



EXPEND, A GAUSS PROGRAMME FOR NON-LINEAR GMM  
ESTIMATION OF *EXPONENTIAL* MODELS WITH *ENDOGENOUS*  
REGRESSORS FOR CROSS SECTION AND PANEL DATA

---

*Frank Windmeijer*

THE INSTITUTE FOR FISCAL STUDIES  
DEPARTMENT OF ECONOMICS, UCL  
**cemmap** working paper CWP14/02

# *ExpEnd*, a Gauss Programme for Non-Linear GMM Estimation of *Exponential* Models with *Endogenous* Regressors for Cross Section and Panel Data\*

Frank Windmeijer  
Centre for Microdata Methods and Practice  
Institute for Fiscal Studies  
7 Ridgmount Street  
London WC1E 7AE  
f.windmeijer@ifs.org.uk

August 2002

## **Abstract**

*ExpEnd* is a Gauss programme for non-linear generalised method of moments (GMM) estimation of exponential models with endogenous regressors for cross section and panel data. The estimators included in this package are simple Poisson pseudo ML; GMM for cross section data using moment conditions based on multiplicative or additive errors; within groups fixed effects Poisson for panel data; GMM estimation using quasi-differenced moment conditions eliminating unobserved heterogeneity and allowing for predetermined or endogenous regressors; and quasi-differenced GMM for a dynamic linear feedback model. This manual describes in detail the various estimators, the data and software requirements, and the programme commands. The programme can be downloaded from <http://cemmap.ifs.org.uk/wps/expend.zip>.

**Key Words:** Generalised Method of Moments, Count Data, Panel data

**JEL Classification:** C13, C21, C23

---

\*I am grateful to Rachel Griffith and Marcos Vera-Hernandez for helpful comments. This programme has been developed while working on a project for the Centre for Economic Evaluation at IFS, a member of the ESRC Evidence Network. Financial support of the ESRC, grant no. H141251024, is gratefully acknowledged.

## 1. Introduction

*ExpEnd* is a Gauss programme for non-linear GMM estimation of *Exponential* models with *Endogenous* regressors for cross section and panel data. The estimators included in this package are simple Poisson pseudo ML; GMM for cross section data using moment conditions based on multiplicative or additive errors; within groups fixed effects Poisson for panel data; GMM estimation using quasi-differenced moment conditions eliminating unobserved heterogeneity and allowing for predetermined or endogenous regressors; and quasi-differenced GMM for a dynamic linear feedback model. The models and estimation methods are described in detail in Chamberlain (1992), Wooldridge (1991, 1997), Mullahy (1997), Windmeijer and Santos Silva (1997), Windmeijer (2000a) and Blundell, Griffith and Windmeijer (2002).

The programme can be downloaded from <http://cemmap.ifs.org.uk/wps/expend.zip>.

The EXPEND.ZIP file contains the following files:

EXPEND.PRG is the main programme file for use with MAXLIK 4.0

EXPENDOP.PRG is the main programme file for use with OPTMUM

EXPENDNM.PRG is the main programme file if MAXLIK 4.0 and OPTMUM are not available

EXPEND.RUN is the run file from which the main programme is called, using MAXLIK 4.0

EXPENDOP.RUN is the run file from which the main programme is called, using OPTMUM

EXPENDNM.RUN is the run file if MAXLIK 4.0 and OPTMUM are not available

GDATA.DAT and .DHT, an example gauss synthetic data set

AUXGDATA.DAT and .DHT, the auxiliary file accompanying the data set

Extract these files to a directory of choice. If the run file is in a different location from the programme file, the `#include` statement at the end of the run file should contain the path to the programme file. Estimation is done by editing the run file and executing it in Gauss.

