

Confidence Intervals on Variance Components
in Multiple Regression Model with
One-fold Nested Error Structure

박 동 준
Dong Joon Park

부산수산대학교 응용수학과

Abstract

Regression model with nested error structure interval estimations about variability on different stages are proposed. This article derives an approximate confidence interval on the variance in the first stage and an exact confidence interval on the variance in the second stage in two stage regression model. The approximate confidence interval is based on Ting et al.(1990) method. Computer simulation is provided to show that the approximate confidence interval maintains the stated confidence coefficient.

1. Introduction

This article considers the multiple regression model where the responses are correlated. In particular, we consider the multiple regression model with one-fold nested error structure, i.e., two stage regression model. This model includes two error terms. One is associated with the first-stage sampling unit and the other with the second-stage sampling unit. These two error terms are independent and normally distributed with zero means and constant variances. However, this error structure gives correlated response variables. Aitken and Longford(1986) showed ignoring the nesting structure is not appropriate to estimate regression coefficients.

Park and Burdick(1993) proposed the confidence intervals on the variance components in simple regression model with one-fold nested error structure. Tsubaki et al.(1995) proposed methods to estimate regression coefficients. This

article derives the confidence intervals on variance components associated with primary and secondary sampling units in two stage regression model.

2. Multiple regression model with one-fold nested error structure

Multiple regression model with one-fold nested error structure is written as

$$Y_{ij} = \beta_0 + \beta_1 X_{ij1} + \dots + \beta_{p_1} X_{ijp_1} + \delta_i + \gamma_1 X_{ij1} + \dots + \gamma_{p_2} X_{ijp_2} + \varepsilon_{ij} \quad (2.1)$$

$$i_1 = 1, \dots, \lambda_1; \dots; i_{p_1} = 1, \dots, \lambda_{p_1}$$

$$i = 1, \dots, l_1; j = 1, \dots, l_2$$

where Y_{ij} is the j th observation in the i th cell(group), β_0 is an intercept term, $\beta_1, \dots, \beta_{p_1}$ are unknown parameters associated with primary units, $X_{i1}, \dots, X_{i, p_1}$ are fixed predictor variables in the primary unit, $\gamma_1, \dots, \gamma_{p_2}$ are unknown parameters associated with secondary units, $X_{ij1}, \dots, X_{ijp_2}$ are fixed predictor variables in the secondary unit, δ_i is a random error term in the primary unit, ε_{ij} is a random

error term in the secondary unit, δ_i and ε_{ij} are jointly independent normal random variables with zero means and variances σ_δ^2 and σ_ε^2 , respectively. The index l_1 is the number of different combinations (cells) of levels among X_{ij} 's, i.e., $l_1 = \lambda_1 \times \lambda_2 \times \dots \times \lambda_{p_1}$ and l_2 is the number of repetitions within an i th cell. We consider the balanced case where l_2 's are same for all i 's. Since β 's, γ 's, X_{ij} 's, and X_{ijk} 's are fixed, and δ_i and ε_{ij} are random, model (2.1) is a mixed model. The model (2.1) is written in matrix notation,

$$\underline{Y} = ZX_1\beta + X_2\gamma + Z\delta + \varepsilon \quad (2.2.1)$$

$$= Z\underline{U} + X_2\gamma + \varepsilon \quad (2.2.2)$$

$$= X\underline{\alpha} + \xi \quad (2.2.3)$$

where

$$\underline{U} = X_1\beta + \delta, \quad X = (ZX_1 \ X_2),$$

$$\underline{\alpha} = \begin{pmatrix} \beta \\ \gamma \end{pmatrix}, \quad \text{and} \quad \xi = Z\delta + \varepsilon$$

where \underline{Y} is an $l_1 l_2 \times 1$ vector of observations, Z is an $l_1 l_2 \times l_1$ design matrix with 0's and 1's, i.e., $Z = \bigoplus_{i=1}^{l_1} \mathbf{1}_{l_2}$ where $\mathbf{1}_{l_2}$ is an $l_2 \times 1$ column vector of 1's and \bigoplus is the direct sum operator, X_1 is an $l_1 \times (p_1 + 1)$ matrix of known values with a column of 1's in the first column and p_1 columns of X_{ij} 's from the second column to the p_1 th column, β is a $(p_1 + 1) \times 1$ vector of parameters associated with X_{ij} 's, X_2 is an $l_1 l_2 \times p_2$ matrix of known values with p_2 columns of X_{ijk} 's from the first column to the p_2 th column, γ is a $p_2 \times 1$ vector of parameters associated with X_{ijk} 's, δ is an $l_1 \times 1$ vector of random error terms, and ε is an $l_1 l_2 \times 1$ vector of random error terms. In particular, From (2.2.3), the variance-covariance matrix of \underline{Y} is

$$\text{Var}(\underline{Y}) = \sigma_\delta^2 ZZ' + \sigma_\varepsilon^2 I_{l_1 l_2}, \quad (2.3)$$

since $\delta \sim N(\mathbf{0}, \sigma_\delta^2 I_{l_1})$ and

$\varepsilon \sim N(\mathbf{0}, \sigma_\varepsilon^2 I_{l_1 l_2})$ where I_{l_1} is an $l_1 \times l_1$ identity matrix. From the assumptions in (2.1) and equation (2.3),

$$\underline{Y} \sim N(X\underline{\alpha}, \sigma_\delta^2 ZZ' + \sigma_\varepsilon^2 I_{l_1 l_2}). \quad (2.4)$$

