

THE STATA JOURNAL

Editor

H. Joseph Newton
Department of Statistics
Texas A & M University
College Station, Texas 77843
979-845-3142; FAX 979-845-3144
jnewton@stata-journal.com

Executive Editor

Nicholas J. Cox
Department of Geography
University of Durham
South Road
Durham City DH1 3LE UK
n.j.cox@stata-journal.com

Associate Editors

Christopher Baum
Boston College

Rino Bellocco
Karolinska Institutet

David Clayton
Cambridge Inst. for Medical Research

Mario A. Cleves
Univ. of Arkansas for Medical Sciences

William D. Dupont
Vanderbilt University

Charles Franklin
University of Wisconsin, Madison

Joanne M. Garrett
University of North Carolina

Allan Gregory
Queen's University

James Hardin
University of South Carolina

Stephen Jenkins
University of Essex

Ulrich Kohler
WZB, Berlin

Jens Lauritsen
Odense University Hospital

Stanley Lemeshow
Ohio State University

J. Scott Long
Indiana University

Thomas Lumley
University of Washington, Seattle

Roger Newson
King's College, London

Marcello Pagano
Harvard School of Public Health

Sophia Rabe-Hesketh
University of California, Berkeley

J. Patrick Royston
MRC Clinical Trials Unit, London

Philip Ryan
University of Adelaide

Mark E. Schaffer
Heriot-Watt University, Edinburgh

Jeroen Weesie
Utrecht University

Nicholas J. G. Winter
Cornell University

Jeffrey Wooldridge
Michigan State University

Stata Press Production Manager

Lisa Gilmore

Copyright Statement: The Stata Journal and the contents of the supporting files (programs, datasets, and help files) are copyright © by StataCorp LP. The contents of the supporting files (programs, datasets, and help files) may be copied or reproduced by any means whatsoever, in whole or in part, as long as any copy or reproduction includes attribution to both (1) the author and (2) the Stata Journal.

The articles appearing in the Stata Journal may be copied or reproduced as printed copies, in whole or in part, as long as any copy or reproduction includes attribution to both (1) the author and (2) the Stata Journal.

Written permission must be obtained from StataCorp if you wish to make electronic copies of the insertions. This precludes placing electronic copies of the Stata Journal, in whole or in part, on publicly accessible web sites, file servers, or other locations where the copy may be accessed by anyone other than the subscriber.

Users of any of the software, ideas, data, or other materials published in the Stata Journal or the supporting files understand that such use is made without warranty of any kind, by either the Stata Journal, the author, or StataCorp. In particular, there is no warranty of fitness of purpose or merchantability, nor for special, incidental, or consequential damages such as loss of profits. The purpose of the Stata Journal is to promote free communication among Stata users.

The *Stata Journal*, electronic version (ISSN 1536-8734) is a publication of Stata Press, and Stata is a registered trademark of StataCorp LP.

Confidence intervals for the variance component of random-effects linear models

Matteo Bottai
Arnold School of Public Health
University of South Carolina
800 Sumter Street
Columbia SC 29208

Nicola Orsini
Division of Nutritional Epidemiology
Institute of Environmental Medicine
Karolinska Institutet
Box 210, SE-171 77 Stockholm, Sweden

Abstract. We present the postestimation command `xtvc` to provide confidence intervals for the variance components of random-effects linear regression models. This command must be used after `xtreg` with option `mle`. Confidence intervals are based on the inversion of a score-based test (Bottai 2003).

Keywords: `st0077`, `xtvc`, variance components, confidence intervals, score test, random-effects linear models

1 Introduction

The random-effects linear model has been widely applied to different areas of data analysis (see, among many others, Breslow and Clayton 1993; Diggle, Liang, and Zeger 1994; Snijders and Bosker 1999; McCulloch and Searle 2001; Skrondal and Rabe-Hesketh 2004). The Stata `xtreg` command fits the random-effects linear regression model, which can be written as

$$y_{it} = \mathbf{x}_{it}\beta + u_i + e_{it}, \quad u_i \sim N(0, \sigma_u^2), e_{it} \sim N(0, \sigma_e^2) \quad (1)$$

where y_{it} is the t th observation taken on some random variable Y for the i th unit and $i = 1, \dots, m$, $t = 1, \dots, T_i$; \mathbf{x}_{it} is a covariate vector and β is a parameter vector of fixed effects; u_i is a unit-specific normal random effect with zero mean and variance σ_u^2 that is assumed to be non-negative; and e_{it} is the normal residual error with variance σ_e^2 that is assumed to be strictly positive. Also, u_i and e_{it} are assumed to be independent. Units can refer to individuals on whom repeated observations are taken, families whose members are sampled, or otherwise-defined groups within which observations may be correlated.