## 2. Software Requirements

Gauss Version 3.2 for DOS, 3.5 and higher for Windows or 3.6 for UNIX, preferably with MAXLIK 4.0 or OPTMUM, although a simple optimisation routine is provided by the programme. If MAXLIK 4.0 is available, use the EXPEND.RUN

file, together with the EXPEND.PRG file. If OPTMUM is available use the EXPENDOP.RUN file, together with the EXPENDOP.PRG file. If MAXLIK 4.0 and OPTMUM are not available, use the EXPENDNM.RUN file, together with the EXPENDNM.PRG file.

### 3. Data Setup

The programme does not allow for missing values in any of the variables used. Further, when estimating a dynamic model, or when using instruments that are lags of the variables, there should be no "gaps" in the individual time series. The data set has to include an indicator of time, like year if the observations are annual. The programme does allow for unbalancedness of the data. In that case the data will have to be sorted in such a way that all individual observations with 1 time series observation come before all individuals with 2 time series observations etc. For example for three individuals, two with 2 and one with 3 observations observed in specific years, the data will have to be set up as follows:

$y$	$x$	$year$
$y_{1,84}$	$x_{1,84}$	1984
$y_{1,85}$	$x_{1,85}$	1985
$y_{2,84}$	$x_{2,84}$	1984
$y_{2,85}$	$x_{2,85}$	1985
$y_{3,83}$	$x_{3,83}$	1983
$y_{3,84}$	$x_{3,84}$	1984
$y_{3,85}$	$x_{3,85}$	1985

The auxiliary file contains two columns. The first column indicates the number of time series observations (1,2,3 etc.) and the second column indicates the number of individuals with these number of time series observations (500, 1000, 300, etc). The auxiliary file for the example above reads

$T$	$N$
2	2
3	1

For a cross section data set with 2500 observations the auxiliary data set has the entries 1 and 2500.

## 4. Models and Moment Conditions<sup>1</sup>

Let  $y_{it}$  denote the discrete count variable to be explained for subject  $i$ ,  $i = 1, \dots, N$ , at time  $t$ ,  $t = 1, \dots, T$ ; and let  $x_{it}$  denote a vector of explanatory variables. The exponential model

$$\begin{aligned} y_{it} &= \exp(x'_{it}\beta) + u_{it} \\ &= \mu_{it} + u_{it}, \end{aligned} \quad (4.1)$$

where  $\mu_{it} = \exp(x'_{it}\beta)$ , is commonly used for count data. If  $x_{it}$  is exogenous, such that  $E(u_{it}|x_{it}) = 0$ , then the moment estimator that solves

$$\frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T x_{it} (y_{it} - \mu_{it}) = 0, \quad (4.2)$$

is consistent and equivalent to the Poisson ML estimator. When  $x_{it}$  is endogenous, but there are valid instruments  $z_{it}$  available, then the GMM estimator for  $\beta$  that minimises

$$\left( \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T z_{it} (y_{it} - \mu_{it}) \right)' W_N \left( \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T z_{it} (y_{it} - \mu_{it}) \right) \quad (4.3)$$

is consistent. These moment conditions are referred to as additive moment conditions for the model in levels.

Multiplicative moment conditions for the model in levels when  $x_{it}$  is exogenous are given by

$$\frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T x_{it} \left( \frac{y_{it} - \exp(x'_{it}\beta)}{\exp(x'_{it}\beta)} \right) = 0, \quad (4.4)$$

the extension to endogenous  $x_{it}$  and instruments  $z_{it}$  is straightforward. These moment conditions were originally proposed by Mullahy (1997). A discussion of multiplicative versus additive moment conditions can be found in Windmeijer and Santos Silva (1997).

An important feature in panel data applications is unobserved heterogeneity or individual fixed effects. For count data models these effects are generally modelled multiplicatively as

$$\begin{aligned} y_{it} &= \exp(x'_{it}\beta + \eta_i) + u_{it} \\ &= \mu_{it}\nu_i + u_{it}, \end{aligned} \quad (4.5)$$

where  $\nu_i = \exp(\eta_i)$  is a permanent scaling factor for the individual specific mean. In general, it is likely that the unobserved fixed components  $\eta_i$  are correlated with the explanatory variables,  $E(x_{it}\eta_i) \neq 0$ , and therefore standard random effects estimators will be inconsistent.

