2 Q_3 \quad (2.6)$$

여기서  $\sigma_2^2 = T\sigma_\nu^2 + \sigma_\varepsilon^2$ ,  $\sigma_3^2 = NT\sigma_\mu^2 + T\sigma_\nu^2 + \sigma_\varepsilon^2$ 이고,  $Q_1 = (I_M \otimes I_N \otimes E_T)$ ,  $Q_2 = (I_M \otimes E_N \otimes \bar{J}_T)$ ,  $Q_3 = (I_M \otimes \bar{J}_N \otimes \bar{J}_T)$ 이다. 또한  $\sigma_\varepsilon^2$ ,  $\sigma_2^2$ 와  $\sigma_3^2$ 는 각각  $\Omega$ 의  $MN(T-1)$ ,  $M(N-1)$ ,  $M$ 개의 서로 다른 고유값이 되며, 각  $Q_i, i = 1, 2, 3$ 는 대칭멱등행렬이고 서로 직교하며 각 행렬의 합은 단위행렬이 된다. 이와 같은 분광분해(spectral decomposition)(Wansbeek과 Kapteyn (1982, 1983))를 이용하면  $\Omega^p$ 은 다음과 같이 표현된다.

$$\Omega^p = (\sigma_\varepsilon^2)^p Q_1 + (\sigma_2^2)^p Q_2 + (\sigma_3^2)^p Q_3 \quad (2.7)$$

$p$ 는 임의의 상수이고, 식 (2.7)을 이용하면  $\Omega$ 의 역행렬과 변환행렬  $\Omega^{-1/2}$ 를 다음과 같이 구할 수 있다.

$$\begin{aligned} \Omega^{-1} &= (\sigma_\varepsilon^2)^{-1} Q_1 + (\sigma_2^2)^{-1} Q_2 + (\sigma_3^2)^{-1} Q_3 \\ \Omega^{-1/2} &= \frac{1}{\sigma_\varepsilon} Q_1 + \frac{1}{\sigma_2} Q_2 + \frac{1}{\sigma_3} Q_3 \end{aligned} \quad (2.8)$$

### 3. 회귀계수에 대한 추정량들

#### 3.1. OLS, WITHIN, GLS 추정량

모형 (2.3)에서  $\beta$ 에 대한 보통최소제곱추정량(OLS추정량)은 다음과 같이 구해진다.

$$\hat{\beta}_{OLS} = (X'X)^{-1}X'y \quad (3.1)$$

식 (3.1)의 OLS추정량은 불편추정량이며 일치추정량이지만 분산성분들의 존재를 무시하였기 때문에 OLS추정량의 효율은 떨어질 것이다. 또한 OLS 추정에 의하여 추정된 표준 오차는 편향된다는 사실이 잘 알려져 있다(Moulton (1986)). OLS추정량에 의한 잔차는  $\hat{u}_{OLS} = y - X\hat{\beta}_{OLS}$ 이다.

회귀계수에 대한 내부변환추정량(within transformation estimator, 이하 WTN추정량)은 식 (2.3)의 원모형에 대하여  $Q_1$ 을 곱하여 그룹효과와 내포된 그룹효과를 제거하고, 변환된 모형에 OLS를 적용하여 얻어지는 추정량으로, 이는 다음과 같이 구해진다(Hsiao (1986)).

$$\tilde{\beta}_s = (X'_s Q_1 X_s)^{-1} X'_s Q_1 y \quad (3.2)$$

여기서  $X_s$ 는 설명변수행렬  $X$ 에서 상수항을 제거한  $(MNT \times k)$  행렬이고  $\beta_s$ 는 상수항을 제외한  $(k \times 1)$  회귀계수벡터이다. 그러므로 식 (3.3)의 WTN 추정량은  $(y_{ijt} - \bar{y}_{ij.})$ 에 대한  $(X_{ijt} - \bar{X}_{ij.})$ 의 회귀에 의하여 추정되는 추정량이다. 여기서  $\bar{y}_{ij.} = \sum_{t=1}^T y_{ijt}/T$ ,  $\bar{X}_{ij.} = \sum_{t=1}^T X_{ijt}/T$ 이다.  $\beta' = (\alpha, \beta'_s)$ 에서 상수항  $\alpha$ 는  $\tilde{\alpha} = \bar{y}_{...} - \bar{X}'_{...s} \tilde{\beta}_s$ 로 추정된다. 여기서

$\bar{y}_{...} = \sum \sum \sum y_{ijt}/MNT$ ,  $\bar{X}_{...s} = \sum \sum \sum X_{ijts}/MNT$ 이다. 그러므로 WTN추정량에 의한 내부변환잔차(WTN잔차)는 다음과 같이 구해진다.

$$\tilde{u}_{WTN} = y - X_s \tilde{\beta}_s - (\bar{y}_{...} - \bar{X}_{...s} \tilde{\beta}_s) \quad (3.3)$$

만일 식 (2.6)에서 분산성분들이 알려져 있다면,  $\beta$ 에 대한 GLS추정량은 다음과 같이 구해진다.

$$\hat{\beta}_{GLS} = (X' \Omega^{-1} X)^{-1} X' \Omega^{-1} y \quad (3.4)$$

그러나, 현실적으로 분산성분들은 알려지지 않은 경우가 대부분이므로 GLS 추정량은 단지 이론적인 추정량에 불과하지만 다른 추정량과의 비교시 기준이 되는 추정량으로 사용될 수 있다.