In such models, it is often of interest to make inference not only about the fixed and random effects but also about the variance components. In particular, testing homogeneity across units is equivalent to testing the null hypothesis

$$H_0: \sigma_u^2 = 0 \quad (2)$$

In general, testing whether a variance parameter is zero implies testing a parameter value on the boundary of the parameter space, the variance being non-negative. Several

authors suggest using the large-sample likelihood-ratio test that adjusts for the boundary condition. In fact, under this irregular scenario, the asymptotic distribution of the usual likelihood-ratio test statistic follows a distribution that is a 50:50 mixture of a $\chi^2_{(1)}$ and the constant zero (Self and Liang 1987). The Stata command `xtreg` provides the upper-tail probability of the appropriate asymptotic distribution of the likelihood-ratio test statistic (Gutierrez, Carter, and Drukker 2001).

However, such a method cannot be used to construct confidence intervals for the variance of the random effect, σ_u^2 . Besides, the confidence intervals provided for the random-effect variance by `xtreg`, based on a Wald-type test, can be shown to be asymptotically wrong. To the best of our knowledge, no published work has provided methods for constructing likelihood-based confidence regions for the variance component that are asymptotically correct.

It can be shown that inference about the variance component σ_u^2 can be accommodated within the irregular problems of singular information. Such a connection had been noted several years ago (Chesher 1984; Lee and Chesher 1986), but only recently a general theory was developed for the singular-information case (Rotnitzky et al. 2000). Using the results derived for the singular-information problem (Bottai 2003), a method is implemented in the Stata command `xtvc` that is based on the inversion of a score-type test, which provides asymptotically correct confidence intervals. Also, when testing the hypothesis of homogeneity across units (2), the proposed method is shown to have better small-sample properties than the one based on the likelihood-ratio test statistic.

The rest of the article is organized as follows: section 2 introduces the syntax of the command `xtvc`; section 3 provides an example in which the command `xtvc` is applied to real data; section 4 reports the observed rejection proportions of the confidence intervals generated by `xtvc` on simulated data; and some final remarks are presented in section 5.

2 The `xtvc` command

2.1 Syntax

The `xtvc` command is to be used after the `xtreg` command with the `mle` option for maximum likelihood estimation. The syntax of `xtvc` is as follows:

```
xtvc [ , llevel( # ) h0( # ) ]
```

2.2 Options

`level(#)` specifies the confidence level, as a percentage, for the confidence interval of the variance component. The default is `level(95)` or as set by `set level`; see [U] 23.6 Specifying the width of confidence intervals.

`h0(#)` performs the score-based test for the null hypothesis $H_0: \text{sigma}_u = \#$. The default null value is 0.

2.3 Saved Results

`xtvc` saves all the results of `xtreg` plus the following:

```
Scalars
  e(score)      score test statistic      e(pval)      p-value
  e(suuppbb)    upper bound of  $\sigma_u$   e(sulowb)    lower bound of  $\sigma_u$ 

Macros
  e(pcmd)       xtvc
```