---

<sup>1</sup>This section draws heavily from Section 2 in Blundell, Griffith and Windmeijer (2002).

#### 4.1. Strictly Exogenous Regressors

When the  $x_{it}$  are *strictly exogenous*, the conditional mean of  $y_{it}$  satisfies

$$E(y_{it} | \nu_i, x_{it}) = E(y_{it} | \nu_i, x_{i1}, \dots, x_{iT}). \quad (4.6)$$

For this case, Hausman, Hall and Griliches (1984) use the Poisson conditional maximum likelihood estimator (CMLE), conditioning on  $\sum_{t=1}^T y_{it}$ , which is the sufficient statistic for  $\eta_i$ . However, the Poisson maximum likelihood estimator (MLE) for  $\beta$  in a model with separate individual specific constants does not suffer from the incidental parameters problem, and is therefore consistent and the same as the CMLE, see Blundell, Griffith and Windmeijer (1997), and Lancaster (1997). The associated first order conditions imply that the Poisson MLE for  $\beta$  is equivalent to a moment estimator in a model where the ratio of individual, or within group, means are used to approximate the individual specific effects. The moment conditions for this within group mean scaling estimator are given by

$$\frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T x_{it} \left( y_{it} - \mu_{it} \frac{\bar{y}_i}{\mu_i} \right) = 0. \quad (4.7)$$

#### 4.2. Predetermined Regressors

A regressor is predetermined when it is not correlated with current and future shocks, but it is correlated with past shocks:

$$\begin{aligned} E(x_{it} u_{it+j}) &= 0, \quad j \geq 0 \\ E(x_{it} u_{it-s}) &\neq 0, \quad s \geq 1. \end{aligned}$$

With predetermined regressors, the within group mean scaling estimator is no longer consistent. Chamberlain (1992) has proposed transformations that eliminate the fixed effect from the multiplicative model and generate orthogonality conditions that can be used for consistent estimation in count data models with predetermined regressors. The quasi-differencing transformation is

$$s_{it} = y_{it} \frac{\mu_{it-1}}{\mu_{it}} - y_{it-1} = u_{it} \frac{\mu_{it-1}}{\mu_{it}} - u_{it-1}. \quad (4.8)$$

Let  $x_i^{t-1} = (x_{i1}, \dots, x_{it-1})$ . When  $x_{it}$  is predetermined, the following moment conditions hold:

$$E(s_{it} | x_i^{t-1}) = 0. \quad (4.9)$$

Wooldridge (1991) proposed the following quasi-differencing transformation

$$q_{it} = \frac{y_{it}}{\mu_{it}} - \frac{y_{it-1}}{\mu_{it-1}},$$

with moment conditions

$$E(q_{it}|x_i^{t-1}) = 0.$$

It is clear that a variable in  $x_{it}$  can not have only non-positive or non-negative values, as then the corresponding estimate for  $\beta$  is infinity. A way around this problem is to transform  $x$  in deviations from overall means, see Windmeijer (2000a).

### 4.3. Endogenous Regressors

Regressors are endogenous when  $E(x_{it}u_{it}) \neq 0$ . In this case, the Chamberlain transformation can not be used. Use of the Wooldridge transformation leads to the following moment conditions

$$E(q_{it}|x_i^{t-2}) = 0,$$

where  $x_i^{t-2} = (x_{i1}, \dots, x_{it-2})$ , see Windmeijer (2000a).

