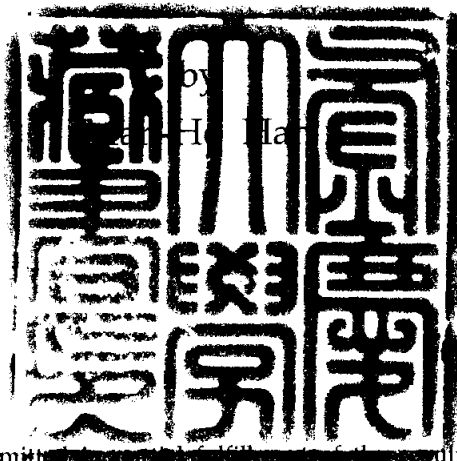# Comparison of Confidence Intervals on Variance Component in a Simple Linear Regression Model with Unbalanced Nested Error Structure

## 불균형 중첩 오차구조를 갖는 단순선형 회귀모형의 분산의 신뢰구간의 비교

Advisor : Dong Joon Park

A thesis submitted in partial fulfillment of the requirements
for the degree of

Master of Education

in the Department of Mathematics Education Graduate School of Education
Pukyong National University

August 2003

# 韓萬浩의 教育學碩士 學位論文을 認准함

## 2003年 6月 日

主　　審　理學博士　　張　大　興　㊞

委　　員　理學博士　　朴　瑢　範　㊞

委　　員　理學博士　　朴　東　俊　㊞

# TABLE OF CONTENTS

# LIST OF TABLES

# 불균형 중첩 오차구조를 갖는 단순선형 회귀모형의 분산의 신뢰구간의 비교

한 만 호


부경대학교 교육대학원 수학교육학과

요        약

불균형 중첩 오차구조를 갖는 단순선형 회귀모형에 나타나는 주 샘플링 단위의 분산 $\sigma_A^2$ 에 대한 세 가지 신뢰구간을 구하였다. 그 중에서 두 가지 신뢰구간은 ANOVA 방법으로부터 구한 기대평균제곱을 $\sigma_A^2$에 관한 식을 구하여 Ting et al.(1990) 방법을 적용하여 구하였다. 나머지 한 방법은 Khuri et al.(1998)이 제안한 일반화 p-값을 활용하여 $\sigma_A^2$애 관한 신뢰구간을 구하였다. 세 가지 신뢰구간을 비교하기 위하여 시뮬레이션을 시행하고 실제 자료에 적용하여 시뮬레이션의 결과와 일관성 있는 결과를 보이는 것을 확인하였다.

# 1. INTRODUCTION

## 1.1 Experimental Design

Statisticians often use experimental design to compare the differences of treatments in their experimental research. Experimental design is to design and conduct experiments to obtain best information at a minimum cost. The relationships of the effects of different levels in an experimental design are expressed as a linear model. The linear model of an experimental design is divided into three kinds: fixed effects model, random effects model, and mixed effects model. The fixed effects model with $I$ populations and a random sample of size $J$ from each population is written as

$$Y_{ij} = \mu + \alpha_i + E_{ij} \tag{1.1}$$

$$i = 1, ..., I; \quad j = 1, ..., J$$

$$\sum_{i=1}^{I} \alpha_i = 0$$

where $Y_{ij}$ is the $j$th observed sample value from the $i$th population, the quantities $\mu$ and $\alpha_i$ are unobservable fixed constants called parameters, $E_{ij}$ is a random error term with mean zero and variance $\sigma_E^2$. The objective of a fixed effects model is to make inferences about all treatment levels included in the experiment.

Suppose that four specific types of training methods are used by a company to train operators to fill bottles with vegetable oil. The purpose of the experiment is to compare four training methods. Model (1.1) can be used for this experiment where $Y_{ij}$ is the weight of the $j$th bottle filled with oil by the $i$th method, $\mu$ is a constant representing the average bottle weight filled with oil,

$\alpha_i$ is the $i$th training method effect, and $E_{ij}$ is an independent normal random variable with mean zero and variance $\sigma_E^2$. Specifically, $I = 4$ represents the number of training methods, $J$ is the number of bottles filled with oil by each training method, and $\sigma_E^2$ is a measure of variability of bottle weight filled with oil for a particular method. The investigator is interested in making inferences about four training methods in this experiment.

Consider a simple experiment in which treatment levels for a factor are randomly selected from a large population. A random effects model is written as

$$Y_{ij} = \mu + A_i + E_{ij} \tag{1.2}$$

$$i = 1, ..., I; \ \ j = 1, ..., J$$

where $A_i$ and $E_{ij}$ are random variables with means of zero and variances $\sigma_A^2$ and $\sigma_E^2$, respectively. In this model we are interested in the variability of the $A_i$ as measured by $\sigma_A^2$. The objective of a random effects model is inferences concerning functions of variances. The random effects model is also referred to as a variance component model.

As an example, consider machines in a large plant that are used to fill bottles with vegetable oil. Five machines are selected at random from which a sample of bottles are filled and weighted. The purpose of the experiment is to determine how much the weight variability in the bottles is attributed to variability among machines. The experimental model is represented as (1.2.) where $Y_{ij}$ represents the weight of the $j$th bottle filled with oil by the $i$th machine, $\mu$ is a constant representing the average bottle weight filled with oil, and $A_i$ and $E_{ij}$ are mutually independent normal random variables with means of zero and variances $\sigma_A^2$ and $\sigma_E^2$, respectively. In the context example, $I = 5$ represents

the number of sampled machines, and $J$ is the number of bottles filled with oil by each machine. The variance component $\sigma_A^2$ is a measure of variability of bottle weight filled with oil across machines and $\sigma_E^2$ is a measure of weight variability for any particular machine. The investigator is primarily interested in determining the amount of variability among machines in the population.

Suppose that a second factor with specific treatment levels of interest is add to model (1.2). If each level of factor A is crossed with each level of factor B and each combination is replicated $k$ times, the model is expressed as

$$Y_{ijk} = \mu + A_i + \beta_j + E_{ijk} \tag{1.3}$$

$$i = 1, ..., I; \quad j = 1, ..., J; \quad k = 1, 2, ... K$$

$$\sum_{j=1}^{J} \beta_j = 0$$

where $Y_{ijk}$ is the $k$th observed sample value in the $i$th level of factor $A$ and the $j$th level of factor $B$, $A_i$ is independently distributed as $N(0, \sigma_A^2)$, $\beta_j$ is the effect of the $j$th level of factor $B$, $E_{ijk}$ is independently distributed as $N(0, \sigma_E^2)$ and $A_i$ and $E_{ijk}$ are independent. Factor $B$ in (1.3) is a fixed effect where the selected treatment levels in the experiment are of interest. That is, an investigator is interested in estimating functions of $\beta_j$. Model (1.3) includes both a random effect (factor $A$) and a fixed effect (factor $B$). This type of model is called a mixed effects model. The objective of a mixed effects model is to make inferences concerning functions of variance $\sigma_A^2$ and inferences of $\beta_j$.

Suppose that two types of training courses are used by the company to train operators to use the filling machines. After five machines are randomly selected from the population of machines, three bottles are filled by operators using each training course. The problem of interest is to not only determine weight

variability among the machines but to also determine the difference between the two training courses. Model (1.3) can be used for this experiment where $Y_{ijk}$ is the weight of the $k$th bottle filled on the $i$th machine by the $j$th training course, $\mu$ is a constant representing the average bottle weight filled with oil, $A_i$ and $E_{ij}$ are mutually independent normal random variables with zero means and variances $\sigma_A^2$ and $\sigma_E^2$, respectively and $\beta_j$ is the $j$th training course effect. Additionally, $I = 5$ represents the number of machines for each training course, $J = 2$ represents the number of training courses, and $K = 3$ represents the number of bottles filled with oil by each combination of machine and training course.

## 1.2 Statistical Inferences

Statistical inferences are largely divided into estimation and test of hypotheses. Estimation includes point estimation and interval estimation. The selection of a function of the sample values that will best represent the parameter of interest is concerned with point estimation.

Interval estimation is generally more informative than point estimation because it is not enough to obtain a single value for the parameter under investigation and a point estimate has no information about confidence and bound of error. An interval estimation provides this information. Let $\theta$ represent a parameter of interest. A confidence interval is a random interval whose endpoints L and U, where $L \leq U$ are functions of the sample values such that $P[L \leq \theta \leq U] = 1 - \alpha$. The term $1 - \alpha$ is the confidence coefficient and is selected prior to data collection. A confidence interval [L, U] that satisfies $P[L \leq \theta \leq U] = 1 - \alpha$ is called an exact two-sided $1 - \alpha$ confidence interval. Often exact $1 - \alpha$ confidence intervals do not exist and $P[L \leq \theta \leq U]$ is only approximately equal to $1 - \alpha$. These intervals are referred to as approximate intervals. An approximate interval is conservative if $P[L \leq \theta \leq U] > 1 - \alpha$ and liberal if $P[L \leq \theta \leq U] < 1 - \alpha$.

