Constrained Maximum Likelihood Estimation

Information in this document is subject to change without notice and does not represent a commitment on the part of Aptech Systems, Inc. The software described in this document is furnished under a license agreement or nondisclosure agreement. The software may be used or copied only in accordance with the terms of the agreement. The purchaser may make one copy of the software for backup purposes. No part of this manual may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying and recording, for any purpose other than the purchaser's personal use without the written permission of Aptech Systems, Inc. ©Copyright 1994-2000 by Aptech Systems, Inc., Maple Valley, WA. All Rights Reserved.

GAUSS, GAUSS Engine, GAUSS Light are trademarks of Aptech Systems, Inc. All other trademarks are the properties of their respective owners.

Documentation Version: August 16, 2001

Part Number: 001860

Contents

1	Inst	allation 1				
	1.1	UNIX	1			
		1.1.1 Download	1			
		1.1.2 Floppy	1			
		1.1.3 Solaris 2.x Volume Management	2			
	1.2	Windows/NT/2000	3			
		1.2.1 Download	3			
		1.2.2 Floppy	3			
	1.3	Differences Between the UNIX and Windows/NT/2000 Versions \dots .	3			
2	Con	trained Maximum Likelihood Estimation	5			
	2.1	Getting Started	5			
		2.1.1 README Files	5			
		2.1.2 Setup	5			
		2.1.3 Converting MAXLIK Command Files	6			
	2.2	The Log-likelihood Function	7			
	2 3	Algorithm	Q			

	2.3.1	Derivatives	9
	2.3.2	The Secant Algorithms	10
	2.3.3	Line Search Methods	10
	2.3.4	Weighted Maximum Likelihood	12
	2.3.5	Active and Inactive Parameters	12
2.4	Manag	ging Optimization	12
	2.4.1	Scaling	13
	2.4.2	Condition	13
	2.4.3	Starting Point	13
	2.4.4	Diagnosis	13
2.5	Const	raints	14
	2.5.1	Linear Equality Constraints	14
	2.5.2	Linear Inequality Constraints	15
	2.5.3	Nonlinear Equality	15
	2.5.4	Nonlinear Inequality	15
	2.5.5	Bounds	16
	2.5.6	Example	16
2.6	Gradie	ents	18
	2.6.1	Analytical Gradient	18
	2.6.2	User-Supplied Numerical Gradient	19
	2.6.3	Analytical Hessian	20
	2.6.4	User-Supplied Numerical Hessian	22
	2.6.5	Analytical Nonlinear Constraint Jacobians	22
2.7	Inforo	neo	93

		2.7.1	Covariance Matrix of the Parameters	24
		2.7.2	Testing Constraints	26
		2.7.3	Heteroskedastic-consistent Covariance Matrix	26
		2.7.4	Confidence Limits	27
		2.7.5	Bootstrap	28
		2.7.6	Profiling	29
	2.8	Run-T	Time Switches	30
	2.9	Error	Handling	31
		2.9.1	Return Codes	31
		2.9.2	Error Trapping	32
	2.10	Refere	ences	32
3	Con	straine	d Maximum Likelihood Reference	35
3				
3	CM	L		36
3	CM	L LBlimit	ts	36 51
3	CM: CM:	L LBlimit LClimit	ts	36 51 52
3	CM CM CM	L LBlimit LClimit LTlimit	ts	36 51 52 56
3	CM CM CM CM	L LBlimit LClimit LTlimit LBoot	ts	36 51 52 56 57
3	CMCCMCCMCCMCCMCCMCCMCCMCCMCCMCCMCCMCCMC	L LBlimit LClimit LTlimit LBoot LDensit	ts	36 51 52 56 57
3	CMCCMCCMCCMCCMCCMCCMCCMCCMCCMCCMCCMCCMC	L LBlimit LClimit LTlimit LBoot LDensit LHist	ts	36 51 52 56 57 59 61
3	CMCCMCCMCCMCCMCCMCCMCCMCCMCCMCCMCCMCCMC	L LBlimit LClimit LTlimit LBoot LDensit LHist LProfile	ts	36 51 52 56 57 59 61 63
3	CMCCMCCMCCMCCMCCMCCMCCMCCMCCMCCMCCMCCMC	L LBlimit LClimit LTlimit LBoot LDensit LHist LProfile LSet .	ts	36 51 52 56 57 59 61 63 66
3	CMCCMCCMCCMCCMCCMCCMCCMCCMCCMCCMCCMCCMC	L LBlimit LClimit LTlimit LBoot LDensit LHist LProfile LSet . LPrt .	ts	36 51 52 56 57 59 61 63