### 4.4. Linear Feedback Model

Blundell, Griffith and Windmeijer (2002) propose use of a linear feedback model for modelling dynamic count panel data process. The linear feedback model (LFM) of order  $p$  is defined as

$$\begin{aligned} y_{it} &= \sum_{j=1}^p \gamma_j y_{it-j} + \exp(x_{it}'\beta + \eta_i) + u_{it} \\ &= \sum_{j=1}^p \gamma_j y_{it-j} + \mu_{it}\nu_i + u_{it}, \end{aligned} \quad (4.10)$$

where lags of the dependent variable enter the model linearly. The LFM has its origins in the Integer-Valued Autoregressive (INAR) process

Even when the  $x_{it}$  are strictly exogenous, the within group mean scaling estimator will be inconsistent for small  $T$ , as the lagged dependent variable is a predetermined variable. For estimation by GMM, the Chamberlain quasi-differencing transformation for the LFM model is given by

$$s_{it} = (y_{it} - \sum_{j=1}^p \gamma_j y_{it-j}) \frac{\mu_{it-1}}{\mu_{it}} - \left( y_{it-1} - \sum_{j=1}^p \gamma_j y_{it-1-j} \right) \quad (4.11)$$

For predetermined  $x_{it}$  the following moment conditions hold

$$E(s_{it}|y_i^{t-2}, x_i^{t-1}) = 0. \quad (4.12)$$





where

$$C(\hat{\theta}_1) = \frac{1}{N} \sum_{i=1}^N \frac{\partial Z_i' s_i(\theta)}{\partial \theta} \Big|_{\hat{\theta}_1}.$$

The asymptotic variance of the efficient two-step GMM estimator is computed as

$$v\hat{a}r(\hat{\theta}_2) = \frac{1}{N} \left( C(\hat{\theta}_2)' W_N(\hat{\theta}_1) C(\hat{\theta}_2) \right)^{-1}. \quad (5.2)$$

Note that the two-step asymptotic standard errors can be severely biased downwards in small samples (small  $N$ ), see Windmeijer (2000b).

The Sargan test for overidentifying restrictions is given by

$$N \left( \frac{1}{N} \sum_{i=1}^N s_i(\hat{\theta}_2)' Z_i \right) W_N(\hat{\theta}_1) \left( \frac{1}{N} \sum_{i=1}^N Z_i' s_i(\hat{\theta}_2) \right), \quad (5.3)$$

which is asymptotically chi-squared distributed with  $k_Z - k_\theta$  degrees of freedom if the moment conditions are valid, where  $k_Z$  is the number of instruments and  $k_\theta$  is the number of parameters.

## 6. Programme Commands

The estimation of the various models is done via the run file. This run file has to be edited by the user and executed in Gauss to obtain the estimation results. EXPEND.RUN contains the following lines:

```

new;
cls;
speed = hsec

dataset = "gdata";      /*** name of data set          ***/
auxset = "auxgdata";    /*** name of auxiliary data set      ***/
                        /*** add paths if in different directory ***/

let yvar = y;           /*** name of dependent variable       ***/
let xvar = x x;        /*** names of explanatory variables   ***/
let lx = 0 1;          /*** lag lengths of explanatory variables ***/

model = 2;             /*** 0 = levels model                 ***/
                        /*** 1 = within group mean scaling    ***/
                        /*** 2 = quasi differenced model      ***/

addit = 1;            /*** 0 = multiplicative moment conditions ***/
                        /*** 1 = additive moment conditions    ***/
                        /*** only active when model=0          ***/

qdif = 1;             /*** 0 = Wooldridge moment conditions  ***/
                        /*** 1 = Chamberlain moment conditions ***/
                        /*** only active when model=2          ***/

devvar = 0;          /*** set to 1 for taking deviations from overall ***/
                        /*** means of explanatory variables     ***/
                        /*** use this for Wooldridge quasi-differencing ***/

lfmy = 1;            /*** 0 = no linear feedback model      ***/
lfmlag = 1;          /*** # of lags of dep. var.           ***/

seqz = 1             /*** 0 = no sequential instruments     ***/
let seqzvar = y x;   /*** names of sequential instruments   ***/
let lseqz1 = 2 1;    /*** lag length of seq. instrs. begin ***/
let lseqz2 = 4 3;    /*** lag length of seq. instrs. end   ***/