Hypothesis testing refers to the process of trying to decide the truth or falsity of hypotheses on the basis of experimental evidence. Confidence intervals and tests of hypotheses are procedures for making statistical inferences that attach measures of uncertainty to the inferences. It is almost always the case, however, that confidence intervals are "uniformly more informative" than tests of hypotheses for making decisions based on parametric values. Thus, tests of hypotheses are seldom needed if confidence intervals are available.

## 1.3 Literature Review

If the number of observations in cells is not equal, then experimental designs are unbalanced. The unbalanced one-fold nested design model is written as

$$Y_{ij} = \mu + A_i + E_{ij} \qquad (1.4)$$

$$i = 1, ..., I; \quad j = 1, ..., J_i$$

where $\mu$ is an unknown constant, $A_i$ and $E_{ij}$ are mutually independent normal random variables with means of zero and variances $\sigma_A^2$ and $\sigma_E^2$, respectively, $I \geq 2$, $J_i \geq 1$, and $J_i > 1$ for at least one value of $i$. The analysis of variance table for the one-fold nested design is given in Table 1.1.

### TABLE 1.1 ANOVA for One-fold Nested Design

| SV | DF | MS | EMS |
|---|---|---|---|
| Among Groups | $n_1 = I - 1$ | $S_1^2$ | $\theta_1 = \sigma_E^2 + c_1 \sigma_A^2$ |
| Within Groups | $n_2 = N - I$ | $S_2^2$ | $\theta_2 = \sigma_E^2$ |
| Total | $N - 1$ | | |

$$N = \sum_{i=1}^{I} J_i, \qquad (1.5)$$

$$S_1^2 = \sum_{i=1}^{I} \sum_{j=1}^{J_i} (\bar{Y}_{i.} - \bar{Y}_{..})^2,$$

$$S_2^2 = \sum_{i=1}^{I} \sum_{j=1}^{J_i} (Y_{ij} - \bar{Y}_{i.})^2, \qquad \text{and}$$

$$c_1 = \frac{N - \sum_{i=1}^{I} J_i^2 / N}{I - 1}.$$

In the balanced design where all $J_i = J$, $N = IJ$, and $c_1 = J$, $n_1 S_1^2/\theta_1$ and $n_2 S_2^2/\theta_2$ are independent chi-squared random variables with $n_1$ and $n_2$ degrees of freedom, respectively. In the unbalanced design, $S_1^2$ and $S_2^2$ are still independent and $n_2 S_2^2/\theta_2$ has a chi-squared distribution with $n_2$ degrees of freedom. However, unless $\sigma_A^2 = 0$, $n_1 S_1^2/\theta_1$ no longer has a chi-squared distribution.

Thomas and Hultquist (1978) proposed a statistic that can be used for constructing a confidence interval on $\sigma_A^2$ in the unbalanced model. The proposed statistic is

$$\frac{n_1 S_{1U}^2}{\theta_{1U}} \tag{1.6}$$

where

$$n_1 S_{1U}^2 = \sum_{i=1}^{I} \bar{Y}_{i.}^2 - \left(\frac{1}{I}\right)\left(\sum_{i=1}^{I} \bar{Y}_{i.}\right)^2,$$

$$\bar{Y}_{i.} = \sum_{j=1}^{J_i} \frac{Y_{ij}}{J_i},$$

$$\theta_{1U} = E(S_{1U}^2) = \sigma_A^2 + \left(\frac{1}{h}\right)\sigma_E^2, \qquad \text{and}$$

$$h = \frac{I}{\sum_{i=1}^{I} \frac{1}{J_i}}.$$

The term $n_1 S_{1U}^2$ means the unweighted sums of squares of the treatment means and $h$ represents the harmonic mean of the $J_i$ values. They showed that the moment generation function on $n_1 S_{1U}^2/\theta_{1U}$ approaches that of a chi-squared

random variable with $n_1$ degrees of freedom as all $J_i$ approach a constant or if either $\lambda_A = \sigma_A^2/\sigma_E^2$ or all $J_i$ approach infinity. They showed that $n_1 S_1^2/\theta_{1U}$ is well approximated by a chi-squared random variables when $\lambda_A > 0.25$. In situations where the Thomas-Hultquist approximation works well, an interval on $\sigma_A^2$ can be formed by replacing $S_1^2$ with $hS_{1U}^2$ and $J$ with $h$ in the balanced design equations.

In extremely unbalanced designs where $\lambda_A < 0.25$, the chi-squared approximation for $n_1 S_{1U}^2/\theta_{1U}$ is not good and this substitution can yield a liberal confidence interval. A method that works well over the entire range of $\lambda_A$ was developed by Burdick and Eickman (1986). The Burdick-Eickman approximate $100(1-\alpha)\%$ confidence interval is

$$
[\frac{hS_{1U}^2 L^*}{F_{\alpha_{11}:n_1,\infty}(1+hL^*)} ; \frac{hS_{1U}^2 U^*}{F_{1-\alpha_{21}:n_1,\infty}(1+hU^*)}] \tag{1.7}
$$

where

$$
L^* = \frac{S_{1U}^2}{(F_{\alpha_{12}:n_1,n_2} S_2^2)} - \frac{1}{m},
$$

$$
U^* = \frac{S_{1U}^2}{(F_{1-\alpha_{22}:n_1,n_2} S_2^2)} - \frac{1}{M},
$$

$$
m = \min(J_1, J_2, ..., J_I),
$$

$$
M = \max(J_1, J_2, ..., J_I), \qquad \text{and}
$$

$$
\alpha_{11} + \alpha_{21} = \alpha_{12} + \alpha_{22} = \alpha.
$$

Burdick and Eickman conducted a simulation study to show (1.7) is generally conservative. They showed that their method can always be recommended over the Thomas-Hultquist approximation. The average interval lengths of these two

methods never differed by more than 5% and the Burdick and Eickman method maintains its confidence coefficient over a wider range of unbalanced designs than does the Thomas-Hultquist method.

The variance component model with one explanatory variable is

$$Y_{ij} = \mu + \beta X_{ij} + A_i + E_{ij} \tag{1.8}$$

$$i = 1, ..., I; \ \ j = 1, ..., J_i$$

where $A_i$ is the cluster effect and assumed to be a random sample from $N(0, \sigma_A^2)$, and $E_{ij}$ is an observational error within a cluster and assumed to be a random sample from $N(0, \sigma_E^2)$. The random variables $A_i$ and $E_{ij}$ are independent. This model is also referred to as a simple regression model with nested error structure. In the balanced case where $J_i = J$, $\hat{\beta}$ is the ordinary least squares estimator of $\beta$. Several methods for point estimation of the regression coefficients have been proposed for model (1.8) and its various extensions.

Researches have also been done in confidence intervals and tests of hypothesis in the variance component model with one explanatory variable in the balanced case where $J_i = J$. In this case the model is also called simple regression model with balanced nested error structure. Tong and Cornelius (1989) compared four estimators of regression coefficient $\beta$ in the model with respect to their mean squared error in a Monte Carlo simulation study. Tong and Cornelius (1991) investigated properties of tests of hypothesis for regression coefficient $\beta$ in the model and compared with respect to type I error rate and power of test in a Monte Carlo simulation study. Guven (1995) derived explicit maximum likelihood estimators of regression coefficient $\beta$ in the model.

Park and Burdick (1993) derived three approximate confidence intervals on $\sigma_A^2$ using distributional results for sums of squares associated with the model.

Park and Burdick (1994) proposed several confidence intervals on the regression coefficient $\beta$ in the model and the intervals were compared using computer simulation. Park and Hwang (2002) derived exact and approximate confidence intervals for the mean response for a given level of the independent variable in the simple linear regression model with nested error structure. Yu and Burdick (1995) extended the model and considered confidence intervals on the variance components in regression models with balanced (Q-1)-fold nested error structure. They used a method proposed by Ting, Burdick, Graybill, Jeyaratnam, and Lu (1990). That is, the regression model with two-fold nested error structure was first considered and then results were generalized to the (Q-1)-fold nested error structure.

