

GENERALIZED LEAST SQUARES
REGRESSION WITH PARTIALLY CORRELATED DATA

Andrew J. R. GILLESPIE¹, Syracuse University Research Consulting Services

ABSTRACT

One assumption of the ordinary least squares (OLS) regression model is that the observations on the response variable Y are statistically uncorrelated. In data sets obtained from cluster sampling (usually a more cost-effective sampling method than simple random sampling), this assumption is violated; observations may be uncorrelated between clusters, but correlated within clusters. Application of OLS methods to such data can lead to serious underestimates of the precision of estimated regression coefficients.

Many methods have been proposed for dealing with such intraclass correlation, including generalized least squares (GLS), random coefficient regressions, and jackknife estimators of variance. This paper describes two algorithms written in the SAS Interactive Matrix Language (SAS/IML™) to fit GLS models based on work by Scott and Holt (1984) and Fuller and Battese (1973). The algorithms pass through the data twice. The first pass is used to obtain an estimate of intraclass correlation; the second pass transforms the data into a form suitable for analysis by the common OLS procedures.

INTRODUCTION

The general linear model, expressed in matrix form, is commonly written as

$$Y = X'\beta + \epsilon \quad (1)$$

where Y is an observation on some response variable of interest, X is a p by 1 vector of observations on p predictor variables used to estimate Y , β is a p by 1 vector of unknown population parameters, ϵ is a random error expressing the difference between $X'\beta$ and Y , and X' denotes the transposed vector X .

This model may be used to derive estimates of the response variable Y as a function of a set of predictor variables X_k and an constant parameter vector β . Under certain assumptions, the ordinary least squares (OLS) estimate b_0 of β is unbiased and of minimum variance among all unbiased estimators. These assumptions are:

1. The predictor variables X_k ($k = 1$ to p) are fixed and known without error;
2. The expected value of ϵ given X is 0, i.e. $E(\epsilon|X) = 0$;
3. The conditional variance of ϵ given X is a constant parameter, i.e. $\text{Var}(\epsilon|X) = \sigma^2$;
4. The ϵ_i are statistically uncorrelated, i.e. $\text{Cov}(\epsilon_i, \epsilon_j; i \neq j) = 0$.

To estimate the variance of Y for hypothesis testing, it is necessary to assume that the ϵ_i are independent random variables, or are from a Normal (Gaussian) probability distribution. Assumptions 3 and 4 together imply a specific variance-covariance (or "dispersion") matrix structure for the vector ϵ of residuals, namely

$$D(\epsilon) = \sigma^2 I \quad (2)$$

where I = an n by n identity matrix. However, the assumptions described above are seldom simultaneously satisfied in real life; rather, the assumptions are taken to be approximately true, and departures are assumed small enough to be unimportant.

One common violation of these assumptions occurs as a result of cluster sampling. In cluster sampling, a population is first subdivided into overlapping or non-overlapping clusters. For example, a forest may be conceptually divided into an infinite number of plots as determined by locations of points as plot centers. A sample of m clusters is selected, then some or all of the elements in each cluster are measured to obtain the sample observations. The result is that observations within clusters will tend to be more similar to each other than observations from different clusters; i.e., there exists a certain amount of intracluster correlation which violates assumption 4 above. Instead of the dispersion structure of equation (2), the covariance between observations from the same cluster is non-zero, so the dispersion matrix for a sample consisting of m clusters takes a

block diagonal form:

$$D(e) = \begin{bmatrix} \mathbf{V}_1 & 0 & 0 & \dots & 0 \\ 0 & \mathbf{V}_2 & 0 & \dots & 0 \\ 0 & 0 & \mathbf{V}_3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \mathbf{V}_m \end{bmatrix}$$

where \mathbf{V}_i is an n_i by n_i matrix which represents the correlation structure within cluster i , n_i is the number of sample elements in cluster i , and 0 is a matrix of 0's of convenient order which implies the absence of correlations between clusters. \mathbf{V}_i has the form

$$\mathbf{V}_i = \sigma^2 \begin{bmatrix} 1 & \rho & \rho & \dots & \rho \\ \rho & 1 & \rho & \dots & \rho \\ \rho & \rho & 1 & \dots & \rho \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \rho & \dots & 1 \end{bmatrix}$$

where ρ = a measure of the intracluster correlation among the observed residuals $e_{ij} = (Y_{ij} - \hat{Y}_{ij})$. Note that ρ measures correlation after accounting for the effects of X on Y .