3 Example: the NLSY data

`xtvc` is applied to the longitudinal data from a subsample of the NLSY data (Center for Human Resource Research 1989) described in many of the [XT] `xt` entries and available on the Stata Press web page (<http://www.stata-press.com/data/r8/xt/>). In this example, we fit a random-effects linear model for the variable `ln_wage` as a function of several variables as was done in the `xtreg` example; see [XT] `xtreg`.

```
. webuse nlswork, clear
(National Longitudinal Survey.  Young Women 14-26 years of age in 1968)
. iis idcode
. xtreg ln_w grade age ttl_exp tenure not_smsa south, mle
Fitting constant-only model:
(output omitted)
Fitting full model:
(output omitted)
Random-effects ML regression      Number of obs      =      28091
Group variable (i): idcode        Number of groups   =       4697
Random effects u_i ~ Gaussian     Obs per group: min =         1
                                   avg   =         6.0
                                   max   =         15
                                   LR chi2(6)      =      6861.27
Log likelihood = -9218.9773        Prob > chi2        =       0.0000
```

ln_wage	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
grade	.0691186	.0017232	40.11	0.000	.0657412 .072496
age	-.0033869	.0006491	-5.96	0.000	-.0051412 -.0025967
ttl_exp	.030151	.0011135	27.08	0.000	.0279687 .0323334
tenure	.013591	.0008454	16.08	0.000	.0119341 .0152478
not_smsa	-.1299789	.0071709	-18.13	0.000	-.1440337 -.1159242
south	-.0941264	.0071354	-13.19	0.000	-.1081115 -.0801413
_cons	.7566548	.0267764	28.26	0.000	.7041741 .8091355
/sigma_u	.2503043	.003531	70.89	0.000	.2433837 .2572249
/sigma_e	.2959207	.0013704	215.94	0.000	.2932348 .2986065
rho	.4170663	.0074739			.4024786 .4317692

Likelihood-ratio test of `sigma_u=0`: `chibar2(01)= 7277.75 Prob>=chibar2 = 0.000`

We then use the `xtvc` command:

```
. xtvc
```

ln_wage	ML Estimate	[95% Conf. Interval]	
/sigma_u	.2503043	.2488335	.2630834

```
Score test of sigma_u=0: chi2(1)= 39399.39 Prob>=chi2 = 0.000
```

The point estimate for the random-effects standard deviation σ_u is exactly the same as the one given by `xtreg`, but the confidence interval provided by `xtvc` is slightly shifted to include greater values. Both the score-type test provided by `xtvc` and the likelihood-ratio test provided by `xtreg` reject the null hypothesis that the standard deviation σ_u is equal to zero. With the `h0` option of the `xtvc` command, it is also possible to test any value for the standard deviation σ_u , not only zero. For example, we can test the value $\sigma_u = 0.25$, which is included in the 95% confidence interval.

```
. xtvc, h0(0.25)
```

ln_wage	ML Estimate	[95% Conf. Interval]	
/sigma_u	.2503043	.2488335	.2630834

```
Score test of sigma_u=0.25: chi2(1)= 2.63 Prob>=chi2 = 0.105
```

4 Simulated data

The `xtvc` command was applied to simulated data. Three thousand samples were pseudo-randomly generated for model (1) under a grid of values for the random-effect standard deviation $\sigma_u = 0, 0.01, \dots, 0.09, 0.10, 10$, and for different numbers of units or groups $m = 10, 100, 1000$. The residual-error standard deviation σ_e was set constant to the value one for all the simulations. Two covariates were pseudo-randomly generated from a `uniform(-1,1)` and a `uniform(0,2)` distribution, respectively, with $\beta = (1, 2)^T$. The observed rejection proportions over the simulated samples of the 95% confidence intervals provided by `xtvc` are shown in table 1. For the samples generated under the value $\sigma_u = 0$, the observed rejection proportion of the adjusted likelihood-ratio test at the 5% level provided by `xtreg` is also reported.

(Continued on next page)

Table 1: Observed rejection proportions of `xtvc` and `xtreg` (using `chibar2(01)`) among 3,000 simulated samples generated under different values of σ_u and number of units or groups for the random-effects linear model (1) (simulation error $\pm 0.78\%$).

σ_u	$m=10$	$m=100$	$m=1000$
xtvc			
0.00	5.20	5.23	4.63
0.01	5.17	5.43	5.37
0.02	5.03	5.23	4.93
0.03	5.33	5.60	4.57
0.04	5.30	5.07	5.63
0.05	4.73	5.63	5.00
0.06	5.77	5.17	4.93
0.07	5.30	5.63	5.30
0.08	5.27	5.40	4.53
0.09	5.47	5.43	5.30
0.10	4.80	5.20	4.07
10.0	4.57	5.03	4.90
xtreg			
0.00	2.43	4.13	4.27

Regardless of the number of units or groups, m , the observed rejection proportion is uniformly close to its nominal level of 5% across the values of the standard deviation σ_u . Although based on a large-sample test, `xtvc` shows acceptable behavior in small samples as well.