# 2. A SIMPLE REGRESSION MODEL WITH AN UNBALANCED ONE-FOLD NESTED ERROR STRUCTURE

The simple regression model with an unbalanced one-fold nested error structure is written as

$$Y_{ij} = \mu + \beta X_{ij} + A_i + E_{ij} \tag{2.1}$$

$$i = 1, ..., I; \quad j = 1, ..., J_i$$

where $Y_{ij}$ is the $j$th observation in the $i$th primary level, $\mu$ and $\beta$ are unknown constants, $X_{ij}$ is a fixed predictor variable, and $A_i$ and $E_{ij}$ are jointly independent normal random variables with zero means and variances $\sigma_A^2$ and $\sigma_E^2$, respectively, $I \geq 2$, $J_i \geq 1$, and $J_i > 1$ for at least one value of $i$. $A_i$ is an error term associated with the first-stage sampling unit and $E_{ij}$ is an error term associated with the second-stage sampling unit. Model (2.1) is unbalanced since the number of observations in cells are not all equal. This model is referred to as either a single-factor covariance model with one covariate or a variance component model with one explanatory variable. Since the $X_{ij}$ and $\beta$ are fixed, model (2.1) is a mixed model. This error structure yields response variables that are correlated. That is,

$$Cov(Y_{ij}, Y_{i'j'}) = \begin{cases} \sigma_A^2 + \sigma_E^2 & \text{if} \quad i = i', j = j'; \\ \sigma_A^2 & \text{if} \quad i = i', j \neq j'; \\ 0 & \text{if} \quad i \neq i'. \end{cases} \tag{2.2}$$

In order to form confidence intervals on linear functions of the variance components, an appropriate set of sums of squares is needed. One possible partitioning of model (2.1) is shown in Table 2.1.

## TABLE 2.1 ANOVA for Model (2.1)

| SV | DF | SS |
|---|---|---|
| Mean | 1 | $J.\bar{Y}_{..}^2$ |
| Covariate after mean | 1 | $\hat{\beta}_L^2(S_{wxxa} + S_{wxxe})$ |
| Primary units adjusted for regression | $I - 1$ | $R_{WB} + R_L$ |
| Residual | $J. - I - 1$ | $R_T$ |
| Total | $J.$ | $\sum\limits_{i=1}^{I} \sum\limits_{j=1}^{J_i} Y_{ij}^2$ |

The notation in Table 2.1 is defined as

$$J. = \sum_{i=1}^{I} J_i,$$

$$\bar{X}_{i.} = \frac{\sum\limits_{j=1}^{J_i} X_{ij}}{J_i},$$

$$\bar{Y}_{i.} = \frac{\sum\limits_{j=1}^{J_i} Y_{ij}}{J_i},$$

$$\bar{X}_{..} = \frac{\sum\limits_{i=1}^{I} \sum\limits_{j=1}^{J_i} X_{ij}}{J.} = \frac{\sum\limits_{i=1}^{I} \bar{X}_{i.} J_i}{J.},$$

$$\bar{Y}_{..} = \frac{\sum\limits_{i=1}^{I} \sum\limits_{j=1}^{J_i} Y_{ij}}{J_{.}} = \frac{\sum\limits_{i=1}^{I} \bar{Y}_{i.} J_i}{J_{.}},$$

$$S_{wxxa} = \sum\limits_{i=1}^{I} (\bar{X}_{i.} - \bar{X}_{..})^2 J_i,$$

$$S_{wxxe} = \sum\limits_{i=1}^{I} \sum\limits_{j=1}^{J_i} (X_{ij} - \bar{X}_{i.})^2,$$

$$S_{wxya} = \sum\limits_{i=1}^{I} (\bar{X}_{i.} - \bar{X}_{..})(\bar{Y}_{i.} - \bar{Y}_{..}) J_i,$$

$$S_{wxye} = \sum\limits_{i=1}^{I} \sum\limits_{j=1}^{J_i} (X_{ij} - \bar{X}_{i.})(Y_{ij} - \bar{Y}_{i.}),$$

$$S_{wyya} = \sum\limits_{i=1}^{I} (\bar{Y}_{i.} - \bar{Y}_{..})^2 J_i,$$

$$S_{wyye} = \sum\limits_{i=1}^{I} \sum\limits_{j=1}^{J_i} (Y_{ij} - \bar{Y}_{i.})^2,$$

$$\hat{\beta}_{WB} = \frac{S_{wxya}}{S_{wxxa}},$$

$$\hat{\beta}_{L} = \frac{(S_{wxya} + S_{wxye})}{(S_{wxxa} + S_{wxxe})},$$

$$\hat{\beta}_{T} = \frac{S_{wxye}}{S_{wxxe}},$$

$$R_{WB} = S_{wyya} - \hat{\beta}_{WB}^2 S_{wxxa},$$

$$R_{L} = \hat{\beta}_{WB}^2 S_{wxxa} + \hat{\beta}_{T}^2 S_{wxxe}$$

$$- \hat{\beta}_{L}^2 (S_{wxxa} + S_{wxxe}), \quad \text{and}$$

$$R_{T} = S_{wyye} - \hat{\beta}_{T}^2 S_{wxxe}.$$

Model (2.1) is written in matrix notation,

$$\underline{Y} = \mathbf{X}\underline{\alpha} + \mathbf{Z}\underline{U} + \underline{E} \tag{2.3}$$

where $\underline{Y}$ is a $J \times 1$ vector of observations, $\mathbf{X}$ is a $J \times 2$ matrix of known values with a column of 1's in the first column and a column of $X_{ij}$'s in the second column, $\underline{\alpha}$ is a $2 \times 1$ vector of parameters with $\mu$ and $\beta$ as elements, $\mathbf{Z}$ is a $J \times I$ design matrix with 0's and 1's, i.e. $\mathbf{Z} = \overset{I}{\underset{i=1}{\oplus}} \underline{1}_{J_i \times 1}$, $\underline{U}$ is an $I \times 1$ vector of random effects, and $\underline{E}$ is a $J \times 1$ vector of random error terms. In particular,

$$
\underline{Y} = \begin{pmatrix} Y_{11} \\ \vdots \\ Y_{1J_1} \\ Y_{21} \\ \vdots \\ Y_{2J_2} \\ \vdots \\ Y_{I1} \\ \vdots \\ Y_{IJ_I} \end{pmatrix}, \quad
\mathbf{X} = \begin{pmatrix} 1 & X_{11} \\ \vdots & \vdots \\ 1 & X_{1J_1} \\ 1 & X_{21} \\ \vdots & \vdots \\ 1 & X_{2J_2} \\ \vdots & \vdots \\ 1 & X_{I1} \\ \vdots & \vdots \\ 1 & X_{IJ_I} \end{pmatrix}, \quad
\underline{\alpha} = \begin{pmatrix} \mu \\ \beta \end{pmatrix}, \quad
\underline{U} = \begin{pmatrix} A_1 \\ \vdots \\ A_I \end{pmatrix},
$$

and

$$
\mathbf{Z} = \overset{I}{\underset{i=1}{\oplus}} \underline{1}_{J_i} = \begin{pmatrix} \underline{1}_{J_1} & \underline{0}_{J_1} & \cdots & \underline{0}_{J_1} \\ \underline{0}_{J_2} & \underline{1}_{J_2} & \cdots & \underline{0}_{J_2} \\ \vdots & \vdots & \ddots & \vdots \\ \underline{0}_{J_I} & \underline{0}_{J_I} & \cdots & \underline{1}_{J_I} \end{pmatrix}
$$

where $\underline{1}_{J_i}$ is a $J_i \times 1$ column vector of 1's and $\oplus$ is the direct sum operator. The direct sum of two matrices $\mathbf{P}$ and $\mathbf{Q}$ is defined as

$$\mathbf{P} \oplus \mathbf{Q} = \begin{pmatrix} \mathbf{P} & 0 \\ 0 & \mathbf{Q} \end{pmatrix}.$$

By the assumptions in (2.1) the response variables have a multivariate normal distribution

$$\underline{Y} \sim N(\mathbf{X}\underline{\alpha}, \sigma_A^2 \mathbf{Z}\mathbf{Z}' + \sigma_E^2 \mathbf{D}_{J.}) \tag{2.4}$$

where $\mathbf{D}_{J.}$ is a $J. \times J.$ identity matrix.