The regression sums of squares of model (2.1) are now investigated. The reductions in sums of squares of the model are attributable to fitting the primary and secondary fixed variables and are expressed into the quadratic forms. Let $G_1 = (X^* X^*)^{-1}$ and $G_2 = (\overline{X}_2' \overline{X}_2)^{-1}$

$$\text{where } X^* = (X_1 \ X_2), \ X_2^* = \frac{Z'}{l_2} X_2,$$

$$\overline{X}_2 = WX_2, \text{ and } W = I_{l_1 l_2} - ZZ'/l_2. \text{ Define}$$

$H_1 = X^* G_1 X^*$ and $H_2 = \overline{X}_2' G_2 \overline{X}_2$. Now consider the quadratic forms

$$R_1 = \underline{Y}' \frac{Z'}{l_2} (I_{l_1} - H_1) \frac{Z'}{l_2} \underline{Y} \quad \text{and}$$

$$R_2 = \underline{Y}' W' (I_{l_1 l_2} - H_2) W \underline{Y}. \text{ The quadratic}$$

form R_1 is determined by computing the regression of $\overline{Y}_{i.}$ on X_{ij} and $\overline{X}_{i.k}$ where

$$\overline{Y}_{i.} = \sum_{j=1}^{l_2} Y_{ij} / l_2 \text{ and } \overline{X}_{i.k} = \sum_{j=1}^{l_2} X_{ijk} / l_2.$$

The quadratic form R_2 is calculated by the regression of Y_{ij} on the secondary fixed variables, X_{ijk} , and grouping variables. Under the distributional assumptions in (2.1), the quadratic forms $R_1 / (\sigma_\delta^2 + \frac{\sigma_\varepsilon^2}{l_2})$ and $R_2 / \sigma_\varepsilon^2$

are chi-squared random variables with $l_1 - p_1 - p_2 - 1$ and $l_1 l_2 - l_1 - p_2$ degrees of freedom, respectively. In addition, the quadratic forms $R_1 / (\sigma_\delta^2 + \frac{\sigma_\varepsilon^2}{l_2})$ and $R_2 / \sigma_\varepsilon^2$ are

independent (see Park(1996)). That is,

$$\frac{R_1}{\sigma_\delta^2 + \frac{\sigma_\varepsilon^2}{l_2}} \sim \chi_{l_1 - p_1 - p_2 - 1}^2 \quad (2.5)$$

$$\text{and} \quad \frac{R_2}{\sigma_\varepsilon^2} \sim \chi_{l_1 l_2 - l_1 - p_2}^2. \quad (2.6)$$

3. Confidence Intervals on σ_δ^2 and σ_ε^2

Define $S_\delta^2 = R_1 / n_1$ and $S_\varepsilon^2 = R_2 / n_2$,

where $n_1 = l_1 - p_1 - p_2 - 1$ and

$n_2 = l_1 l_2 - l_1 - p_2$. Using (2.5) and (2.6), the expected mean squares are

$$E(S_\delta^2) = \sigma_\delta^2 + \frac{\sigma_\varepsilon^2}{l_2} = \theta_\delta \quad (3.1)$$

$$E(S_\varepsilon^2) = \sigma_\varepsilon^2 = \theta_\varepsilon. \quad (3.2)$$

Since $R_2/\sigma_\varepsilon^2 \sim X_{n_2}^2$, an exact confidence interval on σ_ε^2 exists. This exact $1-2\alpha$ two-sided confidence interval on σ_ε^2 is

$$\left[\frac{S_\varepsilon^2}{F_{\alpha:n_2,\infty}} ; \frac{S_\varepsilon^2}{F_{1-\alpha:n_2,\infty}} \right] \quad (3.3)$$

where $F_{\delta:v_1,v_2}$ is the $1-\delta$ th percentile F -value with v_1 and v_2 degrees of freedom.

The variance component σ_δ^2 is represented by the mean squares in (3.1) and (3.2). From (3.1) and (3.2),

$$\sigma_\delta^2 = \theta_\delta - \frac{\theta_\varepsilon}{l_2} \quad (3.4)$$

Confidence intervals on σ_δ^2 can be constructed using the method of Ting et al.(1990). The $1-2\alpha$ two-sided confidence interval on σ_δ^2 using (3.4) is

$$\left[S_\delta^2 - \frac{S_\varepsilon^2}{l_2} - (U_1^2 S_\delta^4 + U_2^2 \frac{S_\varepsilon^4}{l_2^2} + U_{12} S_\delta^2 \frac{S_\varepsilon^2}{l_2})^{\frac{1}{2}} ; S_\delta^2 - \frac{S_\varepsilon^2}{l_2} + (V_1^2 S_\delta^4 + V_2^2 \frac{S_\varepsilon^4}{l_2^2} + V_{12} S_\delta^2 \frac{S_\varepsilon^2}{l_2})^{\frac{1}{2}} \right] \quad (3.5)$$