4	Con	straine	d Event Count and Duration Regression	71
	4.1	Gettir	ng Started	72
		4.1.1	README Files	72
		4.1.2	Setup	72
	4.2	About	the CONSTRAINED COUNT Procedures	73
		4.2.1	Inputs	74
		4.2.2	Outputs	75
		4.2.3	Global Control Variables	75
		4.2.4	Adding Constraints	78
		4.2.5	Statistical Inference	78
		4.2.6	Problems with Convergence	80
	4.3	Annot	ated Bibliography	82
5	CM	LCount	Reference	85
	CM	LCount	Prt	86
	CM	LCount	Prt	87
	CM	LCount	Set	88
	CM	LExpga	um	89
	CM	LExpor	1	94
	CM	LHurdle	ep	98
	CM	LNegbi	n	102
	CM	LPareto)	108
	CM	LPoisso	n	113
	CM	LSuprei	me	118
	CI I	r Cunno	$\mathrm{me}2$	199

Chapter 1

Installation

1.1 UNIX

If you are unfamiliar with UNIX, see your system administrator or system documentation for information on the system commands referred to below. The device names given are probably correct for your system.

1.1.1 Download

- 1. Copy the .tar.gz file to /tmp.
- 2. Unzip the file.

```
gunzip appxxx.tar.gz
```

3. cd to the **GAUSS** or **GAUSS Engine** installation directory. We are assuming /usr/local/gauss in this case.

```
cd /usr/local/gauss
```

4. Untar the file.

```
tar xvf /tmp/appxxx.tar
```

1.1.2 Floppy

1. Make a temporary directory.

```
mkdir /tmp/workdir
```

2. cd to the temporary directory.

cd /tmp/workdir

3. Use tar to extract the files.

tar xvf $device_name$

If this software came on diskettes, repeat the tar command for each diskette.

4. Read the README file.

more README

5. Run the install.sh script in the work directory.

./install.sh

The directory the files are install to should be the same as the install directory of **GAUSS** or the **GAUSS Engine**.

6. Remove the temporary directory (optional).

The following device names are suggestions. See your system administrator. If you are using Solaris 2.x, see Section 1.1.3.

Operating System	3.5-inch diskette	1/4-inch tape	DAT tape
Solaris 1.x SPARC	/dev/rfd0	/dev/rst8	
Solaris 2.x SPARC	/dev/rfd0a (vol. mgt. off)	/dev/rst12	/dev/rmt/1l
Solaris 2.x SPARC	/vol/dev/aliases/floppy0	/dev/rst12	/dev/rmt/1l
Solaris 2.x x86	/dev/rfd0c (vol. mgt. off)		/dev/rmt/1l
Solaris 2.x x86	/vol/dev/aliases/floppy0		/dev/rmt/1l
HP-UX	/dev/rfloppy/c20Ad1s0		/dev/rmt/0m
IBM AIX	/dev/rfd0	/dev/rmt.0	
SGI IRIX	/dev/rdsk/fds0d2.3.5hi		

1.1.3 Solaris 2.x Volume Management

If Solaris 2.x volume management is running, insert the floppy disk and type

volcheck

to signal the system to mount the floppy.

The floppy device names for Solaris 2.x change when the volume manager is turned off and on. To turn off volume management, become the superuser and type

/etc/init.d/volmgt off

To turn on volume management, become the superuser and type

/etc/init.d/volmgt on

1. INSTALLATION

$1.2 \quad \text{Windows/NT/2000}$

1.2.1 Download

Unzip the .zip file into the GAUSS or GAUSS Engine installation directory.