### 3.2. FGLS 추정량들

만일 식 (2.6)에서 분산성분들이 알려져 있지 않으면, 분산성분들을 먼저 추정하고 이들을 식 (3.4)에 대입하여 회귀계수를 추정하는 방법을 이단계 추정방법이라하며 이와 같은 방법으로 얻어지는 추정량을 FGLS추정량이 한다.

$$\hat{\beta}_{FGLS} = (X' \hat{\Omega}^{-1} X)^{-1} X' \hat{\Omega}^{-1} y \quad (3.5)$$

여기서  $\hat{\Omega} = \hat{\sigma}_\epsilon^2 Q_1 + \hat{\sigma}_2^2 Q_2 + \hat{\sigma}_3^2 Q_3$ 이다. FGLS추정량은 점근적으로 GLS추정량과 동일한 성질을 갖게되지만 일반적으로 소표본인 경우에는 정확한 성질이 알려져 있지 않다(Greene (1997)).

분산성분들을 추정하는 방법으로 Balestra (1973)는 이원오차성분모형(two-way error component model)에서 분산성분에 대한 최량이차불편추정량(best quadratic unbiased estimator, BQUE)을 유도하였다. Balestra의 결과를 본 모형에 적용하면 다음과 같은 추정량들을 얻을 수 있다.

$$\begin{aligned} \hat{\sigma}_\epsilon^2 &= \frac{u' Q_1 u}{MN(T-1)} = \frac{1}{MN(T-1)} \sum_{i=1}^M \sum_{j=1}^N \sum_{t=1}^T (u_{ijt} - \bar{u}_{ij.})^2 \\ \hat{\sigma}_2^2 &= \frac{u' Q_2 u}{M(N-1)} = \frac{T}{M(N-1)} \sum_{i=1}^M \sum_{j=1}^N (\bar{u}_{ij.} - \bar{u}_{i..})^2 \\ \hat{\sigma}_3^2 &= \frac{u' Q_3 u}{M} = \frac{NT}{M} \sum_{i=1}^M \bar{u}_{i..}^2 \end{aligned} \quad (3.6)$$

그러나, 식 (3.6)에서 제시한 분산성분 추정량은 알려지지 않은 오차항  $u$ 의 함수이므로 추정이 가능하지 않게 된다. 이에 대한 대안으로 본 절에서는 이원오차성분모형에서 제안되었던 여러가지 분산성분 추정방법을 확장, 수정하여 그에 따라 얻어지는 FGLS추정량들을 제안하려 한다.

1) 수정된 Wallace와 Hussain추정량(WH추정량)

Wallace와 Hussain (1969)은 이원오차성분모형에서 분산성분들을 추정하기 위하여 식 (3.6)에서  $u$  대신 OLS잔차  $\hat{u}_{OLS}$ 를 사용할것을 제안하였다. 이러한 제안을 적용하면 다음과 같이 수정된 분산성분의 추정량들을 얻을 수 있다.

$$\begin{aligned} \hat{\sigma}_\varepsilon^2 &= \frac{\hat{u}'_{OLS} Q_1 \hat{u}_{OLS}}{MN(T-1)} = \frac{1}{MN(T-1)} \sum_{i=1}^M \sum_{j=1}^N \sum_{t=1}^T (\hat{u}_{ijt} - \bar{\hat{u}}_{ij})^2 \\ \hat{\sigma}_2^2 &= \frac{\hat{u}'_{OLS} Q_2 \hat{u}_{OLS}}{M(N-1)} = \frac{T}{M(N-1)} \sum_{i=1}^M \sum_{j=1}^N (\hat{u}_{ij.} - \bar{\hat{u}}_{i..})^2 \\ \hat{\sigma}_3^2 &= \frac{\hat{u}'_{OLS} Q_3 \hat{u}_{OLS}}{M} = \frac{NT}{M} \sum_{i=1}^M \hat{u}_{i..}^2 \end{aligned} \quad (3.7)$$

식 (3.7)의 추정된 분산성분들을 식 (3.5)에 대입하면 FGLS추정량을 구할 수 있다. 이를 WH추정량이라 하자.

2) 수정된 Amemiya-Type추정량(AM추정량)

Amemiya (1971)는 분산성분들을 추정하기 위하여 식 (3.6)의 오차항 대신에 WTN잔차를 사용할것을 제안하였는데, 이를 적용하면 다음과 같은 분산성분들의 추정량들을 얻을 수 있다.

$$\begin{aligned} \hat{\sigma}_\varepsilon^2 &= \frac{\tilde{u}'_{WTN} Q_1 \tilde{u}_{WTN}}{MN(T-1)} \\ \hat{\sigma}_2^2 &= \frac{\tilde{u}'_{WTN} Q_2 \tilde{u}_{WTN}}{M(N-1)} \\ \hat{\sigma}_3^2 &= \frac{\tilde{u}'_{WTN} Q_3 \tilde{u}_{WTN}}{M} \end{aligned} \quad (3.8)$$

식 (3.8)을 식(3.5)에 대입하여 얻어지는 FGLS추정량을 AM추정량이라 하도록 한다.