$$E(\underline{Y}) = E(\mathbf{X}\underline{\alpha} + \mathbf{Z}\underline{U} + \underline{E})$$

$$= \mathbf{X}\underline{\alpha} + E(\mathbf{Z}\underline{U}) + E(\underline{E})$$

$$= \mathbf{X}\underline{\alpha}, \qquad \text{and}$$

$$V(\underline{Y}) = V(\mathbf{X}\underline{\alpha} + \mathbf{Z}\underline{U} + \underline{E})$$

$$= V(\mathbf{Z}\underline{U}) + V(\underline{E})$$

$$= \mathbf{Z}V(\underline{U})\mathbf{Z}' + \sigma_E^2 \mathbf{D}_{J.}$$

$$= \sigma_A^2 \mathbf{Z}\mathbf{Z}' + \sigma_E^2 \mathbf{D}_{J.}$$

In order to define unweighted sums of squares, the vector of means of response variables of primary level and associated variance component matrix are needed. These are defined in matrix notation as

$$\mathbf{M}\underline{Y} = [\bar{Y}_{1.}, \bar{Y}_{2.}, ..., \bar{Y}_{I.}]' = \underline{Y}_M \tag{2.5}$$

where

$$\mathbf{M} = \overset{I}{\underset{i=1}{\oplus}} [J_i^{-1} \underline{1}'_{J_i \times 1}] = \begin{pmatrix} \frac{1}{J_1} \underline{1}_{J_1} & \underline{0} & \cdots & \underline{0} \\ \underline{0} & \frac{1}{J_2} \underline{1}_{J_2} & \cdots & \underline{0} \\ \vdots & \vdots & \ddots & \vdots \\ \underline{0} & \underline{0} & \cdots & \frac{1}{J_I} \underline{1}_{J_I} \end{pmatrix}$$

and $\bar{Y}_{i.}$ is the mean of response variables of the $i$th primary level. The expectation and variance of vector of means of response variables of the primary level are

$$E(\underline{Y}_M) = E(\mathbf{M}\underline{Y})$$

$$= \mathbf{M}\mathbf{X}\underline{\alpha}$$

$$= \mathbf{X}_M \underline{\alpha}, \qquad \text{and}$$

where $\mathbf{X}_M = \mathbf{M}\mathbf{X}$

$$V(\underline{Y}_M) = V(\mathbf{M}\underline{Y}) \tag{2.6}$$

$$= \mathbf{M}(\sigma_A^2 \mathbf{Z}\mathbf{Z}' + \sigma_E^2 \mathbf{D}_{J.})\mathbf{M}'$$

$$= \sigma_A^2 \mathbf{M}\mathbf{Z}\mathbf{Z}'\mathbf{M}' + \sigma_E^2 \mathbf{M}\mathbf{D}_{J.}\mathbf{M}'$$

$$= \sigma_A^2 \mathbf{D}_I + \sigma_E^2 \mathbf{M}\mathbf{M}'$$

since $\mathbf{M}\mathbf{M}' = diag[J_i^{-1}]$ and $\mathbf{M}\mathbf{Z} = \mathbf{D}_I$ where $\mathbf{D}_I$ is an $I \times I$ identity matrix. Thus, the vector of means of response variables of primary level has a multivariate normal distribution

$$\underline{Y}_M \sim N(\mathbf{X}_M \underline{\alpha}, \ \mathbf{V}_M) \tag{2.7}$$

where $\mathbf{V}_M = \sigma_A^2 \mathbf{D}_I + \sigma_E^2 \mathbf{M}\mathbf{M}'$.

# 3. DISTRIBUTIONAL PROPERTY OF ERROR SUMS OF SQUARES

In this section we report distributional results used to derive confidence intervals. Four regression coefficient estimators are considered. Consider weighted between regression coeficient estimator $\hat{\beta}_{WB}$ that is obtained from least squares regression of $\bar{Y}_{i.}$ on $\bar{X}_{i.}$ with weigh $J_i$ for each primary level $i$ and $\hat{\beta}_{WB}$ is written as

$$\hat{\beta}_{WB} = \frac{S_{wxya}}{S_{wxxa}}.$$

The weighted between regression coefficient estimator is the second element of the vector

$$(\mathbf{X}_M' \mathbf{W} \mathbf{X}_M)^{-1} \mathbf{X}_M' \mathbf{W} \underline{Y}_M \tag{3.1}$$

where $\mathbf{W} = diag[J_i]$. The error sum of squares

$$R_{WB} = S_{wyya} - \hat{\beta}_{WB}^2 S_{wxxa}$$

$$= \underline{Y}_M' \mathbf{A}_W \underline{Y}_M \tag{3.2}$$

where $\mathbf{A}_W = \mathbf{W} - \mathbf{W} \mathbf{X}_M (\mathbf{X}_M' \mathbf{W} \mathbf{X}_M)^{-1} \mathbf{X}_M' \mathbf{W}$

Unweighted between regression coeffcient estimator considering primary level's means and their unwighted mean is used as an alternative of between regression coefficient estimator. Unweighted between regression coeficient estimator $\hat{\beta}_{UB}$ is obtained from the least squares regression of $\bar{Y}_{i.}$ on $\bar{X}_{i.}$ and $\hat{\beta}_{UB}$ is written as

$$\hat{\beta}_{UB} = \frac{S_{uxya}}{S_{uxxa}}$$

where,

$$S_{uxya} = \sum_{i=1}^{I} (\bar{X}_{i.} - \bar{\bar{X}}_{..})(\bar{Y}_{i.} - \bar{\bar{Y}}_{..})$$

$$S_{uxxa} = \sum_{i=1}^{I} (\bar{X}_{i.} - \bar{\bar{X}}_{..})^2,$$

$$\bar{\bar{X}}_{..} = \frac{\sum_{i=1}^{I} \bar{X}_{i.}}{I}, \quad \text{and}$$

$$\bar{\bar{Y}}_{..} = \frac{\sum_{i=1}^{I} \bar{Y}_{i.}}{I}.$$

The unweighted between regression coefficient estimator is the second element of the vector

$$(\mathbf{X}'_M \mathbf{X}_M)^{-1} \mathbf{X}'_M \underline{Y}_M \qquad (3.3)$$

The error sum of squares $R_A$ associated with this regression model is

$$R_{UB} = S_{uyya} - \hat{\beta}^2_{UB} S_{uxxa}$$

$$= \underline{Y}'_M \mathbf{A}_U \underline{Y}_M \qquad (3.4)$$

where $\mathbf{A}_U = \mathbf{D}_I - \mathbf{X}_M (\mathbf{X}'_M \mathbf{X}_M)^{-1} \mathbf{X}'_M$.

The within regression coefficient estimator $\hat{\beta}_T = S_{wxye}/S_{wxxe}$ is obtained from the least squares regression of $Y_{ij}$ on $X_{ij}$ and the grouping variables. The point estimator $\hat{\beta}_T$ is the second element of the vector

$$(\mathbf{X}^{*'} \mathbf{X}^*)^{-} \mathbf{X}^{*'} \underline{Y} \qquad (3.5)$$

where $\mathbf{X}^* = [\mathbf{X} \ \mathbf{Z}]$, and $(\mathbf{X}^{*'}\mathbf{X}^*)^-$ is a generalized inverse of $\mathbf{X}^{*'}\mathbf{X}^*$. The error sum of squares $R_T$ associated with this regression model is

$$R_T = S_{wyye} - \hat{\beta}_T^2 S_{wxxe}$$

$$= \underline{Y}'\mathbf{T}\underline{Y} \tag{3.6}$$

where $\mathbf{T} = \mathbf{D}_{J.} - \mathbf{X}^*(\mathbf{X}^{*'}\mathbf{X}^*)^-\mathbf{X}^{*'}$. and $\mathbf{D}_{J.}$ is an identity matrix of order $J.$.

Finally, the total regression coefficient estimator

$$\hat{\beta}_L = \frac{(S_{wxya} + S_{wxye})}{(S_{wxxa} + S_{wxxe})}.$$

is obtained from the least squares regression of $Y_{ij}$ on $X_{ij}$. The point estimator $\hat{\beta}_L$ is the second element of the vector

$$(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\underline{Y}. \tag{3.7}$$

The error sum of squares $R_L$ associated with this regression model is

$$R_L = (S_{wyya} + S_{wyye}) - \hat{\beta}_L^2(S_{wxxa} + S_{wxxe}) - R_{WB} - R_T$$

$$= \underline{Y}'(\mathbf{L} - \mathbf{M}'\mathbf{A}_W\mathbf{M} - \mathbf{T})\underline{Y} \tag{3.8}$$

where $\mathbf{L} = \mathbf{D}_{J.} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$.

**Theorem 1.**

$R_T/\sigma_E^2$ a chi-squared random variable with $J. - I - 1$ degree of freedom.