1.2.2 Floppy

- 1. Place the diskette in a floppy drive.
- 2. Call up a DOS window
- 3. In the DOS window log onto the root directory of the diskette drive. For example:

A:<enter>
cd\<enter>

4. Type: ginstall source_drive target_path

source_drive Drive containing files to install

with colon included

For example: **A:**

target_path Main drive and subdirectory to install

to without a final \

For example: C:\GAUSS

A directory structure will be created if it does not already exist and the files will be copied over.

 $target_path \$ source code files $target_path \$ library files $target_path \$ examples example files

1.3 Differences Between the UNIX and Windows/NT/2000 Versions

• If the functions can be controlled during execution by entering keystrokes from the keyboard, it may be necessary to press *Enter* after the keystroke in the UNIX version.

• On the Intel math coprocessors used by the Windows/NT/2000 machines, intermediate calculations have 80-bit precision, while on the current UNIX machines, all calculations are in 64-bit precision. For this reason, **GAUSS** programs executed under UNIX may produce slightly different results, due to differences in roundoff, from those executed under Windows/NT/2000.

Chapter 2

Constrained Maximum Likelihood Estimation

written by

Ronald Schoenberg

This module contains a set of procedures for the solution of the constrained maximum likelihood problem

2.1 Getting Started

GAUSS 3.2.8+ is required to use these routines.

2.1.1 README Files

The file **README.cml** contains any last minute information on this module. Please read it before using the procedures in this module.

2.1.2 Setup

In order to use the procedures in the *CONSTRAINED MAXIMUM LIKELIHOOD* Module, the **CML** library must be active. This is done by including cml in the **LIBRARY** statement at the top of your program or command file:

library cml,pgraph;

This enables **GAUSS** to find the *CONSTRAINED MAXIMUM LIKELIHOOD* procedures. If you plan to make any right hand references to the global variables (described in the *REFERENCE* section), you also need the statement:

```
#include cml.ext;
```

Finally, to reset global variables in succeeding executions of the command file the following instruction can be used:

```
cmlset;
```

This could be included with the above statements without harm and would insure the proper definition of the global variables for all executions of the command file.

The version number of each module is stored in a global variable:

_cml_version 3×1 matrix, the first element contains the major version number of the $CONSTRAINED\ MAXIMUM\ LIKELIHOOD\ Module$, the second element the minor version number, and the third element the revision number.

If you call for technical support, you may be asked for the version number of your copy of this module.

2.1.3 Converting MAXLIK Command Files

The **CML** module includes a utility for processing command files to change **MAXLIK** global names to **CML** global names and vice versa. This utility is a standalone executable program that is called outside of **GAUSS**. The format is:

chgvar control_file target_directory file...

The *control_file* is an ASCII file containing a list of the symbols to change in the first column and the new symbol names in the second column. The **CML** module comes with three control_files:

cmltoml4 CML to MAXLIK 4.x ml4tocml MAXLIK 4.x to CML ml3tocml MAXLIK 3.x to CML

CHGVAR processes each file and writes a new file with the same name in the target directory.

A common use for **CHGVAR** is translating a command file that had been used before with **MAXLIK 3.x** to one that can be run with **CML**. For example:

mkdir new chgvar ml3tocml new max*.cmd

This would convert every file matching max*.cmd in the current directory and create a new file with the same name in the new directory.

The reverse translation is also possible. However, there are many global names in **CML** that don't have a corresponding global in **MAXLIK**, and in these cases no translation occurs.

Further editing of the file may be necessary after processing by CHGVAR.

You may edit the control files or create your own. They are ASCII files with each line containing a pair of names, the first column being the old name, and the second column the new name.

2.2 The Log-likelihood Function

CONSTRAINED MAXIMUM LIKELIHOOD is a set of procedures for the estimation of the parameters of models via the maximum likelihood method with general constraints on the parameters, along with an additional set of procedures for statistical inference.

 $CONSTRAINED\ MAXIMUM\ LIKELIHOOD\ solves$ the general weighted maximum likelihood problem

$$L = \sum_{i=1}^{N} \log P(Y_i; \theta)^{w_i}$$

where N is the number of observations, w_i is a weight. $P(Y_i, \theta)$ is the probability of Y_i given θ , a vector of parameters, subject to the linear constraints,

 $A\theta = B$

 $C\theta \ge D$

the nonlinear constraints

 $G(\theta) = 0$

 $H(\theta) \ge 0$

and bounds

$$\theta_l < \theta < \theta_u$$

 $G(\theta)$ and $H(\theta)$ are functions provided by the user and must be differentiable at least once with respect to θ .