3) 수정된 Swamy와 Arora추정량(SA추정량)

Swamy와 Arora (1972)는 이원오차성분모형에서 분산성분들을 추정하기 위하여 내부 변환(within transformation)을 이용한 회귀, 개체변이를 이용한 회귀(between individual regression)와 시간변이를 이용한 회귀(between time regression) 등 3개의 회귀모형을 적합하였으며, 각 모형의 평균제곱오차로 분산성분을 추정하였다. Swamy와 Arora의 결과를 모

형 (2.1)에 확장 적용하면 다음과 같은 분산성분들에 관한 추정량들을 유도 할 수 있다.

$$\begin{aligned}\hat{\sigma}_\epsilon^2 &= \frac{y'Q_1y - y'Q_1X_s(X'_sQ_1X_s)^{-1}X'_sQ_1y}{MN(T-1) - k} \\ \hat{\sigma}_2^2 &= \frac{y'Q_2y - y'Q_2X_s(X'_sQ_2X_s)^{-1}X'_sQ_2y}{M(N-1) - k} \\ \hat{\sigma}_3^2 &= \frac{y'Q_3y - y'Q_3X(X'_sQ_3X)^{-1}X'_sQ_3y}{M - k - 1}\end{aligned}\quad (3.9)$$

식 (3.9)를 이용하여 회귀계수를 추정하는 FGLS추정량을 SA추정량이라 하도록 한다.

#### 4) Henderson's method 3 (Fuller와 Battese(1973), HFB추정량)

Fuller와 Battese (1973)는 내포된 오차성분모형에서 '상수적합방법 (fitting constants method)' 을 이용하여 분산성분들을 다음과 같이 추정하였다.

$$\begin{aligned}\hat{\sigma}_\epsilon^2 &= \frac{\tilde{u}'_{WTN}\tilde{u}_{WTN}}{MN(T-1) - k} \\ \hat{\sigma}_\nu^2 &= \frac{\hat{u}_2^*\hat{u}_2^* - (MN(T-1) - k)\hat{\sigma}_\epsilon^2}{MNT - MT - \text{tr}\left[\left(X'_sZ_\nu Z'_\nu Q_2 X_s\right)\left(X'_s(Q_1 + Q_2)X_s\right)^{-1}\right]} \\ \hat{\sigma}_\mu^2 &= \frac{\hat{u}'_{OLS}\hat{u}_{OLS} - (MNT - k - 1)\hat{\sigma}_\epsilon^2 - \left[MNT - \text{tr}\{X'Z_\nu Z'_\nu X(X'X)^{-1}\}\right]\hat{\sigma}_\nu^2}{MNT - \text{tr}\left[(X'Z_\mu Z'_\mu X)(X'X)^{-1}\right]}\end{aligned}\quad (3.10)$$

여기서  $\hat{u}_2^*$ 은  $(y_{ijt} - \bar{y}_{i.})$ 에 대한  $(X_{ijt} - \bar{X}_{i.})$ 의 회귀에 의하여 얻어지는 잔차이다. 식 (3.10)을 이용하여 구해지는 FGLS추정량을 HFB추정량이라 하도록 한다.

### 3.3. 최대우도추정량(ML추정량)

ML 추정방법은 회귀계수와 분산성분들을 동시에 추정한다. Wansbeek과 Kapteyn (1982, 1983)의 분광분해를 이용하면  $|\Omega| = (\sigma_\epsilon^2)^{MN(T-1)}(\sigma_2^2)^{M(N-1)}(\sigma_3^2)^M$ 이 되며, 로그우도함수는 다음과 같다.

$$\begin{aligned}\log L &= -\frac{MNT}{2}\log(2\pi) - \frac{MNT}{2}\log\sigma_\epsilon^2 - \frac{M(N-1)}{2}\log(T\rho_2 + 1) \\ &\quad - \frac{M}{2}\log(NT\rho_1 + T\rho_2 + 1) - \frac{1}{2\sigma_\epsilon^2}u'\Sigma^{-1}u\end{aligned}\quad (3.11)$$

여기서  $\rho_1 = \sigma_\mu^2/\sigma_\epsilon^2$ ,  $\rho_2 = \sigma_\nu^2/\sigma_\epsilon^2$ 이고  $\Sigma = \Omega/\sigma_\epsilon^2$ 이다. 식 (3.11)의 우도함수에 대하여  $\hat{\rho}_1$ 과  $\hat{\rho}_2$ 가 주어졌을 때  $\sigma_\epsilon^2$ 과  $\beta$ 에 대한 일차미분은 대수적인 해를 제공해 주는데, 이는 다음과 같이 구해진다.

$$\begin{aligned}\hat{\beta}_{ML} &= (X'\hat{\Sigma}^{-1}X)^{-1}X'\hat{\Sigma}^{-1}y \\ \hat{\sigma}_\epsilon^2 &= (y - X\hat{\beta}_{ML})'\hat{\Sigma}^{-1}(y - X\hat{\beta}_{ML})/MNT.\end{aligned}\quad (3.12)$$