The adjusted likelihood-ratio test provided by `xtreg` was applied only to the samples simulated under the value $\sigma_u = 0$. In the present simulation, when the number of units or groups $m = 10$, its observed rejection proportion is 2.43%, well below its nominal level of 5%. In other extensive simulation experiments not reported here, we observed that the rejection proportion becomes satisfactorily close to the nominal level only when the number of units or groups is no smaller than a thousand.

The observed rejection proportion of the confidence regions obtained by inverting the Wald-type test, as provided by `xtreg`, is wrong in small samples as well as large samples. Depending on the values of σ_u and m , its rejection probability can be as high as 15% or as low as 0.5%. Besides, its confidence intervals may happen to include negative values, which are out of the feasible space of the variance parameter.

5 Final remarks

The `xtvc` command is the only solution for those seeking to construct confidence intervals for the variance component of a random-effects linear regression model. The method can be extended to more general models, such as generalized linear mixed mod-

els, whose estimation is based on the likelihood function. In the present version, the command `xtvc` only provides interval estimates when the number of units or groups is greater than eight. For balanced data, explicit solutions for the upper and lower bounds of the confidence intervals are available but are not implemented in the command `xtvc`. Instead, in the unbalanced case, the bounds of the confidence intervals are obtained by iterative algorithms. Equations are solved by bisection methods, which usually take little time to converge. In later versions, the Newton–Raphson optimization could be used instead, should the command take too long.

6 References

- Bottai, M. 2003. Confidence regions when the Fisher information is zero. *Biometrika* 90(1): 73–84.
- Breslow, N. E. and D. G. Clayton. 1993. Approximate inference in generalized linear mixed models. *Journal of the American Statistical Association* 88(421): 9–25.
- Center for Human Resource Research. 1989. National Longitudinal Survey of Labor Market Experience, Young Women 14–26 years of age in 1968. Ohio State University.
- Chesher, A. 1984. Testing for neglected heterogeneity. *Econometrica* 52(4): 865–872.
- Diggle, P. J., K. Y. Liang, and S. L. Zeger. 1994. *Analysis of Longitudinal Data*. Oxford: Oxford University Press.
- Gutierrez, R. G., S. Carter, and D. M. Drukker. 2001. sg160: On boundary-value likelihood-ratio tests. *Stata Technical Bulletin* 60: 15–18. In *Stata Technical Bulletin Reprints*, vol. 10, 269–273. College Station, TX: Stata Press.
- Lee, L. F. and A. Chesher. 1986. Specification testing when score test statistics are identically zero. *Journal of Econometrics* 31: 121–149.
- McCulloch, C. E. and S. R. Searle. 2001. *Generalized, Linear, and Mixed Models*. New York: Wiley.
- Rotnitzky, A., D. R. Cox, M. Bottai, and J. M. Robins. 2000. Likelihood-based asymptotic inference with singular information. *Bernoulli* 6(2): 243–284.
- Self, S. and K. Liang. 1987. Asymptotic properties of maximum likelihood estimators and likelihood-ratio test under nonstandard conditions. *Journal of the American Statistical Association* 82(398): 605–610.
- Skrondal, A. and S. Rabe-Hesketh. 2004. *Generalized Latent Variable Modeling: Multilevel, Longitudinal, and Structural Equation Models*. Boca Raton, FL: Chapman & Hall/CRC.
- Snijders, T. A. B. and R. J. Bosker. 1999. *Multilevel Analysis: An Introduction to Basic and Advance Multilevel Modeling*. London: Sage.

About the Authors

Matteo Bottai is Assistant Professor at the Arnold School of Public Health, University of South Carolina, Columbia, SC.

Nicola Orsini is a Ph.D. student at the Institute of Environmental Medicine, Division of Nutritional Epidemiology, Karolinska Institutet, Stockholm, Sweden.