where $U_1 = 1 - 1/F_{\alpha:n_1,\infty}$,

$$U_2 = 1/F_{1-\alpha:n_2,\infty} - 1,$$

$$U_{12} = [(F_{\alpha:n_1,n_2} - 1)^2 - U_1^2 F_{\alpha:n_1,n_2}^2 - U_2^2]$$

$$/ F_{\alpha:n_1,n_2}, V_1 = 1/F_{1-\alpha:n_1,\infty} - 1,$$

$$V_2 = 1 - 1/F_{\alpha:n_2,\infty}, V_{12} = [(1 - F_{1-\alpha:n_1,n_2})^2 - V_1^2 F_{1-\alpha:n_1,n_2}^2 - V_2^2] / F_{1-\alpha:n_1,n_2}.$$
 Since

$\sigma_\delta^2 > 0$, any negative bound is defined to be zero.

4. Simulation Study

Computer simulation was performed to compare the stated confidence coefficient and expected interval lengths. Let $\rho = \sigma_\delta^2 / (\sigma_\delta^2 + \sigma_\varepsilon^2)$. Without loss of generality

$\sigma_\delta^2 = 1 - \sigma_\varepsilon^2$ so that $\rho = \sigma_\delta^2$ and $1 - \rho = \sigma_\varepsilon^2$.

Therefore, $S_\delta^2 \sim ((\rho + \frac{1-\rho}{l_2})/n_1)\chi_{n_1}^2$ and

$S_\varepsilon^2 \sim ((1-\rho)/n_2)\chi_{n_2}^2$. These independent

scaled chi-squared random variables can be generated by the RANGAM routine of the Statistical Analysis System(SAS). Values of ρ are varied from 0 to 1 in increments of 0.1 and simulated 1000 times for each design. Simulated values of S_δ^2 and S_ε^2 are substituted into (3.5) and the intervals are computed.

The average lengths of the two-sided confidence intervals are also calculated. Table 1 reports the results of the simulation for stated 90% confidence intervals and the range of two-sided interval lengths on σ_δ^2 using (3.5) when $p_1 = 2$ and $p_2 = 3$. The proposed interval generally keep the stated confidence coefficients since all simulated confidence coefficients are bigger than 0.8814 and are not too conservative. The interval lengths get smaller as l_1 becomes bigger since it increases n_1 degrees of freedom. In addition, the interval lengths get smaller as l_2 becomes bigger since it increases n_2 degrees of freedom.

TABLE 1. Simulated Confidence Coefficients and Average Interval Lengths for 90% Two-sided Intervals on σ_δ^2

l_1	l_2		Coefficient	Length
8	5	Max	0.916	19.3902
		Min	0.890	3.7052
8	10	Max	0.911	18.7950
		Min	0.895	1.9164
14	5	Max	0.920	2.3854
		Min	0.892	0.3743
14	10	Max	0.910	2.4561
		Min	0.891	0.1848
20	5	Max	0.926	1.5090
		Min	0.896	0.2310
20	10	Max	0.915	1.5130
		Min	0.888	0.1119

5. Conclusions

This paper utilized distributional property of variance components in multiple regression model

with one-fold nested error structure and derived confidence intervals on the variance components by use of independent quadratic forms which are chi-squared distributed. An exact confidence interval on the variability in the second stage of the model was obtained in (3.3) and an approximate confidence interval on the variability in the first stage of the model was proposed in (3.5). The simulations were performed to show that the proposed approximate confidence interval kept the stated confidence coefficients and average interval lengths changed as degrees of freedom of chi-squared random variables increased. The proposed confidence interval is recommended in multiple regression model with one-fold nested error structure.

References

- [1] Aitken, M. and Longford, N. T. "Statistical modelling issues in school effectiveness studies(with discussion)", *Journal of the Royal Statistical Society A.*, 14A, 1-43, 1986.
- [2] Park, D. J. and Burdick, R. K. "Confidence intervals on the among group variance component in a simple linear regression model with a balanced one-fold nested error structure", *Communications in Statistics-theory and methods*, 22(12), 3435-3452, 1993.
- [3] Park, D. J. "The distributions of variance components in two stage regression model", *통계이론방법연구* (forthcoming).
- [4] Searle, S. R. *Linear Models*, New York: John Wiley & Sons inc., 1971
- [5] Statistical Analysis System SAS *User's Guide: Statistics*, Cary, North Carolina: SAS Institute Inc, 1994.
- [6] Ting, N., Burdick, R. K., Graybill, F. A., Jeyaratnam, S., and Lu, T.-F.C. "Confidence intervals on linear combinations of variance components", *Journal of Statistical Computation and Simulation*, 35, 135-143, 1990.
- [7] Tsubaki, M., Tsubaki, H, and Kusumi, M. A new estimation method for a useful class of mixed models, *Proceedings of the 9th Asia Quality Management Symposium - Quality Enhancement for Global Prosperity*, 275-280, 1995.