그러나,  $\rho_1$ 과  $\rho_2$ 에 대한 일차미분은  $\hat{\beta}$ 와  $\hat{\sigma}_\epsilon^2$ 이 주어졌다 하더라도 다음과 같은 비선형의 형태를 보인다.

$$\begin{aligned} \frac{\partial \log L}{\partial \rho_1} &= -\frac{MNT}{2(1+T\rho_2+NT\rho_1)} + \frac{NT}{2\sigma_\epsilon^2(1+T\rho_2+NT\rho_1)^2} u'Q_3u = 0 \\ \frac{\partial \log L}{\partial \rho_2} &= -\left(\frac{M(N-1)T}{2(1+T\rho_2)} + \frac{MT}{2(1+T\rho_2+NT\rho_1)}\right) \\ &\quad + \left(\frac{T}{2\sigma_\epsilon^2(1+T\rho_2)^2} u'Q_2u + \frac{T}{2\sigma_\epsilon^2(1+T\rho_2+NT\rho_1)^2} u'Q_3u\right) = 0 \end{aligned} \quad (3.13)$$

그러므로  $\hat{\rho}_1$ 과  $\hat{\rho}_2$ 에 대해서는 반복에 의한 수치적인 해를 구해야 하는데, 본 연구에서는 Fisher의 'Scoring' 방법을 사용하였다. 이에 필요한 정보행렬은 다음과 같이 구해진다.

$$\begin{aligned} E\left[-\frac{\partial^2 \log L}{\partial \rho_1^2}\right] &= \frac{MN^2T^2}{2(1+T\rho_2+NT\rho_1)^2} \\ E\left[-\frac{\partial^2 \log L}{\partial \rho_1 \partial \rho_2}\right] &= \frac{MNT^2}{2(1+T\rho_2+NT\rho_1)^2} \\ E\left[-\frac{\partial^2 \log L}{\partial \rho_2^2}\right] &= \frac{M(N-1)T^2}{2(1+T\rho_2)^2} + \frac{MT^2}{2(1+T\rho_2+NT\rho_1)^2} \end{aligned} \quad (3.14)$$

적당한 초기값에서 시작하여  $(n+1)$ 번째 단계에서 갱신되는(updated)  $\hat{\rho}_1$ 과  $\hat{\rho}_2$ 는 다음과 같이 구해진다.

$$\begin{bmatrix} \hat{\rho}_1 \\ \hat{\rho}_2 \end{bmatrix}_{n+1} = \begin{bmatrix} \hat{\rho}_1 \\ \hat{\rho}_2 \end{bmatrix}_n + \begin{bmatrix} E\left[-\frac{\partial^2 \log L}{\partial \rho_1^2}\right] & E\left[-\frac{\partial^2 \log L}{\partial \rho_1 \partial \rho_2}\right] \\ E\left[-\frac{\partial^2 \log L}{\partial \rho_1 \partial \rho_2}\right] & E\left[-\frac{\partial^2 \log L}{\partial \rho_2^2}\right] \end{bmatrix}_n^{-1} \begin{bmatrix} \frac{\partial \log L}{\partial \rho_1} \\ \frac{\partial \log L}{\partial \rho_2} \end{bmatrix}_n \quad (3.15)$$

각 단계에서  $\partial \log L / \partial \rho_1$ 과  $\partial \log L / \partial \rho_2$ 는 식 (3.13)에 의하여 구하고,  $\hat{\beta}$ 와  $\hat{\sigma}_\epsilon^2$ 은 식 (3.12)에 의하여 구해지며, 정보행렬은 식 (3.14)에 의하여 구한다. 만일 갱신된  $\rho_1$ 과  $\rho_2$ 의 값이 0보다 작게 되면, 이는 0으로 대체한다.

### 3.4. 제한적 최대우도추정량(REML추정량)

분산성분에 대한 ML추정량은 회귀계수로 인하여 발생하는 자유도의 손실을 보정하지 않는다는 관점에서 비판을 받아왔다. 이에 대한 대안으로 Patterson과 Thompson (1971)은 비정칙변환(singular transformation)을 이용한 제한적 최대우도추정법(restricted ML추정법, REML추정법)을 제안하였다. 이는 원자료에 대하여  $y'[A : \Sigma^{-1}X/\sigma_\epsilon^2]$ 이라는 변환에 근거한 추정법으로, 여기서 행렬  $A$ 는  $A'A = I - X(X'X)^{-1}X'$ 과  $AA' = I$ 를 만족하는 계수  $(MNT - k - 1) \times MNT$ 인 행렬이다. 이 때  $A'y$ 과  $X'\Sigma^{-1}y/\sigma_\epsilon^2$ 는 독립이 되며 다음과 같은 분포를 따르게 된다.

$$\begin{bmatrix} A'y \\ X'\Sigma^{-1}y/\sigma_\epsilon^2 \end{bmatrix} \sim N\left(\begin{bmatrix} 0 \\ X'\Sigma^{-1}X\beta/\sigma_\epsilon^2 \end{bmatrix}, \begin{bmatrix} \sigma_\epsilon^2 A \Sigma A' & 0 \\ 0 & X'\Sigma^{-1}X/\sigma_\epsilon^2 \end{bmatrix}\right). \quad (3.16)$$

