

Model Fitting Tests in the Analysis of Panel Data

Marcel de Toledo Vieira *

Federal University of Juiz de Fora, Juiz de Fora, Brazil marcel.vieira@ice.ufjf.br

Abstract

The interest for fitting models to panel data has grown in the last few decades. Modelling methods have to consider the variation in the response variable across the population as well as across time in this context. We review model fitting statistics under the classical (simple random sampling) approach, while we propose new developments on fitting measures when working under the complex sampling approach. We modify the Wald goodness of fit test in the context of models for covariance structures, which is shown to be equivalent to modifying the scaled test statistics. We also propose a modification for the Wald significance test for nested hypothesis. Goodness of fit indices are also modified in order to be utilised in the complex survey data context.

Keywords: covariance structure, goodness of fit, longitudinal surveys, multistage sampling.

1. Introduction

The interest for fitting models to longitudinal complex survey data has grown in the last few decades. Longitudinal data modelling methods have to consider the variation across the population as well as across time. A wide class of ‘regression-type’ models has found a broad range of useful applications with panel survey data (e.g. Diggle *et al.*, 2002). Furthermore, a number of methods have been developed in the survey sampling literature to take account of complex sampling schemes in the regression analysis of cross section survey data. See Skinner, Holt and Smith (1989), and Chambers and Skinner (2003), for references. Moreover, some previous work on estimation for panel data models under complex designs has been undertaken, for example, by Sutradhar and Kovacevic (2000), Skinner and Holmes (2003), Skinner and Vieira (2007), and Vieira and Skinner (2008). See Vieira (2009), for an overview.

In this paper we extend this broad approach to model fitting statistics in the longitudinal data modelling context, allowing for complex sampling designs. We review model fitting statistics under the classical (simple random sampling - srs) approach, while we propose new developments on fitting measures when working under the complex sampling approach. The models we consider are presented in Section 2. The paper proceeds in Section 3 to consider model testing statistics in both classic and complex survey contexts, and brief remarks are presented in Section 4.

2. Models

Let the finite population be denoted by \mathcal{U} , which is treated as fixed on occasions $1, \dots, T$. Let N represent the size of \mathcal{U} and $N_o = N \cdot T$. Let $\underline{Y}_i = (Y_{i1}, \dots, Y_{iT})'$ be a random vector containing T repeated observations on the study variable for unit $i = 1, 2, \dots, N$ over the T waves of the survey. Moreover, $E(\underline{Y}_i) = \underline{\mu}_i(\underline{\beta})$ is a $T \times 1$ vector, where $\underline{\mu}_i(\underline{\beta}) = [\mu(x_{i1}, \underline{\beta}), \dots, \mu(x_{iT}, \underline{\beta})]'$ and \underline{x}_{it} is a vector of values of the covariates.

A covariance structure model is a model for the $T \times T$ symmetric population variance-covariance matrix of \underline{Y}_i , which is

$$\Sigma = \text{COV}(\underline{Y}_i) = E\{[\underline{Y}_i - \underline{\mu}_i][\underline{Y}_i - \underline{\mu}_i]'\} = \Sigma(\underline{\theta}), \quad (1)$$

where $\text{COV}(\cdot)$ denote population covariance. We assume (1) is the same for each unit i , and that $k = T(T+1)/2$ distinct elements of the variance-covariance matrix $\Sigma(\underline{\theta})$ are constrained to be functions of the $b \times 1$ parameter vector $\underline{\theta}$, with $b < k$.

If we consider a uniform correlation model (Model A in Skinner and Holmes, 2003), $\Sigma(\underline{\theta})$ has diagonal values $\sigma_u^2 + \sigma_v^2$ and off-diagonal values σ_u^2 , where σ_v^2 is the

variation on the same individual, σ_u^2 and is the variance across individuals (Lindsey, 1994), and $\underline{\theta} = (\sigma_u^2, \sigma_v^2)'$, with $b=2$. If we consider a transitory random effects model with a first-order autoregressive process (Model B in Skinner and Holmes, 2003), then $\underline{\theta} = (\sigma_u^2, \sigma_v^2, \gamma)'$, where γ is a regression parameter and $-1 < \gamma < 1$.