**Proof.** Notice that $\mathbf{T}$ is idempotent. It can be shown that $\mathbf{X}^*(\mathbf{X}^{*'}\mathbf{X}^*)^-\mathbf{X}^{*'}\mathbf{X} =$

$\mathbf{X}$ and $\mathbf{X}^*(\mathbf{X}^{*'}\mathbf{X}^*)^-\mathbf{X}^{*'}\mathbf{Z} = \mathbf{Z}$ by Theorem 7.1 in Searle (1987, p. 218). Therefore, as may be easily verified, $\mathbf{TX} = \mathbf{0}$ and $\mathbf{TZ} = \mathbf{0}$. It follows that

$$E(R_T) = E(\underline{Y}'\mathbf{T}\underline{Y})$$

$$= tr(\mathbf{TV}) + \underline{\alpha}'\mathbf{X}'\mathbf{TX}\underline{\alpha}$$

$$= tr((\mathbf{D}_{J.} - \mathbf{X}^*(\mathbf{X}^{*'}\mathbf{X}^*)^-\mathbf{X}^{*'})(\sigma_A^2\mathbf{ZZ}' + \sigma_E^2\mathbf{D}_{J.}))$$

$$= tr(\sigma_E^2(\mathbf{D}_{J.} - \mathbf{X}^*(\mathbf{X}^{*'}\mathbf{X}^*)^-\mathbf{X}^{*'}))$$

$$= tr(\sigma_E^2\mathbf{T})$$

$$= \sigma_E^2 r(\mathbf{T})$$

$$= (J_. - I - 1)\sigma_E^2.$$

The distribution of $R_T$ is determined by writing $R_T/\sigma_E^2 = \underline{Y}'(\mathbf{T}/\sigma_E^2)\underline{Y}$ and noting $(\mathbf{T}/\sigma_E^2)\mathbf{V} = \mathbf{T}(\sigma_A^2\mathbf{ZZ}' + \sigma_E^2\mathbf{D}_{J.})/\sigma_E^2 = \mathbf{T}$. By Theorem 2 in Searle (1971, p. 57) $R_T/\sigma_E^2$ is a che-squared random variable with $J_. - I - 1$ degree of freedom.

**Theorem 2.**

If $\sigma_A^2 = 0$, then $R_{WB}/\sigma_E^2$ is a chi-squared random variable with $I - 2$ degree of freedom.

**Proof.**

Notice that

$$\mathbf{A}_W\mathbf{V}_M = (\mathbf{W} - \mathbf{WX}_M(\mathbf{X}_M'\mathbf{WX}_M)^{-1}\mathbf{X}_M'\mathbf{W})(\sigma_A^2\mathbf{D}_I + \sigma_E^2\mathbf{M}')$$

$$= \sigma_A^2\mathbf{W} - \sigma_A^2\mathbf{WX}_M(\mathbf{X}_M'\mathbf{WX}_M)^{-1}\mathbf{X}_M'\mathbf{W}$$

$$+ \sigma_E^2\mathbf{D}_I - \sigma_E^2\mathbf{WX}_M(\mathbf{X}_M'\mathbf{WX}_M)^{-1}\mathbf{X}_M'$$

since $\mathbf{WMM}' = diag[J_i] \cdot diag[J_i^{-1}] = \mathbf{D}_I$ and

$$tr(\mathbf{WX}_M(\mathbf{X}_M'\mathbf{WX}_M)^{-1}\mathbf{X}_M'\mathbf{W}) = k_1$$

where

$$k_1 = (\sum_{i=1}^{I} J_i \bar{X}_{i.} \sum_{i=1}^{I} J_i^2 - 2\sum_{i=1}^{I} J_i \bar{X}_{i.} \sum_{i=1}^{I} J_i^2 \bar{X}_{i.} + \sum_{i=1}^{I} J_i \sum_{i=1}^{I} J_i^2 \bar{X}_{i.}^2)/(J.S_{xxa}).$$

It follows that

$$
\begin{aligned}
E(R_{WB}) &= E(\underline{Y}_M' \mathbf{A}_W \underline{Y}_M) \\
&= tr(\mathbf{A}_W \mathbf{V}_M) + \underline{\alpha}'\mathbf{X}_M'\mathbf{A}_W\mathbf{X}_M\underline{\alpha} \\
&= \sigma_A^2(tr(\mathbf{W}) - k_1) + \sigma_E^2 tr(\mathbf{D}_I - \mathbf{WX}_M(\mathbf{X}_M'\mathbf{WX}_M)^{-1}\mathbf{X}_M') \\
&= (J. - k_1)\sigma_A^2 + (I - 2)\sigma_E^2
\end{aligned}
$$

since $\mathbf{A}_W\mathbf{X}_M = \mathbf{0}$. The distribution of $R_{WB}$ is determined by writing

$$R_{WB}/\sigma_E^2 = \underline{Y}_M'(\mathbf{A}_W/\sigma_E^2)\underline{Y}_M$$

and noting

$$
\begin{aligned}
(\mathbf{A}_W/\sigma_E^2)\mathbf{V}_M &= (\sigma_A^2/\sigma_E^2)\mathbf{A}_W + \mathbf{A}_W\mathbf{MM}' \\
&= (\sigma_A^2/\sigma_E^2)\mathbf{A}_W + \mathbf{D}_I - \mathbf{WX}_M(\mathbf{X}_M'\mathbf{WX}_M)^{-1}\mathbf{X}_M'
\end{aligned}
$$

since $\mathbf{WMM}' = \mathbf{D}_I$. Note that $\mathbf{D}_I - \mathbf{WX}_M(\mathbf{X}_M'\mathbf{WX}_M)^{-1}\mathbf{X}_M'$ is idempotent. It follows that $R_{WB}/\sigma_E^2$ is a chi-squared random variable with $I - 2$ degrees of freedom if $\sigma_A^2 = 0$.

**Theorem 3.**

If $\sigma_E^2 = 0$, then $R_{UB}/\sigma_A^2$ is a chi-squared random variable with $I - 2$ degrees of freedom.

**Proof.**

Notice that

$$\mathbf{A}_U \mathbf{V}_M = (\mathbf{D}_I - \mathbf{X}_M(\mathbf{X}_M'\mathbf{X}_M)^{-1}\mathbf{X}_M')(\sigma_A^2 \mathbf{D}_I + \sigma_E^2 \mathbf{M}\mathbf{M}')$$

$$= \sigma_A^2 (\mathbf{D}_I - \mathbf{X}_M(\mathbf{X}_M'\mathbf{X}_M)^{-1}\mathbf{X}_M')$$

$$+ \sigma_E^2 (\mathbf{M}\mathbf{M}' - \mathbf{X}_M(\mathbf{X}_M'\mathbf{X}_M)^{-1}\mathbf{X}_M'\mathbf{M}\mathbf{M}'),$$

$$tr(\mathbf{M}\mathbf{M}') = \sum_{i=1}^{I} (1/J_i), \quad \text{and}$$

$$tr(\mathbf{X}_M(\mathbf{X}_M'\mathbf{X}_M)^{-1}\mathbf{X}_M'\mathbf{M}\mathbf{M}') = k_2$$

where

$$k_2 = \frac{\sum_{i=1}^{I} \sum_{k=1}^{I} (\bar{X}_{i.} - \bar{X}_{k.})^2 / J_k}{I \cdot S_{uxxa}}$$

and $\mathbf{A}_U \mathbf{X}_M = \mathbf{0}$. It follows that

$$E(R_{UB}) = E(\underline{Y}_M' \mathbf{A}_U \underline{Y}_M)$$

$$= tr(\mathbf{A}_U \mathbf{V}_M) + \underline{\alpha}' \mathbf{X}_M' \mathbf{A}_U \mathbf{X}_M \underline{\alpha}$$

$$= \sigma_A^2 tr(\mathbf{D}_I - \mathbf{X}_M) + \sigma_E^2 (\sum_{i=1}^{I} (1/J_i) - k_2)$$

$$= (I - 2)\sigma_A^2 + (\sum_{i=1}^{I} (1/J_i) - k_2)\sigma_E^2.$$

The distribution of $R_{UB}$ is determined by writing

$$R_{UB}/\sigma_A^2 = \underline{Y}_M' (\mathbf{A}_U/\sigma_A^2)\underline{Y}_M$$

and noting $(\mathbf{A}_U/\sigma_A^2)\mathbf{V}_M = \mathbf{A}_U + (\sigma_E^2/\sigma_A^2)\mathbf{A}_U\mathbf{M}\mathbf{M}'$. Note that $\mathbf{A}_U$ is idempotent. Thus $R_{UB}/\sigma_A^2$ is a chi-squared random variable with $I - 2$ degrees of freedom if $\sigma_E^2 = 0$.

**Theorem 4.**

$R_{WB}/\sigma_E^2$ and $R_T/\sigma_E^2$ are independent and $R_{UB}/\sigma_A^2$ and $R_T/\sigma_E^2$ are independent.