The CONSTRAINED MAXIMUM LIKELIHOOD procedure CML finds values for the parameters in θ such that L is maximized. In fact CML minimizes -L. It is important to note, however, that the user must specify the log-probability to be maximized. CML transforms the function into the form to be minimized.

CML has been designed to make the specification of the function and the handling of the data convenient. The user supplies a procedure that computes $\log P(Y_i; \theta)$, i.e., the log-likelihood, given the parameters in θ , for either an individual observation or set of observations (i.e., it must return either the log-likelihood for an individual observation or a vector of log-likelihoods for a matrix of observations; see discussion of the global variable **___row** below). **CML** uses this procedure to construct the function to be minimized.

2.3 Algorithm

CONSTRAINED MAXIMUM LIKELIHOOD uses the Sequential Quadratic Programming method. In this method the parameters are updated in a series of iterations beginning with a starting values that you provide. Let θ_t be the current parameter values. Then the succeeding values are

$$\theta_{t+1} = \theta_t + \rho \delta$$

where δ is a $K \times 1$ direction vector, and ρ a scalar step length.

DIRECTION

Define

$$\Sigma(\theta) = \frac{\partial^2 L}{\partial \theta \partial \theta'}$$

$$\Psi(\theta) = \frac{\partial L}{\partial \theta}$$

and the Jacobians

$$\dot{G}(\theta) = \frac{\partial G(\theta)}{\partial \theta}$$
$$\dot{H}(\theta) = \frac{\partial H(\theta)}{\partial \theta}$$

For the purposes of this exposition, and without loss of generality, we may assume that the linear constraints and bounds have been incorporated into G and H.

The direction, δ is the solution to the quadratic program

$$minimize \ \frac{1}{2}\delta'\Sigma(\theta_t)\delta + \Psi(\theta_t)\delta$$

subject to
$$\dot{G}(\theta_t)\delta + G(\theta_t) = 0$$

 $\dot{H}(\theta_t)\delta + H(\theta_t) \ge 0$

This solution requires that Σ be positive semi-definite.

In practice, linear constraints are specified separately from the G and H because their Jacobians are known and easy to compute. And the bounds are more easily handled separately from the linear inequality constraints.

LINE SEARCH

Define the merit function

$$m(\theta) = L + \max \mid \kappa \mid \sum_{j} \mid g_{j}(\theta) \mid -\max \mid \lambda \mid \sum_{\ell} \min(0, h_{\ell}(\theta))$$

where g_j is the j-th row of G, h_ℓ is the ℓ -th row of H, κ is the vector of Lagrangean coefficients of the equality constraints, and λ the Lagrangean coefficients of the inequality constraints.

The line search finds a value of ρ that minimizes or decreases $m(\theta_t + \rho \delta)$.

2.3.1 Derivatives

The SQP method requires the calculation of a Hessian, Σ , and various gradients and Jacobians, Ψ , $\dot{G}(\theta)$, and $\dot{H}(\theta)$. **CML** computes these numerically if procedures to compute them are not supplied.

If you provide a proc for computing Ψ , the first derivative of L, **CML** uses it in computing Σ , the second derivative of L, i.e., Σ is computed as the Jacobian of the gradient. This improves the computational precision of the Hessian by about four places. The accuracy of the gradient is improved and thus the iterations converge in fewer iterations. Moreover, the convergence takes less time because of a decrease in function calls - the numerical gradient requires k function calls while an analytical gradient reduces that to one.

2.3.2 The Secant Algorithms

The Hessian may be very expensive to compute at every iteration, and poor start values may produce an ill-conditioned Hessian. For these reasons alternative algorithms are provided in **CML** for updating the Hessian rather than computing it directly at each iteration. These algorithms, as well as step length methods, may be modified during the execution of **CML**.

Beginning with an initial estimate of the Hessian, or a conformable identity matrix, an update is calculated. The update at each iteration adds more "information" to the estimate of the Hessian, improving its ability to project the direction of the descent. Thus after several iterations the secant algorithm should do nearly as well as Newton iteration with much less computation.