3. Model testing

We assume that: (i) the observations are equally spaced in time; (ii) the sample size is 'large' relative to the number of repeated observations; (iii) the sample is selected on one occasion and then the same sample units are returned to on each of the $T-1$ subsequent waves; and (iv) there is no nonresponse. Estimation procedures for $\underline{\beta}$ are discussed in details by Skinner and Vieira (2007), Vieira and Skinner (2008), and Vieira (2009), where methods on inference about the covariance matrix Σ , for the variance estimation of $\hat{\Sigma}$, and for estimation of the parameter $\underline{\theta}$ are also developed under the complex surveys approach, including generalized least squares (GLS) and pseudo maximum likelihood (PML) methods.

Note that (1) is our covariance structure hypothesis. Model fit measures are used to assist in the evaluation of whether (1) is valid or not, and if not, such measures could assist to calculate the deviation of Σ from $\Sigma(\underline{\theta})$. In this section, let $\hat{\underline{\theta}}$ denote an estimator of $\underline{\theta}$ which minimizes either $F(\underline{\theta})_{ML}$ or $F(\underline{\theta})_{GLS}$, which are maximum likelihood (ML) and GLS fitting functions, respectively, and $F(\underline{\theta})_{PML}$ or $F(\underline{\theta})_{GLSC}$, which are PML and GLS fitting functions, under the complex sampling context.

3.1 Model testing in the classical context

In this sub-section we work under the assumptions of independent and identically distributed observations. Both Σ and $\Sigma(\underline{\theta})$ are unknown population parameters. Thus for calculating model fit measures we would in fact need to consider their estimators S and $\Sigma(\hat{\underline{\theta}})$, respectively, where $\Sigma(\hat{\underline{\theta}})$ is the covariance matrix evaluated at $\hat{\underline{\theta}}$.

In order to perform a goodness of fit test, we may initially define a null and a generic alternative hypothesis as $H_0 : \Sigma = \Sigma(\underline{\theta})$ against $H_1 : \Sigma$, which is an unrestricted covariance matrix (any $T \times T$ positive definite matrix). Let the population residual covariance matrix be denoted by E_p , so that

$$E_p = [\Sigma - \Sigma(\underline{\theta})]. \tag{2}$$

When H_0 is true, E_p is a zero matrix. The sample residual covariance matrix \hat{E} , defined as $\hat{E} = [S - \Sigma(\hat{\underline{\theta}})]$, is the simplest model fit measure. Let $[S_{it'} - \Sigma(\hat{\underline{\theta}})_{it'}]$ be individual sample residual covariances, where $S_{it'}$ and $\Sigma(\hat{\underline{\theta}})_{it'}$ are the it' th elements in S and $\Sigma(\hat{\underline{\theta}})$ respectively.

Furthermore, Jöreskog and Sörbom (1989) proposed the following statistic for summarising the residuals,

$$RMR = \sqrt{2 \cdot \frac{\sum_{t=1}^T \sum_{t'=1}^t [S_{it'} - \Sigma(\hat{\underline{\theta}})_{it'}]^2}{T(T+1)}},$$

where RMR stands for root mean-square residual. This measure may be adopted to compare the fit of two different models for the same data.

Sample residuals are not only affected by differences between Σ and $\Sigma(\underline{\theta})$, but also by the scales of \underline{Y}_i and by sampling fluctuations (errors). A direct solution for the scales issue may be to calculate correlation residuals as (Bollen, 1989) $r_{it'} - \hat{r}_{it'}$, where $r_{it'}$ is the sample correlation between \underline{Y}_{it} and $\underline{Y}_{it'}$, and $\hat{r}_{it'}$ denotes the model predicted correlation, so that

$$\hat{r}_{it'} = \frac{\Sigma(\hat{\theta})_{it'}}{\sqrt{\Sigma(\hat{\theta})_{ii} \cdot \Sigma(\hat{\theta})_{t't}}}$$

Although this correlation residual is allowed to range from -2 to $+2$, we should expect values rather close to zero for models with a reasonably good fit.

Regarding sampling errors, even when H_0 is true the expected amplitude of the individual sample residual covariances depends on n , and $\lim_{n \rightarrow \infty} \hat{E} = E_p$. Therefore, we present below a simultaneous significance test, based on sample residuals. Notice that the H_0 discussed above could be tested via a chi-square test.