식(3.16)에서  $A'y$ 와  $X'\Sigma^{-1}y/\sigma_\epsilon^2$ 의 로그우도함수를 각각  $\log L_1$ 과  $\log L_2$ 라 하면  $\log L_1$ 은 회귀 계수에 의존하지 않게 된다. 그러므로 이 부분을 최대화하면 분산성분들에 대한 REML 추정량을 얻게 된다. Hocking (1985)에 의하면  $A'(A\Sigma A')^{-1}A' = \Sigma^{-1}[I - X(X'\Sigma^{-1}X)^{-1}X'\Sigma^{-1}]$ 라는 사실을 알 수 있으므로  $\log L_1$ 에 대한  $\sigma_\epsilon^2$ ,  $\rho_1$ 과  $\rho_2$ 에 대한 일차미분값은 다음과 같이 구해진다.

$$\begin{aligned}\frac{\partial \log L_1}{\partial \sigma_\epsilon^2} &= -\frac{MNT - k - 1}{2\sigma_\epsilon^2} + \frac{1}{2\sigma_\epsilon^4}y'\Sigma^{-1}Py, \\ \frac{\partial \log L_1}{\partial \rho_1} &= -\frac{1}{2}tr[Z'_\mu\Sigma^{-1}PZ_\mu] + \frac{1}{2\sigma_\epsilon^2}y'[\Sigma^{-1}PZ_\mu Z_\mu\Sigma^{-1}P]y, \\ \frac{\partial \log L_1}{\partial \rho_2} &= -\frac{1}{2}tr[Z'_\nu\Sigma^{-1}PZ_\nu] + \frac{1}{2\sigma_\epsilon^2}y'[\Sigma^{-1}PZ_\nu Z_\nu\Sigma^{-1}P]y,\end{aligned}\quad (3.17)$$

여기서  $P = I - X(X'\Sigma^{-1}X)^{-1}X'\Sigma^{-1}$ 이다. 식 (3.17)의 값이 0이 되는 해가 모두  $\sigma_\epsilon^2$ ,  $\rho_1$ 과  $\rho_2$ 에 대한 REML 추정량이 된다.  $\rho_1$ 과  $\rho_2$ 가 주어졌을 때 식 (3.17)의 첫 번째 식을 풀면  $\sigma_\epsilon^2$ 에 대한 다음과 같은 추정량을 얻을 수 있다.

$$\hat{\sigma}_\epsilon^2 = y'[\hat{\Sigma}^{-1}P]y/(MNT - k - 1).\quad (3.18)$$

그러나,  $\rho_1$ 과  $\rho_2$ 에 대해서는 대수적인 해를 구할 수 없으므로 반복에 의한 수치적인 해를 구해야만 한다. 여기에서도  $\rho_1$ 과  $\rho_2$ 를 추정하기 위하여 Fisher의 'Scoring' 방법이 사용되었다. 이에 필요한 정보행렬은 Harville (1977)의 결과를 이용하면 쉽게 구할 수 있다.

갱신되는  $\hat{\rho}_1$ 과  $\hat{\rho}_2$ 의 값은 식 (3.15)와 같은 식에 의하여 얻어진다. REML 추정량도 ML 추정량과 마찬가지로 갱신된 값이 0보다 작게 되면, 이는 0으로 대체된다.

### 3.5. MIVQUE

Rao (1971a)는 분산성분을 추정하는 일반적인 처리절차로서 최소노오미차불편추정량 (MINQUE)을 제안하였다. MINQUE 처리절차는 오차항이 정규분포를 따른다는 가정하에서는 최소분산이차불편추정량(MIVQUE)와 동일한 결과를 제공해 주게 된다. 본 연구에서는 정규분포의 가정하에서 MIVQUE 처리절차에 초점을 맞추고자 한다. 분산성분에 대한 선형결합인  $p_\mu\sigma_\mu^2 + p_\nu\sigma_\nu^2 + p_\epsilon\sigma_\epsilon^2$ 의 MIVQUE는 다음과 같은 두가지 조건하에서  $var(y'Gy)$ 를 최소화하는 대칭행렬  $G$ 를 찾음으로서 얻어진다.

$$\begin{aligned}(a) \quad &GX = 0 \\ (b) \quad &tr(GZ_\mu Z_\mu) = p_\mu, \quad tr(GZ_\nu Z_\nu) = p_\nu \quad \text{와} \quad tr(G) = p_\epsilon.\end{aligned}\quad (3.19)$$

위의 두가지 조건은  $y'Gy$ 가 회귀계수  $\beta$ 에 무관하게  $p_\mu\sigma_\mu^2 + p_\nu\sigma_\nu^2 + p_\epsilon\sigma_\epsilon^2$ 의 불편추정량을 제공하는데 필요한 필요충분조건이 된다. 먼저 다음과 같은 사항들을 정의해보자.

$$R = \Sigma^{-1}P/\sigma_\epsilon^2,\quad (3.20)$$

$$S = \{s_{ij}\} = \{tr(V_iRV_jR)\}, \quad i, j = 1, 2, 3 \quad (3.21)$$

$$u = \{u_i\} = \{y'RV_iRy\}, \quad i = 1, 2, 3, \quad (3.22)$$



여기서  $P = I - X(X'\Sigma^{-1}X)^{-1}X'\Sigma^{-1}$  이며,  $V_1 = I_{MNT}$ ,  $V_2 = Z_\mu Z'_\mu$  and  $V_3 = Z_\nu Z'_\nu$ 이다. Rao (1971b)는 분산성분에 대한 MIVQUE는 다음과 같이 계산되는것을 보였다.