```

```

nonseqz = 1;          /*** 0 = no non-sequential instruments      ***/
let nseqzvar = q;    /*** names of non-sequential instruments      ***/
let lnseqz = 1;      /*** lag length of non-seq instruments      ***/

let timevar = year;  /*** name of time variable          ***/
                    /*** set equal to 0 for cross section analysis ***/

timedum = 0;        /*** 1 = time dummies in model          ***/
                    /*** estimated coefficients are d_(t)-d_(t-1)      ***/
                    /*** in quasi-differenced model          ***/
timez = 0;          /*** 1 = time dummies included in instrument set ***/
                    /*** when timedum=0          ***/

lagl = 1;           /*** maximum lag length in the model      ***/

llev = 0;           /*** observations with less than lagl+llev+1 ***/
                    /*** (lagl+llev+2 if model=2)          ***/
                    /*** time periods get discarded          ***/

nind = 20;          /*** no. of individual units processed in each read ***/

saveres = 1;        /*** 1 = save one- and two-step parameters and ***/
                    /*** variance matrices as b1, v1, b2 and v2.fmt      ***/

ml = 1;             /*** 1= uses MAXLIK 4 routine of GAUSS      ***/
                    /*** otherwise a simple method of scoring is used ***/
alg = 4;            /*** sets _max_Algorithm, eg. 4 = NEWTON    ***/
                    /*** active if ml=1          ***/

sval = 1;           /*** sval=0 sets all starting values to 0   ***/
let startvaly = 0.4; /*** start values for lagged dep. vars.     ***/
let startvalx = 0.5 0.2; /*** start values for expl. vars.         ***/
startvalc = 0;      /*** start value for constant              ***/
                    /*** start values for time dummies are 0      ***/

output file = expend.out on;

#include expend.prg; /*** add path if in different directory    ***/

```

Most definitions in the programme are self explanatory or explained by the comments above. Below follow some examples of models and estimation methods and the accompanying commands.

### 6.1. Cross Section Estimation

For estimation using cross section data, only the model in levels can be estimated, so  $model = 0$ . All the lag lengths of the explanatory variables have to be set to zero,  $lx = 0...0$ , and  $lfmlag = 0$ . For additive moment conditions as in (4.2) set  $addit = 1$ , for multiplicative moment conditions as in (4.4) set  $addit = 0$ . As there are no sequential instruments, set  $seqz = 0$  and  $nonseqz = 1$ . As there are no time effects,  $timevar = timedum = timez = 0$ . If the instruments are the same as the explanatory variables, the Poisson pseudo ML estimates will be returned (i.e. with robust standard errors) if  $addit = 1$ .

### 6.2. Within Group Mean Scaling Estimation

For within group mean scaling estimation, using the moment conditions as in (4.7) set  $model = 1$ . The Poisson fixed effects pseudo ML estimation results are obtained when only non-sequential instruments are used that are identical to the explanatory variables. Within-groups with the linear feedback model is also possible, but this is of course no longer identical to the Poisson fixed effects estimator, and biased and inconsistent for small  $T$ .

### 6.3. Quasi-Differencing

For quasi-differencing, set  $model = 2$ . When  $x_{it}$  is endogenous, set  $qdif = 0$  for the Wooldridge moment conditions. It is recommended to set  $devvar = 1$  when using the Wooldridge moment conditions, so none of the variables take only non-positive or only non-negative values. Using non-sequential moment conditions, set  $nonseqz = 1$ . Using sequential moment conditions, set  $seqz = 1$ . When the time dimension is large it is recommended not to use all lagged values as sequential instruments. This is controlled by setting  $lseqz1$  and  $lseqz2$ , for example setting  $lseqz1 = 2$  and  $lseqz2 = 4$  results in the following instrument matrix when



---

CHAMBERLAIN MOMENT CONDITIONS

DATASET	gdata		
DATE	2/04/02		
N	1500	NT	7000
LAGL	1	LLEV	0
PERIOD	1987	1991	
DEP. VAR.	Y		
INSTRUMENTS	Y_2_4	X_1_3	
	Q_1		
SARGAN DOF P:	29.4340	24	0.2042