Several simulation studies (Briggs, 1981; Kotimaki 1981; Gillespie, 1985) have shown that application of ordinary least squares procedures to samples containing such intracluster correlation may result in an underestimate of σ_{bb} , the variance-covariance matrix of the estimated least squares parameter vector b_0 . An alternative estimator of β is the generalized least squares estimator

$$\beta_G = (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}(\mathbf{X}'\mathbf{V}^{-1}\mathbf{Y}) \quad (3)$$

where \mathbf{X} , \mathbf{Y} , and \mathbf{V} are as defined above. An operational estimator of β_G requires an estimate of ρ , the intracluster correlation coefficient. This paper discusses two different approaches for dealing with intracluster correlation.

The first approach, discussed by Scott and Holt (1983), works with an estimate of the matrix \mathbf{V} in equation (3) to calculate an estimate b_G of β_G . A modified form of an equation given in Cochran (1977) is used to estimate the intracluster correlation coefficient ρ from a sample with unequal cluster size. The SAS/IML routine first fits the OLS regression of Y on X to obtain the observed

residuals e_{ij} . Next, one must pass through the data one cluster at a time to complete the calculation of $\hat{\rho}$ = the estimate of ρ :

$$\hat{\rho} = \frac{2n \sum_{i=1}^m \sum_{j<u} (e_{ij})(e_{iu})}{\sum_{i=1}^m n_i (n_i-1) (n-p) S_{yy|x}} \quad (4)$$

Theoretically, once ρ is estimated, one could construct \mathbf{Y} , \mathbf{X} , and \mathbf{V} and proceed by the usual matrix manipulations to derive the parameter estimates. However, \mathbf{V} is a square matrix of dimension n = sample size. For an interesting problem where n is large (say over 300), \mathbf{V} might be too large for operations such as inversion or multiplication. One can take advantage of the block diagonal structure of \mathbf{V} and proceed iteratively, cluster by cluster, since it can be shown that

$$\mathbf{X}'\mathbf{V}^{-1}\mathbf{X} = \sum (\mathbf{X}_i'\mathbf{V}_i^{-1}\mathbf{X}_i) \quad (5)$$

$$\mathbf{X}'\mathbf{V}^{-1}\mathbf{Y} = \sum (\mathbf{X}_i'\mathbf{V}_i^{-1}\mathbf{Y}_i) \quad (6)$$

where \mathbf{X}_i , \mathbf{V}_i , and \mathbf{Y}_i correspond to the subsets of \mathbf{X} , \mathbf{V} , and \mathbf{Y} contributed by cluster i ($i = 1$ to m). Thus the second pass through the data consists of constructing the submatrices \mathbf{X}_i , \mathbf{V}_i , and \mathbf{Y}_i , then summing the appropriate products over all m clusters. The resulting products are then used to obtain the following parameter estimates:

$$b_G = [\sum (\mathbf{X}_i'\mathbf{V}_i^{-1}\mathbf{X}_i)]^{-1} [\sum \mathbf{X}_i'\mathbf{V}_i^{-1}\mathbf{Y}_i] \quad (7)$$

an estimate of the GLS regression coefficient vector;

$$S_{yy|x} = \{[\sum (\mathbf{Y}_i'\mathbf{V}_i^{-1}\mathbf{Y}_i)] - (b_G' (\sum \mathbf{X}_i'\mathbf{V}_i^{-1}\mathbf{Y}_i))\} / (n-p) \quad (8)$$

an estimate of the conditional variance of Y given X ; and

$$S_{bb} = S_{yy|x} \{[\sum (\mathbf{X}_i'\mathbf{V}_i^{-1}\mathbf{X}_i)]^{-1}\} \quad (9)$$

an estimate of the covariance matrix of b_G

The second approach, proposed by Fuller and Battese (1983), involves the construction of a transformational matrix \mathbf{T} . This matrix is used to premultiply \mathbf{Y} and \mathbf{X}

to produce transformed matrices, say \mathbf{U} and \mathbf{V} . The standard OLS calculations, applied to the transformed matrices \mathbf{U} and \mathbf{V} , yield parameter estimates which account for the intracluster correlation effect. Like the \mathbf{V} matrix in model 1, \mathbf{T} is block diagonal, where \mathbf{T}_i is constructed for each cluster i . \mathbf{T}_i is defined as

$$\mathbf{T}_i = \mathbf{I}_{n_i} - \left\{ \left(1 - \frac{\sigma_e^2}{(\sigma_e^2 + n_i \sigma_u^2)} \right)^{1/2} \left[\frac{1}{n_i} (\mathbf{J}_{n_i}) \right] \right\} \quad (10)$$

where σ_e^2 is a variance component due to the variation of the sampled individuals, σ_u^2 is a variance component due to the correlation among individuals within a cluster, and \mathbf{J}_{n_i} is an n_i by n_i matrix of 1's.