$$\hat{\theta} = S^{-1}u, \tag{3.23}$$

여기서  $\theta' = (\sigma_\epsilon^2, \sigma_\mu^2, \sigma_\nu^2)$ 이다. 그러나, MIVQUE는 항상 분산성분에 대한 초기추정값을 요구한다. 수많은 MIVQUE추정량이 초기추정값에 선택에 따라 만들어질 수 있다. 이와 같이 만들어진 추정량은 초기추정값을 참 분산성분의 값으로 주지않는한 최소분산의 성질을 갖지 못하게 된다. 본 연구에서는 두가지 형태의 초기값을 사용하여 분산성분에 대한 MIVQUE를 추정하고, 추정된 분산성분을 이용하여 회귀계수에 대한 MIVQUE형태의 FGLS추정량을 구하였다. 첫 번째 MIVQUE는 모든 분산성분의 초기값을 1로 사용하는 방법으로, 이 결과를 이용한 추정량을 MQ1추정량이라 명명한다(Swallow와 Searle (1978)). 두 번째 MIVQUE는 WH추정량에서 추정된 분산성분을 초기값으로 사용하는 방법인데, 이 결과를 이용한 추정량을 MQ2추정량이라 명명한다(Swallow와 Monahan (1984)). MIVQUE추정에서도 분산성분에 대하여 음의 추정값이 발생할 수 있다. 이 경우에는 음의 추정값을 0으로 대치하고 회귀계수에 대한 MIVQUE형태의 FGLS를 구하였다.

## 4. 모의실험

### 4.1. 모의실험 방법

3장에서 유도한 회귀계수에 대한 여러 추정량들의 효율성을 알아보기 위하여 모의실험을 실시하였다. 비교를 위하여 모의실험에 사용된 모형은 다음과 같은 단순패널 회귀모형이다.

$$\begin{aligned} y_{ijt} &= \alpha + \beta x_{ijt} + u_{ijt}, \quad i = 1, 2, \dots, M, \quad j = 1, 2, \dots, N, \quad t = 1, 2, \dots, T, \\ u_{ijt} &= \mu_i + \nu_{ij} + \epsilon_{ijt} \end{aligned} \tag{4.1}$$

모형 (4.1)에서 설명변수  $x_{ijt}$ 는 Nerlove (1971)에 의해 제안된 방법으로 생성하였다. 즉,  $w_{ijt}$ 가 균일분포  $(-0.5, 0.5)$ 를 따르는 난수라고 했을 때  $x_{ijt}$ 는  $x_{ijt} = 0.1t + 0.5x_{ij,t-1} + w_{ijt}$ 의 식에 의하여 반복적으로 생성하였으며, 이 때 초기값  $x_{ij0}$ 은  $5 + 10w_{ij0}$ 에서 구하였다. 모의실험 전체를 통하여  $\alpha = 5$ ,  $\beta = 0.5$ 와  $\sigma^2 = \sigma_\mu^2 + \sigma_\nu^2 + \sigma_\epsilon^2 = 20$ 으로 고정하였다. 각 실험에서는  $w_1 = \sigma_\mu^2/\sigma^2$ ,  $w_2 = \sigma_\nu^2/\sigma^2$ 의 값을 0.0에서 0.8사이에서 0.2단위로 변화시켜가며 항상  $1 - w_1 - w_2$ 의 값이 0보다 크도록 하였으며,  $\mu_i$ ,  $\nu_{ij}$ 와  $\epsilon_{ijt}$ 의 값은 각각  $N(0, w_1\sigma^2)$ ,  $N(0, w_2\sigma^2)$ 와  $N(0, (1 - w_1 - w_2)\sigma^2)$ 에서 생성하였다. 모든 실험은 1000번 독립적으로 반복 실시하였으며, SAS/IML 프로시저를 이용하여 수행하였다. 본 실험에 사용된 표본은  $M = 5, 10$  인 경우에  $(N, T)$ 의 값을 각각 (5,5) (10,5) (20,5) (5,10) (10,10) (20,10) (5,20) (10,20) (20,20)으로 변화시켜가며 실시하였다. 각 실험에서는  $\beta$ 에 대한 GLS추정량의 평균 제곱오차와 3장에서 유도된 추정량들의 평균제곱오차와의 상대비율과 분산성분들에 대한 ML추정량의 평균제곱오차에 대한 다른 추정량들의 평균제곱오차의 상대비율을 계산하였다.