Under multivariate normality of \underline{Y}_i and no covariates, if H_0 holds, if $\underline{\theta}$ is identified, and if

$$\Delta = \frac{\partial \{vech[\Sigma(\underline{\theta})]\}}{\partial \underline{\theta}}, \tag{3}$$

is of full rank b , where (3) is a $k \times b$ matrix of partial derivatives of elements of $vech[\Sigma(\underline{\theta})]$ with respect to the elements of $\underline{\theta}$, it holds that $n \cdot F(\hat{\theta})_{ML} \sim \chi^2_{k-b}$ under H_0 , when $F(\underline{\theta})_{ML}$ is evaluated at the final estimates and the model is true. When \underline{Y}_i is multivariate normally distributed, $n \cdot F(\hat{\theta})_{GLS}$ has the same property as $n \cdot F(\hat{\theta})_{ML}$. Note that $n \cdot F(\hat{\theta})_{GLS}$ is a Wald goodness of fit test statistic, which may be obtained as the minimum value of $F(\underline{\theta})_{GLS}$, when evaluated at $\hat{\theta}_{GLS}$.

Satorra (1989), propose asymptotically equivalent significance tests for the difference in chi-square statistics for nested models: (i) likelihood ratio test (LRT) or chi-square difference test); (ii) Lagrangian multiplier test (LMT) or efficient score test; and (iii) Wald test (WT); which all assume that $F(\underline{\theta})$ is asymptotic optimal, i.e. leads to efficient estimators and chi-square statistics. In general, these types of test aim to compare an ‘initial’ model with a restricted model, which has a sub-vector of parameters that is set to be equal to zero. There are also alternative approaches for performing model selection that have been proposed in the literature. Jöreskog and Sörbom (1989), for example, have proposed a goodness of fit index (GFI) and an adjusted index (AGFI), which penalises the models with more parameters.

4.2 Model testing under complex sampling

In this sub-section, we consider some further developments on covariance structure model fitting statistics when assuming the sample is selected under the complex survey design approach. According to Skinner, Holt and Smith (1989), ignoring the characteristics of the complex samples can lead to invalid statistical tests.

For calculating model fit measures in the present context we adopt S_w (Vieira and Skinner, 2008), the survey weighted sample covariance matrix, as an estimator of Σ , considering that $E(E_p(S_w)) \doteq E(S_N) = \Sigma$, as shown by Vieira (2009). We thus consider model fit measures which are functions of S_w and Σ .

For examining (2), i.e. for identifying components of the variance-covariance matrix that are not well fit, we may adopt the sample weighted residual covariance matrix $\hat{E}_c = S_w - \Sigma(\hat{\theta})$, where the subscript c denotes ‘complex’. Moreover, let $[S_{w, it'} - \Sigma(\hat{\theta})_{it'}]$ be the individual weighted sample residual covariances, where $S_{w, it'}$ is the it' th element in S_w and $\Sigma(\hat{\theta})_{it'}$. The RMR measure may be adapted to

$$RMR_c = \sqrt{2 \cdot \sum_{t=1}^T \sum_{t'=1}^t \frac{[S_{w, it'} - \Sigma(\hat{\theta})_{it'}]^2}{T(T+1)}}$$

which has as special case the RMR measure, when the sampling weights are constant.

Theory developed in the categorical complex survey data analysis and modelling literature (Rao and Scott, 1979) suggests that, under complex sampling, $n \cdot F(\underline{\theta})_{PML}$, would not be asymptotically chi-squared distributed. Simple corrections have been proposed, which may be adapted here to be applied to $n \cdot F(\underline{\theta})_{PML}$ in order to make it approximately chi-squared distributed. We follow an approach proposed by Skinner (1989, Section 3.4).

Remark 1: We initially consider the case of a GLSC type estimator, obtained by minimizing the $F(\underline{\theta})_{GLSC}$ fitting function, with matrix U given by

$$U = 2 \cdot K'(W \otimes W)K, \tag{4}$$

where W is any consistent estimator of Σ . Under this situation, a Wald goodness of fit test statistic is given by (Skinner, 1989)

$$X_{W,srs}^2 = n \cdot \{vech[S_w] - vech[\Sigma(\theta)]\}' U^{-1} \{vech[S_w] - vech[\Sigma(\theta)]\},$$

which implies,

$$X_{W,srs}^2 = n \cdot \left(\frac{1}{2}\right) \cdot tr\{[I - \Sigma(\theta)S_w^{-1}]^2\},$$

when S_w is considered as a choice for W in (4). Note, nevertheless, that according to Skinner (1989), the test statistic $X_{W,srs}^2$ is no longer asymptotically chi-squared distributed, but in fact asymptotically

$$X_{W,srs}^2 \sim \sum_{d=1}^{k-b} \lambda_d \chi_1^2$$

under H_0 , where χ_1^2 are independent chi-squared distributed random variables, and λ_d are non-zero eigenvalues of

$$H = U^{-1}C_c - U^{-1}C_c U^{-1} \cdot \Delta(\Delta' U^{-1} \Delta)^{-1} \Delta',$$

with C_c defined as the asymptotic covariance matrix of $vech[S_w]$. As a one moment approximation (Skinner, 1989)

$$\frac{(k-b) \cdot X_{W,srs}^2}{tr(H)},$$

is asymptotically distributed as a χ_{k-b}^2 , and may be adopted for testing the goodness of fit of a covariance structure model in a complex sampling context. We also consider substituting matrix U by \hat{C}_c , i.e.