**Proof.**

Notice that

$$\mathbf{M}'\mathbf{A}_W\mathbf{M}(\sigma_A^2\mathbf{Z}\mathbf{Z}' + \sigma_E^2\mathbf{D}_{J.})\mathbf{T} = \sigma_A^2\mathbf{M}'\mathbf{A}_W\mathbf{M}\mathbf{Z}\mathbf{Z}'\mathbf{T} + \sigma_E^2\mathbf{M}'\mathbf{A}_W\mathbf{M}\mathbf{T}$$

$$= \mathbf{0}$$

using $\mathbf{M}\mathbf{T} = \mathbf{M}\mathbf{M}'\mathbf{Z}'\mathbf{T} = \mathbf{0}$ since $\mathbf{M} = \mathbf{M}\mathbf{M}'\mathbf{Z}'$ and $\mathbf{Z}'\mathbf{T} = \mathbf{0}$. Accordinly $R_{WB}/\sigma_E^2$ and $R_T/\sigma_E^2$ are independent. Note that

$$\mathbf{M}'\mathbf{A}_U\mathbf{M}(\sigma_A^2\mathbf{Z}\mathbf{Z}' + \sigma_E^2\mathbf{D}_{J.})\mathbf{T} = \sigma_A^2\mathbf{M}'\mathbf{A}_U\mathbf{M}\mathbf{Z}\mathbf{Z}'\mathbf{T} + \sigma_E^2\mathbf{M}'\mathbf{A}_U\mathbf{M}\mathbf{T}$$

$$= \mathbf{0}$$

Thus $R_{UB}/\sigma_A^2$ and $R_T/\sigma_E^2$ are independent.

Olsen et al.(1976), Thomas and Hultquist(1978), and El-Bassiouni(1994) used spectral decomposition method to obtain following statistics. They proposed a statistic $SSM = \mathbf{U}'\mathbf{U}$ which is asymptotically chi-squared distributed. In particular,

$$\frac{\mathbf{U}'\mathbf{U}}{(\sigma_A^2 + \sigma_E^2/\lambda_H)} \to \chi_{(I-1)}^2 \quad \text{as} \quad \sigma_E^2 \to 0$$

where $\mathbf{U} = \mathbf{C}^+\mathbf{Z}'(\mathbf{D}_{J.} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')\mathbf{Z}$, $\lambda_H$ is the harmonic mean of positive eigenvalues, $\lambda_i$, of $C$,

$$\lambda_H = \frac{\sum\limits_{i=1}^{I} r_i}{(\sum\limits_{i=1}^{I} r_i/\lambda_i)},$$

and $r_i$ is the multiplicity of positive eigenvalue $\lambda_i$. Thus,

$$E(\mathbf{U}'\mathbf{U}) \cong (I - 1)(\sigma_A^2 + \frac{\sigma_E^2}{\lambda_H}).$$

It was also shown that $\mathbf{U}'\mathbf{U}/(\sigma_A^2 + \sigma_E^2/\lambda_H)$ and $R_T/\sigma_E^2$ are independent.

If the covariate values within each group are same, this proposed statistics becomes the error sum of squares associated with unweighted between regression coefficient and the total regresion coefficient estimator reduces to the weighted between regression coefficient estimator. That is, if $X_{ij} = X_i$ for all $j$, then $SSM = R_{UB}$ and $\hat{\beta}_L = \hat{\beta}_{WB}$. If group means of the covariate values are all same, i.e., $\bar{X}_{i.} = \bar{X}_{..} = \bar{\bar{X}}_{..}$ for all $i$, then $\mathbf{X}_M$ is linearly dependent and $\hat{\beta}_{WB}$ and $\hat{\beta}_{UB}$ are not defined.

# 4. CONFIDENCE INTERVAL ON $\sigma_A^2$

The expected mean squares are sumarized using the distributional property of error sums of squares.

$$E(S_{WB}^2) = c_1\sigma_A^2 + \quad \sigma_E^2 = \theta_{WB}, \tag{4.1a}$$

$$E(S_{UB}^2) = \quad \sigma_A^2 + c_2\sigma_E^2 = \theta_{UB}, \quad \text{and} \tag{4.1b}$$

$$E(S_T^2) = \quad\quad\quad \sigma_E^2 = \theta_T \tag{4.1c}$$

where

$$S_{WB}^2 = \frac{R_{WB}}{(I-2)},$$

$$S_{UB}^2 = \frac{R_{UB}}{(I-2)},$$

$$S_T^2 = \frac{R_T}{(J-I-1)},$$

$$c_1 = \frac{(J-k_1)}{(I-2)}, \quad \text{and}$$

$$c_2 = \frac{(\sum_i (1/J) - k_2)}{(I-2)}.$$

The mean square errors, $S_{WB}^2$ and $S_{UB}^2$, are independent of $S_T^2$ and they are exactly chi-squared distributed depending on cases where $\sigma_A^2 = 0$ and $\sigma_E^2 = 0$, respectively.

In the case where $\sigma_A^2 \to 0$, $S_{WB}^2$ and $S_T^2$ should be used to construct confidence intervals on $\sigma_A^2$. The variance component $\sigma_A^2$ can be represented by functions of expected mean squares in (4.1a) and (4.1c),

$$\sigma_A^2 = (\theta_{WB} - \theta_T)/c_1$$

An approximate confidence interval on $\sigma_A^2$ can be constructed using the method of Ting et al.(1990). In particular, the $1 - 2\alpha$ two-sided confidence interval for this form of $\sigma_A^2$ is

$$\frac{1}{c_1}[(S_{WB}^2 - S_T^2) - (G_1^2 S_{WB}^4 + G_2^2 S_T^4 + G_{12} S_{WB}^2 S_T^2)^{\frac{1}{2}};$$

$$(S_{WB}^2 - S_T^2) + (H_1^2 S_{WB}^4 + H_2^2 S_T^4 + H_{12} S_{WB}^2 S_T^2)^{\frac{1}{2}}] \tag{4.2}$$

where

$$F_1 = F_{(\alpha:I-2,J.-I-1)},$$

$$F_2 = F_{(1-\alpha:I-2,J.-I-1)},$$

$$G_1 = 1 - \frac{1}{F_{(\alpha:I-2,\infty)}},$$

$$G_2 = \frac{1}{F_{(1-\alpha:J.-I-1,\infty)}} - 1,$$

$$G_{12} = \frac{[(F_1 - 1)^2 - G_1^2 F_1^2 - G_2^2]}{F_1},$$

$$H_1 = \frac{1}{F_{(1-\alpha:I-2,\infty)}} - 1,$$

$$H_2 = 1 - \frac{1}{F_{(\alpha:J.-I-1,\infty)}},$$

$$H_{12} = \frac{[(1 - F_2)^2 - H_1^2 F_2^2 - H_2^2]}{F_2},$$

and $F_{(\delta:n_1,n_2)}$ is the $F$-value for $n_1$ and $n_2$ degrees of freedom with $\delta$ area to the right. Since $\sigma_A^2 > 0$, any negative bound is defined to be zero. Interval (4.2) is referred to as TINGW method.

Another approach is adapting generalized p-values method proposed by Khuri et al.(1998) to construct an approximate confidence interval on $\sigma_A^2$. It was shown in Chapter 3 that $(I-1)S_M^2/(\sigma_A^2 + \sigma_E^2/\lambda_H)$ is chi-squared distributed with $(I-1)$ degrees of freedom as $\sigma_E^2$ approaches zero,

$$\frac{(J_{..} - I - 1)S_T^2}{\sigma_E^2} \sim \chi^2_{(J_{..}-I-1)},$$

and they are independent where $S_M^2 = SSM/(I-1)$. Thus, using this property, the estimators of $\sigma_E^2$ are obtained by $(J_{..} - I - 1)s_T^2/U_1$ where $s_T^2$ is an observed value of $S_T^2$ and $U_1$ has a chi-squared distribution with $(J_{..} - I - 1)$ degrees of freedom. The estimators of $\sigma_{AE}^2$ are obtained by $(I-1)s_M^2/U_2$ where $\sigma_{AE}^2 = \sigma_A^2 + \sigma_E^2/\lambda_H$, $s_M^2$ is an observed value of $S_M^2$, and $U_2$ has a chi-squared distribution with $(I-1)$ degrees of freedom. Thus, a generalized pivotal quantity $\sigma_A^2$ can by represented as

$$\sigma_A^2 = \frac{(I-1)s_M^2}{U_2} - \frac{1}{\lambda_H} \cdot \frac{(J_{..} - I - 1)s_T^2}{U_1}.$$

Accordingly, an approximate $1 - 2\alpha$ two-sided confidence interval for this form of $\sigma_A^2$ is

$$[C_\alpha \quad ; \quad C_{1-\alpha}] \tag{4.3}$$

where $C_\alpha$ is the $\alpha$th percentile of the distribution constructed by the generalized pivotal quantity. Interval (4.3) is referred to as GPQ method.

When $\sigma_E^2$ approaches zero, $S_{UB}^2$ and $S_T^2$ can be used and $\sigma_A^2$ is represented

$$\sigma_A^2 = \theta_{UB} - c_2\theta_T$$

from (4.1b) and (4.1c). The Ting et al. $1 - 2\alpha$ two-sided confidence interval for this form of $\sigma_A^2$ is

$$[S_{UB}^2 - c_2 S_T^2 - (G_1^2 S_{UB}^4 + c_2^2 G_2^2 S_T^4 + c_2 G_{12} S_{UB}^2 S_T^2)^{\frac{1}{2}};$$

$$S_{UB}^2 - c_2 S_T^2 + (H_1^2 S_{UB}^4 + c_2^2 H_2^2 S_T^4 + c_2 H_{12} S_{UB}^2 S_T^2)^{\frac{1}{2}}] \qquad (4.4)$$

Interval (4.4) is referred to as TINGU method.