## 4.2. 모의실험의 결과

### 4.2.1. 회귀계수에 대한 비교

(표 4.1)과 (표 4.2)는  $(M, N, T)$ 가 각각 (5,5,5)인 경우와 (10,20,10)인 경우에  $\beta$ 에 대한 GLS추정량의 평균제곱오차(MSE)에 대한 각 추정량들의 평균제곱오차의 상대비율을  $w_1$ 과  $w_2$ 의 각 수준에서 계산된 값들이다. (표 4.1)과 (표 4.2)의 결과를 요약하면 다음과 같다. 먼저 OLS추정량은  $w_1 = 0$ 과  $w_2 = 0$ 인 경우를 제외한 모든 실험조합에서 효율성이 떨어졌다. 이는 분산성분들을 무시하였기 때문에 발생한 결과로  $w_1$ 과  $w_2$ 의 값이 커질수록 더욱더 비효율적이었다. 특히 본 모형에서는  $w_2$ 의 값이 커지는 것보다  $w_1$ 의 값이 커질때 효율성이 더욱 떨어졌는데, 이는  $w_1$ 은 전체그룹에 영향을 미치는 반면  $w_2$ 는 내포된 그룹에 영향을 미치기 때문인 것으로 판단된다. WTN추정량은  $(M, N, T) = (5, 5, 5)$ 인 경우에 GLS추정량보다 최대 4.4배 큰 평균제곱오차를 보이고 있다. 이 비율은  $M, N, T$ 가 커질수록 감소하여  $(M, N, T) = (10, 20, 10)$ 인 경우에는 GLS추정량보다 최대 2배정도 큰 평균제곱오차를 보이고 있다.  $(M, N, T) = (5, 5, 5)$ 인 경우에 SA추정량과 AM추정량의 평균제곱오차가 다른 WH추정량, HFB추정량이나 ML추정량에 비하여 약간 큰 것으로 나타났으며, WH추정량이나 HFB추정량은 ML추정량, REML추정량이나 MIVQUE추정량과 비슷한 평균제곱오차를 보이고 있다. FGLS추정량(WH, AM, SA, HFB)들은 모든 실험조합에서 GLS보다 최대 11.9%정도 큰 평균제곱오차를 보이는데, 이는  $M, N, T$ 가 증가할수록 감소하는 경향을 보인다.  $(M, N, T) = (10, 20, 10)$ 인 경우에는 GLS와 거의 같은 수준의 평균제곱오차를 보이는 것으로 나타났다. 또한 FGLS추정량들의 평균제곱오차는 계산상으로 좀 더 복잡한 ML, REML추정량이나 MINQUE에 비하여  $M, N, T$ 가 너무 작지않은 경우에는 거의 비슷한 효율성을 보여 주고 있다. 그러므로 회귀계수에 대한 추정에 관한한 평균제곱오차의 기준에서는 계산상으로 복잡한 ML, REML추정량이나 MIVQUE보다는 FGLS추정량들을 사용하더라도 효율성이 떨어지지 않고 있음을 알 수 있다.

### 4.2.2 분산성분에 대한 비교

(표 4.3)과 (표 4.4)는  $(M, N, T) = (10, 20, 10)$ 인 경우에 각각  $\sigma_\mu^2$ 와  $\sigma_v^2$ 의 추정에서 MLE의 평균제곱오차에 대한 각 추정량의 평균제곱오차의 상대 비율을 나타낸다. (표 4.3)과 (표 4.4)의 결과를 살펴보면, 먼저  $\sigma_\mu^2$ 의 추정에 있어서는 SA추정량, HFB추정량, REML추정량, MQ1추정량 및 MQ2추정량의 평균제곱오차가 ML추정량의 평균제곱오차에 비하여 크게 나타났다. 그러나 WH추정량과 AM추정량의 평균제곱오차는 ML추정량과 비교하여 평균제곱오차간에 큰 차이를 보이지않아 대조를 이루었다. 반면 내포된 그룹의 효과를 나타내는  $\sigma_v^2$ 의 추정에 있어서는 고려된 모든 추정량에서 ML추정량의 평균제곱오차와 비교하여 상대적으로 비슷한 평균제곱오차를 보였다. (지면관계로 제시하지 않은  $(M, N, T)$ 의 다른 조합에 대한 회귀계수와 분산성분에 대한 모의실험 결과들은 저자에게 요구할 수 있음니다.)





## 5. 결론

본 연구에서는 내포오차성분을 가정한 패널회귀모형에서 회귀계수에 대한 여러가지 확장된 FGLS추정량들을 유도하였고, OLS, WTN, ML, REML추정량과 MIVQUE등과 모의실험을 통하여 평균제곱오차 판정기준하에서 효율성을 비교하였다. 모의실험 결과, OLS추정량은 그룹효과와 내포된 그룹효과가 모두 존재하지 않는 경우를 제외한 모든 실험조합에서 효율성이 떨어졌으며, WTN추정량도 FGLS추정량이나 ML, REML추정량 및 MIVQUE에 비하여 효율성이 떨어졌다. 그러나 본 연구에서 제안된 FGLS추정량들은 모두 상대 효율성에서 큰 차이를 보이지 않았다. 특히 계산상으로 훨씬 복잡한 연산과정을 요구하는 ML, REML추정량이나 MIVQUE등과 거의 같은 수준의 효율성을 보여주었다. 그러므로, 내포오차성분을 가지는 패널회귀모형에서는 회귀계수의 추정을 위하여 ML, REML추정이나 MIVQUE보다 비교적 쉽게 얻을 수 있는 FGLS추정량들의 이용을 추천한다.

## 감사의 글

본 논문에 대하여 많은 조언을 해주신 심사위원님들께 감사를 드립니다.