$$X_W^2 = n \cdot \{vech[S_w] - vech[\Sigma(\theta)]\}' \hat{C}_c \{vech[S_w] - vech[\Sigma(\theta)]\},$$

where \hat{C}_c , given by

$$aCOV(\hat{\theta}) = n^{-1} (\Delta' U^{-1} \Delta)^{-1} \Delta' U^{-1} C_c U^{-1} \Delta (\Delta' U^{-1} \Delta)^{-1},$$

is a consistent estimator for C_c (Skinner and Holmes, 2003), with $C_c = VAR\{\sqrt{n} \cdot vech[S_w]\}$. In this context, X_W^2 is approximately distributed as χ_{k-b}^2 under H_0 (Skinner, 1989). Note that when U is consistent for C_c ,

$$X_W^2 = X_{W,srs}^2 = \frac{(k-b) \cdot X_{W,srs}^2}{tr(H)}, \text{ as } tr(H) = k-b \text{ in that situation.} \quad \blacksquare$$

We now consider a Wald significance test for nested hypothesis for situations where the PML fitting function is adopted. Note that we shall follow an approach proposed by Skinner (1989, Section 3.4).

Remark 2: We assume that the asymptotic covariance matrix of $vech[S_w]$ and the information matrix of the model are non-singular. Thus, let $\hat{\theta}_{PML,r}$ be the PML estimator for the restrictive (nested) model, and let $\hat{\theta}_{PML,u}$ be the PML estimator for the covariance structure model without constraints. Let θ^R be a $b^* \times 1$ vector when $\theta^R = 0$ corresponds to the restraints imposed to the unrestricted model, where $b^* < b$. We propose a modification to the Wald test, so that

$$WT_c = \hat{\theta}_{PML,u}^R \cdot \left\{ \begin{bmatrix} \frac{\partial \hat{\theta}_{PML,u}^R}{\partial \hat{\theta}_{PML,u}} \\ \frac{\partial \hat{\theta}_{PML,u}^R}{\partial \hat{\theta}_{PML,u}} \end{bmatrix}' \cdot [a \text{cov}(\hat{\theta}_{PML,u})] \cdot \begin{bmatrix} \frac{\partial \hat{\theta}_{PML,u}^R}{\partial \hat{\theta}_{PML,u}} \\ \frac{\partial \hat{\theta}_{PML,u}^R}{\partial \hat{\theta}_{PML,u}} \end{bmatrix}' \right\}^{-1} \cdot \hat{\theta}_{PML,u}^R,$$

where we adopt the approach of Binder (1983), adapted to the covariance structure models context, for calculating $a \text{cov}(\hat{\theta}_{PML,u})$. When $\sqrt{n^{-1}} \cdot (\hat{\theta}_{PML} - \theta)$ is asymptotically normal distributed, then the modified WT_c statistic introduced above is asymptotically $\chi_{b^*}^2$ distributed (Skinner, 1989) under H_0 . ■

We also propose modifying the overall model fit descriptive measures, such as the Jöreskog and Sörbom (1989) goodness of fit indices (GFI). Therefore, let

$$GFI_{c,PML} = 1 - \left\{ \frac{\text{tr}[(\Sigma(\hat{\theta})^{-1}S_w - I)^2]}{\text{tr}[(\Sigma(\hat{\theta})^{-1}S_w)^2]} \right\},$$

be a modified version GFI for complex survey data, when considering $F(\theta)_{PML}$. We also modify Tanaka and Huba (1985) GLS version of goodness of fit indices, so that

$$GFI_{c,GLSC}^1 = 1 - \left\{ \frac{\text{tr}[(I - \Sigma(\hat{\theta})S_w^{-1})^2]}{T} \right\}, \text{ or}$$

$$GFI_{c,GLSC}^2 = 1 - \left\{ \frac{\{vech[S_w] - vech[\Sigma(\hat{\theta})]\}' U^{-1} \{vech[S_w] - vech[\Sigma(\hat{\theta})]\}}{vech[S_w]' U^{-1} vech[S_w]} \right\}.$$