If $I = 3$, then $c_2 = 1/c_1$ and $c_2 A_W = A_U$. Thus $S_{WB}^2/c_1 = c_2 S_{WB}^2 = S_{UB}^2$ and TINGW and TINGU methods are same.

# 5. SIMULATION AND EXAMPLE

## 5.1 Simulation Study

The methods proposed in Chapter 4 are now compared using simulation study. The criteria for analyzing the performance of the methods are ; 1) their ability to maintain stated confidence coefficient, and 2) the average length of two-sided confidence intervals. Although shorter average interval lengths are preferable, it is necessary that the methods first maintain the stated confidence coefficient. Four unbalanced patterns were selected for simulation study and are shown in Table 5.1

### TABLE 5.1 Unbalanced Patterns Used in Simulation

| Pattern | I | $J_i$ |
|---------|----|------------------------------|
| 1 | 3 | 3 5 10 |
| 2 | 5 | 1 3 5 7 10 |
| 3 | 7 | 1 2 4 6 8 10 |
| 4 | 10 | 1 1 1 5 5 5 5 10 10 10 |

Let $\rho = \sigma_A^2/(\sigma_A^2 + \sigma_E^2)$. Without loss of generality $\sigma_A^2 = 1 - \sigma_E^2$ so that $\rho = \sigma_A^2$ and $1 - \rho = \sigma_E^2$. $A_i$ and $E_{ij}$ are independently generated from normal populations with zero means and variance $\rho$ and $1 - \rho$, respectively, using RAN-NOR routines of SAS. Values of $\mu$ and $\beta$ are respectively varied from -3 to 3 in increments of 1 so that 49 different combinations of $\mu$ and $\beta$ are used. Any fixed values of $X_{ij}$'s are given. Then $Y_{ij}$'s are calculated according to model (2.1) and $R_{WB}$, $R_T$, $SSM$, and $R_{UB}$ are computed as shown in Chapter 3. Simulated values for $S_{WB}^2$, $S_T^2$, $S_M^2$, and $S_{UB}^2$ are substituted into appropriate formula and the intervals are computed. Values of $\rho$ are varied from 0.001 to 0.999 in

increments of 0.1. Each value of $\rho$ is simulated 2000 times for each pattern. Two-sided intervals are computed based on equal tailed $F$-values. Confidence coefficients are determined by counting the number of the intervals that contain $\sigma_A^2$. Using the normal approximation to the binomial, if the true coefficient is 0.90, there is less than a 2.5% chance that an estimated confidence coefficient based on 2000 replications will be less than 0.8866. The average lengths of the two-sided confidence intervals are also calculated.

Table 5.2 and 5.3 present the results of the simulation for stated 90% confidence intervals on $\sigma_A^2$. The numbers in the body of Table 5.2 and 5.3 respectively report range of simulated confidence coefficients and average interval lengths and minimum and maximum values for the range as $\rho$ ranges from 0.001 to 0.999. Different combinations of $\mu$ and $\beta$ do not change the trend of simultion results and the change of minimum values of stated confidence coefficients is at most 0.012.

**TABLE 5.2 90% Range of Simulated Confidence Coefficients**

| Pattern | 1 | | | 2 | | |
|---|---|---|---|---|---|---|
| $\rho$ | TINGW | GPQ | TINGU | TINGW | GPQ | TINGU |
| 0.001 | 0.9005 | 0.9035 | 0.9005 | 0.8905 | 0.894 | 0.87 |
| 0.1 | 0.908 | 0.9045 | 0.908 | 0.8925 | 0.894 | 0.883 |
| 0.2 | 0.898 | 0.9085 | 0.898 | 0.8875 | 0.8955 | 0.884 |
| 0.3 | 0.893 | 0.904 | 0.893 | 0.89 | 0.896 | 0.898 |
| 0.4 | 0.9095 | 0.905 | 0.9095 | 0.8845 | 0.8985 | 0.907 |
| 0.5 | 0.9015 | 0.905 | 0.9015 | 0.898 | 0.8985 | 0.905 |
| 0.6 | 0.897 | 0.905 | 0.897 | 0.88 | 0.897 | 0.9045 |
| 0.7 | 0.899 | 0.9065 | 0.899 | 0.8655 | 0.896 | 0.8905 |
| 0.8 | 0.898 | 0.907 | 0.898 | 0.87 | 0.897 | 0.893 |
| 0.9 | 0.8975 | 0.905 | 0.8975 | 0.8635 | 0.896 | 0.8935 |
| 0.999 | 0.902 | 0.905 | 0.902 | 0.872 | 0.896 | 0.8925 |
| MAX | 0.9095 | 0.9085 | 0.9095 | 0.898 | 0.8985 | 0.907 |
| MIN | 0.893 | 0.9035 | 0.893 | 0.8635 | 0.894 | 0.87 |
| Pattern | 3 | | | 4 | | |
| 0.001 | 0.9 | 0.908 | 0.854 | 0.897 | 0.899 | 0.8135 |
| 0.1 | 0.8965 | 0.9095 | 0.866 | 0.901 | 0.8995 | 0.853 |
| 0.2 | 0.8955 | 0.907 | 0.888 | 0.8845 | 0.8985 | 0.868 |
| 0.3 | 0.8865 | 0.906 | 0.8955 | 0.8885 | 0.9015 | 0.8765 |
| 0.4 | 0.863 | 0.9065 | 0.883 | 0.869 | 0.905 | 0.8915 |
| 0.5 | 0.871 | 0.904 | 0.884 | 0.862 | 0.902 | 0.884 |
| 0.6 | 0.865 | 0.905 | 0.882 | 0.8715 | 0.903 | 0.891 |
| 0.7 | 0.8645 | 0.9055 | 0.895 | 0.862 | 0.9025 | 0.913 |
| 0.8 | 0.8735 | 0.901 | 0.907 | 0.857 | 0.902 | 0.898 |
| 0.9 | 0.858 | 0.9005 | 0.895 | 0.841 | 0.902 | 0.899 |
| 0.999 | 0.8685 | 0.8995 | 0.9045 | 0.856 | 0.9 | 0.904 |
| MAX | 0.9 | 0.9095 | 0.907 | 0.901 | 0.905 | 0.913 |
| MIN | 0.858 | 0.8995 | 0.854 | 0.841 | 0.8985 | 0.8135 |

**TABLE 5.3 90% Range of Average Interval Lengths**

| Pattern | 1 | | | 2 | | |
|---|---|---|---|---|---|---|
| $\rho$ | TINGW | GPQ | TINGU | TINGW | GPQ | TINGU |
| 0.001 | 44.670385 | 4.7037761 | 44.670385 | 1.6963909 | 1.7565652 | 2.4211901 |
| 0.1 | 59.306305 | 6.203999 | 59.306305 | 2.3232264 | 2.1452094 | 2.9775824 |
| 0.2 | 88.336569 | 7.7152441 | 88.336569 | 2.9723129 | 2.5340456 | 3.5968303 |
| 0.3 | 107.57395 | 9.2208541 | 107.57395 | 3.7117547 | 2.9165519 | 4.2278956 |
| 0.4 | 120.93906 | 10.721291 | 120.93906 | 4.1990992 | 3.2901384 | 4.6446833 |
| 0.5 | 142.92995 | 12.216601 | 142.92995 | 4.8447865 | 3.6554658 | 5.2186535 |
| 0.6 | 168.20474 | 13.706938 | 168.20474 | 5.6089767 | 4.0127826 | 5.7815406 |
| 0.7 | 185.53583 | 15.192167 | 185.53583 | 6.0088214 | 4.3633191 | 6.2783312 |
| 0.8 | 216.90426 | 16.673825 | 216.90426 | 7.0531375 | 4.7080909 | 7.0924004 |
| 0.9 | 245.96681 | 18.151819 | 245.96681 | 7.5738493 | 5.0491171 | 7.6091444 |
| 0.999 | 246.10563 | 19.612088 | 246.10563 | 8.4534132 | 5.3857702 | 8.5034603 |
| MAX | 246.10563 | 19.612088 | 246.10563 | 8.4534132 | 5.3857702 | 8.5034603 |
| MIN | 44.670385 | 4.7037761 | 44.670385 | 1.6963909 | 1.7565652 | 2.4211901 |
| Pattern | 3 | | | 4 | | |
| 0.001 | 0.8841862 | 1.2104748 | 1.4681376 | 0.3847124 | 0.7396307 | 0.7652631 |
| 0.1 | 1.3228924 | 1.5056313 | 1.8549788 | 0.6214983 | 0.9254198 | 0.9536777 |
| 0.2 | 1.848266 | 1.8000139 | 2.327806 | 0.8308489 | 1.1062396 | 1.1387063 |
| 0.3 | 2.2514865 | 2.0867147 | 2.6258684 | 1.0543185 | 1.2748048 | 1.3411049 |
| 0.4 | 2.7162211 | 2.3635822 | 3.127743 | 1.2549925 | 1.4296433 | 1.5096349 |
| 0.5 | 3.1438896 | 2.6303032 | 3.4974074 | 1.4559445 | 1.5735124 | 1.6897315 |
| 0.6 | 3.568139 | 2.8878889 | 3.7790143 | 1.5924883 | 1.7080275 | 1.7819612 |
| 0.7 | 4.0057759 | 3.1385151 | 4.1079702 | 1.806329 | 1.8377487 | 1.9402982 |
| 0.8 | 4.5275753 | 3.3852496 | 4.5825573 | 1.9948288 | 1.9655106 | 2.0982791 |
| 0.9 | 4.7874478 | 3.6297314 | 4.8973924 | 2.2045574 | 2.0937172 | 2.2359168 |
| 0.999 | 5.1810588 | 3.8715458 | 5.1928903 | 2.4126914 | 2.2214677 | 2.3882471 |
| MAX | 5.1810588 | 3.8715458 | 5.1928903 | 2.4126914 | 2.2214677 | 2.3882471 |
| MIN | 0.8841862 | 1.2104748 | 1.4681374 | 0.3847124 | 0.7396307 | 0.7652631 |