The authors suggested estimating the variance components σ_e^2 and σ_u^2 by the "fitting of constants" method outlined by Searle (1971). The first step is to regress the Y deviations from cluster means ($Y_{ij} - Y_{i.}$) on the analogous X deviations ($X_{ijk} - X_{i.k}$), where $Y_{i.}$ and $X_{i.k}$ denote the means in cluster i of Y and X_k respectively. This regression includes only those X variables for which the deviations are not 0 (i.e., ignoring the column of ones included for an intercept, since $1-1=0$ in each case). Then σ_e^2 is unbiasedly estimated by

$$\hat{\sigma}_e^2 = \mathbf{e}'\mathbf{e} / (n - m - p + \delta) \quad (11)$$

where \mathbf{e} = the vector of observed residuals from the above regression, and δ = the number of predictor variables for which all deviations are 0 (ie, the number of predictor variables ignored in the regression of residuals). This computation requires one pass through the data, to calculate the cluster residuals. The regression can then be done with normal matrix procedures.

The second variance component, σ_u^2 , is unbiasedly estimated by

$$\hat{\sigma}_u^2 = \frac{\mathbf{u}'\mathbf{u} - (n - p)\hat{\sigma}_e^2}{n - \text{trace}[(\mathbf{X}'\mathbf{X})^{-1} \sum (n_i^2 \mathbf{x}_{i.}' \mathbf{x}_{i.})]} \quad (12)$$

where \mathbf{u} is the residual vector from the least squares regression of Y on X , and $\mathbf{x}_{i.}$ is the mean vector of the observations on X_{ijk} in cluster i .

At this point, one has enough information to estimate \mathbf{T}_i for each cluster. The original cluster matrices \mathbf{Y}_i and \mathbf{X}_i are premultiplied by \mathbf{T}_i and concatenated, yielding transformed vectors $\mathbf{U} = \mathbf{T}\mathbf{Y}$ and $\mathbf{V} = \mathbf{T}\mathbf{X}$. This requires a second pass through the data. Finally, the generalized least squares model is fit by applying the OLS procedures to the transformed matrices \mathbf{U} and \mathbf{V} .

SUMMARY

Two algorithms have been developed using the SAS Interactive Matrix Language (SAS/IML) to fit Generalized Least Squares (GLS) models for data collected by a cluster sample. Such models are better suited to data from complex samples than the common Ordinary Least Squares (OLS) model, because the GLS model accounts for the effect of intracluster correlation among observations from the same cluster.

Copies of the SAS/IML code used to estimate the GLS models are available on request from Syracuse University Research Computing Services, 732 Ostrom Avenue, Syracuse, NY 13244. Bitnet address: RESCON@SUVM. SAS/IML™ is a registered trademark of SAS Institute Inc., Cary, NC, USA.

LITERATURE CITED

- Briggs, E.F. 1980. "Biomass table construction by regression modified for cluster sampling." M.S. Thesis, SUNY College of Environmental Science and Forestry, Syracuse, NY. 326 pp.
- Cochran, W.G. 1977. Sampling Techniques. 3rd. edition. New York: John Wiley & Sons. 428 pp.
- Gillespie, A.G. 1985. "Estimation of biomass tables by cluster sampling: Results of a simulation study." M.S. Thesis, SUNY College of Environmental Science and Forestry, Syracuse, NY. 314 pp.
- Fuller, W.A., and G.E. Battese. 1973. "Transformations for estimation of linear models with nested-error structure". Journal of the American Statistical Association 68:626-632.

Scott, A.J., and D. Holt. 1982.
"The effect of two stage
sampling on ordinary least
squares methods". Journal of
the American Statistical
Association 77:848-854.

Searle, S.R. 1971. "Topics in
variance component estimation".
Biometrics 27:1-76.

¹Current address: Institute of
Tropical Forestry, Box 25000, Rio
Piedras, Puerto Rico 00928-2500.