## 참고문헌

- [1] Amemiya, T. (1971). The estimation of variances in a variance components model. *International Economic Review*. Vol. 12, 1-13.
- [2] Balestra, P. (1973). Best quadratic unbiased estimators of the variance-covariance matrix in normal regression. *Journal of Econometrics*. Vol. 2, 17-28.
- [3] Baltagi, B.H. (1995). *Econometric Analysis of Panel Data*. Wiley, New York.
- [4] Burdick, R.K. and Graybill, F.A. (1985). Confidence intervals on the total variance in an unbalanced two-fold nested classification with equal subsampling. *Communications in Statistics -Theory and Methods*. Vol. 14, 761-774.
- [5] Fuller, W.A. and Battese, G.E. (1973). The transformations for estimation of linear models with nested error structure. *Journal of the American Statistical Association*. Vol 68, 626-632.
- [6] Graybill, F.A. and Wang, C.M. (1979). Confidence intervals for proportions of variability in two-factor nested variance component models. *Journal of the American Statistical Association*. Vol 75, 869-873.
- [7] Greene, W.H. (1997). *Econometric Analysis*. Macmillan, New York.

- [8] Harville, D.A. (1977). Maximum likelihood approaches to variance component estimation and to related problems. *Journal of the American Statistical Association*. Vol 72, 320-340.
- [9] Hocking, R.R. (1985). *The Analysis of Linear Model*. Brooks/Cole Company, Monterey, California.
- [10] Hsiao, C. (1986). *Analysis of Panel Data*. Cambridge University Press, Cambridge.
- [11] Moulton, B.R. (1986). Random group effects and the precision of regression estimates. *Journal of Econometrics*. Vol. 32, 385-397.
- [12] Nerlove, M. (1971). Further evidence on the estimation of dynamic economic relations from a time-series of cross-sections. *Econometrica*. Vol. 39, 359-382.
- [13] Pantula, S.G. and Pollock, K.H. (1985). Nested analysis of variance with autocorrelated error. *Biometrics*. Vol 41, 909-920.
- [14] Patterson, H.D. and Thompson, R. (1971). Recovery of inter-block information when block sizes are unequal. *Biometrika*. Vol 58, 545-554.
- [15] Rao, C.R. (1971a). Estimation of variance and covariance components - MINQUE theory. *Journal of Multivariate Analysis*. Vol 1, 257-275.
- [16] Rao, C.R. (1971b). Minimum variance quadratic unbiased estimation of variance components. *Journal of Multivariate Analysis*. Vol 1, 445-456.
- [17] Sahai, H. (1976). A comparison of estimators of variance components in the balanced three-stage nested random effects model using mean squared error criterion. *Journal of the American Statistical Association*. Vol 71, 435-444.
- [18] Swallow, W.H. and Monahan, J.F. (1984). Monte Carlo comparison of ANOVA, MIVQUE, REML, and ML estimators of variance components. *Technometrics*. Vol 26, 47-57.
- [19] Swallow, W.H. and Searle, S.R. (1978). Minimum variance quadratic unbiased estimation(MIVQUE) of variance components. *Technometrics*. Vol 20, 265-272.
- [20] Swamy, P.A.V.B. and Arora, A.A. (1972). The exact finite sample properties of the estimators of coefficients in the error components regression models. *Econometrica*. Vol. 40, 261-275.
- [21] Wallace, T.D. and Hussain, A. (1969). The use of error components models in combining cross-section and time-series data. *Econometrica*. Vol 37, 55-72.

- [22] Wang, C.M. and Graybill, F.A. (1981). Confidence intervals on a ratio of variances in the two-factor nested components of variance model. *Communications in Statistics -Theory and Methods*. Vol. 10, 1357-1368.
- [23] Wansbeek, T.J. and Kapteyn, A. (1982). A simple way to obtain the spectral decomposition of variance components models for balanced data. *Communications in Statistics -Theory and Methods*. Vol. 11, 2105-2112.
- [24] Wansbeek, T.J. and Kapteyn, A. (1983). A note on spectral decomposition and maximum likelihood estimation in ANOVA models with balanced data. *Statistics and Probability Letters*. Vol. 1, 213-215.

[ 1998년 11월 접수, 1999년 5월 채택 ]

## A Comparison of Estimation Procedures in a Nested Error Components Regression Model

Seuck-Heun Song<sup>1)</sup> Myoungshic Jhun<sup>2)</sup> Byoung-Cheol Jung<sup>3)</sup>

### ABSTRACT

This paper considers a linear regression model with nested error components and investigates the performances of various estimation procedures for the regression coefficient. Through simulation study some of the basic finding is following: For the regression coefficients, the computationally simple FGLS estimation methods perform reasonably well when compared with the computationally involved MLE, RMLE, MIVQUE methods.

*Keywords:* Panel Regression Model; FGLS; ML; REML; MIVQUE.

---

1) Assistant Professor, Dept. of Statistics, Korea University. E-mail: ssong@kucncx.korea.ac.kr

2) Professor, Dept. of Statistics, Korea University. E-mail: jhun@kucncx.korea.ac.kr

3) Postdoctoral Researcher, Institute of Statistics, Korea University. E-mail: bcjung@kustat.korea.ac.kr