Simulation results are consistent with our study since TINGW method improves as $\rho$ approaches zero while TINGU method performs well as $\rho$ is closed to one across all values of $\rho$ for patterns 1. However, only GPQ method keeps the stated confidence coefficients for all $\rho$ values of four patterns. The average interval lengths of three methods generate wider intervals as $\rho$ increases for all four patterns. For smaller $\rho$ value, say $\rho \leq 0.1$, in pattern 3 and 4, TINGW method has shortest interval lengths. For other values of $\rho$ in four patterns, GPQ method has shortest interval length.

## 5.2 Numerical Example

The results of the simulation study are applied to a data set. Scheffe (1959, p216) wrote a data set of 94 observations for seven types of starch film and the data set was reproduced with permission of the author and publisher from Industrial Statistics by Freeman (1942). The dependent variable in the data set is the breaking strength in grams and the independent variable is the thickness in $10^{-4}$ inch from tests of starch film. The data set was constructed by selecting three types of starch, Potato, Canna, and Wheat. Three observations are selected from Potato, five from Canna, and ten from Wheat. This data set has the form of pattern 1 in Table 5.1 and is used to fit the simple linear regression model of the breaking strength on the thickness of starch film assuming an unbalanced nested error structure.

The selected data set is listed in Table 5.4 In order to apply the methods proposed in Chapter 4 to the data set a SAS code was programmed and 90% confidence intervals on $\sigma_A^2$ were calculated. The resulting intervals were given in Table 5.5 From SAS output the estimators $\hat{\sigma}_A^2$ and $\hat{\sigma}_E^2$ are computed as 8479.97 and 3063.89, respectively. Therefore, the estimate of the ratio of variance in primary unit to total variance $\hat{\rho}$ is 0.7345. GPQ should be used because it keeps the stated confidence level and generates the shortest interval length among three methods in patern 1 of Tables 5.2 and 5.3. The calculated interval lengths in Table 5.5 are consistent with the results in Table 5.3.

**TABLE 5.4 The Data Set Used for The Example**

| Type | Potato | | Canna | | Wheat | |
|------|--------|--------|-------|-------|-------|-------|
| Obs. | X | Y | X | Y | X | Y |
| 1 | 13.0 | 983.3 | 7.7 | 791.7 | 5.0 | 263.7 |
| 2 | 13.3 | 958.8 | 6.3 | 610.0 | 3.5 | 130.8 |
| 3 | 10.7 | 747.8 | 8.6 | 710.0 | 4.7 | 382.9 |
| 4 | | | 11.8 | 940.7 | 4.3 | 302.5 |
| 5 | | | 12.4 | 990.0 | 3.8 | 213.3 |
| 6 | | | | | 3.0 | 132.1 |
| 7 | | | | | 4.2 | 292.0 |
| 8 | | | | | 4.5 | 315.5 |
| 9 | | | | | 4.3 | 262.4 |
| 10 | | | | | 4.1 | 314.4 |

**TABLE 5.5 90% Confidence Intervals on $\sigma_A^2$**

| Methods | Lower bound | Upper bound | Length |
|---------|-------------|-------------|--------|
| TINGW(TINGU) | 2702.9 | 3359262.4 | 3356559.5 |
| GPQ | 915.6 | 216887.0 | 215971.4 |

# 6. CONCLUSIONS

Three approximate confidence intervals on the variance component of the primary level in a simple linear regression model with unbalanced nested error structure were proposed. The simulation study was conducted to compare the proposed intervals on the selected unbalanced patterns in Table 5.1 From Tables 5.2 and 5.3 if $\rho < 0.1$ in pattern 2 and $\rho \square 0.1$ in patterns 3 and 4, TINGW method is recommended because it keeps the stated confidence coefficients as well as shortest average interval lengths. For other values of $\rho$ in four patterns GPQ method is recommended.

# REFERENCES

[1] Burdick, R. K. and Eickman, J. (1986), "Confidence Intervals on The Among Group Variance Component in The Unbalanced One-fold Nested Design," *J. Stat. Comput. Simul.*, 26, 205-219.

[2] El-Bassiouni, M. Y.(1994), "Short Confidence Intervals for Variance Components," *Comm. Stat. -Theor. Meth.*, 23(7), 1915-1933.

[3] Freeman, H. A.(1942), *Industrial Statistics*, John Wiley & Sons, New York.

[4] Guven, Bilgehan (1995), "Maximum Likelihood Estimation in Simple Linear Regression with One-fold Nested Error," *Comm. Stat. -Theor. Meth.*, 24(1), 121-130.

[5] Khuri, A. I., Mathew, T., and Sinha, B. (1998), *Statistical Tests for Mixed Linear Models*, John Wiley & Sons, New York.

[6] Olsen, A., Seely, J., and Birkes, D.(1976), "Invariant Quadratic Unbiased Estimation for Two Varianc Components," *Ann. Stat.*, 4(5), 878-890.

[7] Scheffe, Henry.(1959), *The Analysis of Variance*, John Wiley & Sons, New York.

[8] Searle, S. R. (1971), *Linear Models*, John Wiley & Sons, New York.

[9] Searle, S. R.(1987), *Linear Models for Unbalanced Data*, John Wiley & Sons, New York.

[10] Thomas, J. D. and Hultquist, R. A.(1978), "Interval Estimaion for The Unbalanced Case of The One-way Random Effects Model," *Ann. Stat.*, 6(3), 582-587.

[11] Ting, N., Burdick, R.K., Graybill, F. A., Jeyaratnam, S., and Lu, T.-F. C.(1990), "Confidence Intervals on Linear Combinations of Variance Components," *J. Stat. Comput. Simul.*, 35, 135-143.

[12] Tong, L. I. and Cornelius, P. L. (1989), "Studies on The Estimation of The Slope Parameter in The Simple Linear Regression Model with One-fold Nested Error Structure," *Comm. Stat. - Simul., Comput.*, 18, 201-225.

[13] Tong, L. I. and Cornelius, P. L. (1991), "Studies on The Hypothesis Testing The Slope Parameter in The Simple Linear Regression Model with One-fold Nested Error Structure," *Comm. Stat. -Theor. Meth.*, 20(7), 2023-2043.

[14] Park, D. J. and Burdick, R. K. (1993), "Confidence Intervals on The Among Group Variance Component in A Simple Linear Regression Model with Nested Error Structure," *Comm. Stat. -Theor. Meth.*, 22(12), 3435-3452.

[15] Park, D. J. and Burdick, R. K. (1994), "Confidence Intervals on The Regression Coefficient in A Simple Linear Regression Model with A Balanced One-fold Nested Error Structure," *Comm. Stat. - Simul., Comput.*, 23(1), 43-58.

[16] Park, D. J. and H. M. Hwang (2002), "Confidence Intervals for the Mean Response in the Simple Linear Regression Model with Balanced Nested Error Structure," *Comm. Stat. -Theor. Meth.*, 31(1).

[17] Yu, Qing-Ling and Burdick, Richard K. (1995), " Confidence Intervals on Variance Components in Regression Models with Balanced (Q-1)-fold Nested Error Structure," *Comm. Stat.-Theor. Meth.*, 23(5), 1151-1167.