---

ONE-STEP

# ITERATIONS 10

	coeff	rob se	t-ratio	p-value
Y_1	0.3982	0.0908	4.3849	0.0000
X_	0.3414	0.1150	2.9681	0.0030
X_1	0.1300	0.1091	1.1913	0.2335

TESTS FOR SERIAL CORRELATION AND P-VALUES

M1	-5.1561	0.0000
M2	0.1652	0.8688

TWO-STEP

# ITERATIONS 5

	coeff	rob se	t-ratio	p-value
Y_1	0.3184	0.0435	7.3180	0.0000
X_	0.3708	0.0376	9.8642	0.0000
X_1	0.2287	0.0264	8.6585	0.0000

TESTS FOR SERIAL CORRELATION AND P-VALUES

M1	-7.3128	0.0000
M2	-0.5290	0.5968

Execution time is 0.32 minutes

The DATE format is month/day/year.

$N$  is the number of individual units,  $NT$  is the total number of observations  
Sequential instruments have a double subscript, non-sequential instruments  
have a single subscript.

Sargan is the test for overidentifying restrictions as in (5.3), based on the  
two-step results.

The standard errors for the one-step estimation results are calculated from the  
robust asymptotic variance (5.1).

The two-step standard errors are calculated from the robust asymptotic vari-  
ance (5.2).

The tests for serial correlation test for serial correlation of the residuals  $s_{it}$  as in  
(4.11) in this case and is an extension of the tests for serial correlation in Arellano  
and Bond (1991). Under the null of no serial correlation, these test statistics  
are asymptotically  $N(0, 1)$  distributed. For the Chamberlain and Wooldridge  
residuals one expects first-order, but not second-order autocorrelation if the model  
is well specified.

## References

- [1] Arellano, M., and S. Bond (1991), Some Tests of Specification for Panel Data: Monte Carlo Evidence and an Application to Employment Equations, *Review of Economic Studies*, 58, 277-98.
- [2] Blundell, R., R. Griffith and F. Windmeijer (1997), Individual Effects and Dynamics in Count Data Models, mimeo, Institute for Fiscal Studies.
- [3] Blundell, R., R. Griffith and F. Windmeijer (2002), Individual Effects and Dynamics in Count Data Models, *Journal of Econometrics* 108, 113-131.
- [4] Hansen, L.P. (1982), Large Sample Properties of Generalized Method of Moments Estimators, *Econometrica*, 50, 1029-1054.
- [5] Hausman, J., B. Hall, and Z. Griliches (1984), Econometric Models for Count Data and an Application to the Patents-R&D Relationship, *Econometrica*, 52, 909-938.
- [6] Lancaster, T. (1997), Orthogonal Parameters and Panel Data, mimeo, Brown University.
- [7] Mullahy, J. (1997), Instrumental Variable Estimation of Poisson Regression Models: Application to Models of Cigarette Smoking Behavior, *Review of Economics and Statistics* 79, 586-593.

- [8] Windmeijer, F., and J.M.C. Santos Silva (1997), Endogeneity in Count Data Models: An Application to Demand for Health Care, *Journal of Applied Econometrics* 12, 281-294.
- [9] Windmeijer, F. (2000a), Moment Conditions for Fixed Effects Count Data Models with Endogenous Regressors, *Economics Letters* 68, 21-24.
- [10] Windmeijer, F. (2000b), A Finite Sample Correction for the Variance of Linear Two-Step GMM Estimators, Institute for Fiscal Studies Working Paper Series No. W00/19, London, <http://www.ifs.org.uk/workingpapers/wp0019.pdf>.
- [11] Wooldridge, J.M. (1991), Multiplicative Panel Data Models without the Strict Exogeneity Assumption, MIT Working Paper No. 574.
- [12] Wooldridge, J.M. (1997), Multiplicative Panel Data Models without the Strict Exogeneity Assumption, *Econometric Theory*, 13, 667-678.