Moreover, modified adjusted fit indices (AGFI_c), could be calculated as

$$AGFI_c = 1 - \left(\frac{k}{df} \right) \cdot (1 - GFI_c).$$

5 Concluding remarks

Model testing is an important step in any model fit procedure. According to Menard (1991), in longitudinal models there is an increase to problems comparably to cross-sectional models in this regard. Eltinge (1999) acknowledged the need to improve and to develop new techniques for model assessment and diagnostic in the complex survey data context. Classic measures that are often used in model testing are appropriate for situations where data is obtained from a srs design.

We propose some new developments on model fitting statistics when working with longitudinal data in a complex survey design framework. Following Rao and Scott's conceptions (Rao and Scott, 1979), we propose modifying the Wald goodness of fit test in the context of models for covariance structures. Furthermore, we also propose a modification for the Wald significance test for nested hypothesis, following an approach suggested by Skinner (1989, Section 3.4). Goodness of fit indices proposed by Jöreskog and Sörbom (1989) are also modified in order be utilised in the complex survey data context.

Further research could involve evaluating the consequences of using standard model fit test techniques without consideration to the complexity of the sample, by extending the simulation study performed in Vieira and Skinner (2008).

Acknowledgement: The author thanks FAPEMIG for grant number CEX-APQ-02071-12 (Universal), and for supporting the attendance at the 59th WSC of the ISI.

References

- Binder, D. A. (1983) On the Variances of Asymptotically Normal Estimators from Complex Surveys. *International Statistical Review*, 51, 279-292.
- Bollen, K. A. (1989) *Structural Equations with Latent Variables*. New York, John Wiley & Sons.
- Chambers, R. L. and Skinner, C. J. eds. (2003) *Analysis of Survey Data*. Chichester, John Wiley & Sons.
- Diggle, P. J., Heagerty, P., Liang, K. & Zeger, S. L. (2002) *Analysis of Longitudinal Data*. 2nd ed. Oxford, Oxford University Press.
- Eltinge, J. L. (1999) *Assessment of Information Capacity and Sensitivity in the Analysis of Complex Surveys*. Bulletin of the International Statistical Institute. 52nd Session Proceedings.
- Jöreskog, K. G. and Sörbom, D. (1989) *Lisrel 7: A Guide to the Program and Applications*. Chicago, SPSS publications.
- Lindsey, J. K. (1994) *Models for Repeated Measurements*. Oxford, Oxford University Press.
- Rao, J. N. K. and Scott (1979) Chi-squared Tests for Analysis of Categorical Data from Complex Surveys. *Proceedings of the Survey Research Methods Section*, American Statistical Association.
- Satorra, A. (1989) Alternative Test Criteria in Covariance Structure Analysis: A Unified Approach. *Psychometrika*, Vol. 54, N. 1, p. 131-151.
- Skinner, C. J. (1989) Domain Means, Regression and Multivariate Analysis. In Skinner, C. J., Holt, D. and Smith, T. M. F. eds. *Analysis of Complex Surveys*. Chichester, John Wiley & Sons.
- Skinner, C. J. and Holmes, D. (2003) Random Effects Models for Longitudinal Survey Data. In Chambers, R. L. and Skinner, C. J. eds. *Analysis of Survey Data*. Chichester, John Wiley & Sons.
- Skinner, C. J., Holt, D. and Smith, T. M. F. eds. (1989) *Analysis of Complex Surveys*. Chichester, John Wiley & Sons.
- Skinner, C. J. and Vieira, M. D. T. (2007) Variance estimation in the analysis of clustered longitudinal survey data. *Survey Methodology*, Vol. 33, No. 1, pp. 3-12.
- Tanaka, J. S., and Huba, G. J. (1985) A fit index for covariance structure models under arbitrary GLS estimation. *British Journal of Mathematical and Statistical Psychology*, Vol. 38, 197-201.
- Vieira, M. D. T. and Skinner, C. J. (2008) Estimating Models for Panel Survey Data under Complex Sampling. *Journal of Official Statistics*, 24, 343-364.
- Vieira, M. D. T. (2009) *Analysis of Longitudinal Survey Data*. 1. ed. Saarbrücken: VDM Verlag Dr. Müller.