There are two basic types of secant methods, the BFGS (Broyden, Fletcher, Goldfarb, and Shanno), and the DFP (Davidon, Fletcher, and Powell). They are both rank two updates, that is, they are analogous to adding two rows of new data to a previously computed moment matrix. The Cholesky factorization of the estimate of the Hessian is updated using the functions **CHOLUP** and **CHOLDN**.

Secant Methods (BFGS and DFP)

BFGS is the method of Broyden, Fletcher, Goldfarb, and Shanno, and DFP is the method of Davidon, Fletcher, and Powell. These methods are complementary (Luenberger 1984, page 268). BFGS and DFP are like the NEWTON method in that they use both first and second derivative information. However, in DFP and BFGS the Hessian is approximated, reducing considerably the computational requirements. Because they do not explicitly calculate the second derivatives they are sometimes called *quasi-Newton* methods. While it takes more iterations than the NEWTON method, the use of an approximation produces a gain because it can be expected to converge in less overall time (unless analytical second derivatives are available in which case it might be a toss-up).

The secant methods are commonly implemented as updates of the *inverse* of the Hessian. This is not the best method numerically for the BFGS algorithm (Gill and Murray, 1972). This version of **CML**, following Gill and Murray (1972), updates the Cholesky factorization of the Hessian instead, using the functions **CHOLUP** and **CHOLDN** for BFGS. The new direction is then computed using **CHOLSOL**, a Cholesky solve, as applied to the updated Cholesky factorization of the Hessian and the gradient.

2.3.3 Line Search Methods

Given a direction vector d, the updated estimate of the parameters is computed

$$\theta_{t+1} = \theta_t + \rho \delta$$

where ρ is a constant, usually called the *step length*, that increases the descent of the function given the direction. **CML** includes a variety of methods for computing ρ . The value of the function to be minimized as a function of ρ is

$$m(\theta_t + \rho \delta)$$

Given θ and d, this is a function of a single variable ρ . Line search methods attempt to find a value for ρ that decreases m. STEPBT is a polynomial fitting method, BRENT and HALF are iterative search methods. A fourth method called ONE forces a step length of 1. The default line search method is STEPBT. If this, or any selected method, fails, then BRENT is tried. If BRENT fails, then HALF is tried. If all of the line search methods fail, then a random search is tried (provided **_cml_RandRadius** is greater than zero).

STEPBT

STEPBT is an implementation of a similarly named algorithm described in Dennis and Schnabel (1983). It first attempts to fit a quadratic function to $m(\theta_t + \rho \delta)$ and computes an ρ that minimizes the quadratic. If that fails it attempts to fit a cubic function. The cubic function more accurately portrays the F which is not likely to be very quadratic, but is, however, more costly to compute. STEPBT is the default line search method because it generally produces the best results for the least cost in computational resources.

BRENT

This method is a variation on the golden section method due to Brent (1972). In this method, the function is evaluated at a sequence of test values for ρ . These test values are determined by extrapolation and interpolation using the constant, $(\sqrt{5}-1)/2=.6180...$ This constant is the inverse of the so-called "golden ratio" $((\sqrt{5}+1)/2=1.6180...$ and is why the method is called a golden section method. This method is generally more efficient than STEPBT but requires significantly more function evaluations.

HALF

This method first computes m(x+d), i.e., sets $\rho=1$. If m(x+d) < m(x) then the step length is set to 1. If not, then it tries m(x+.5d). The attempted step length is divided by one half each time the function fails to decrease, and exits with the current value when it does decrease. This method usually requires the fewest function evaluations (it often only requires one), but it is the least efficient in that it is not very likely to find the step length that decreases m the most.

BHHHSTEP

This is a variation on the golden search method. A sequence of step lengths are computed, interpolating or extrapolating using a golden ratio, and the method exits when the function decreases by an amount determined by **_cml_Interp**.

2.3.4 Weighted Maximum Likelihood

Weights are specified by setting the **GAUSS** global, **__weight** to a weighting vector, or by assigning it the name of a column in the **GAUSS** data set being used in the estimation. Thus if a data matrix is being analyzed, **__weight** must be assigned to a vector.

CML assumes that the weights sum to the number of observations, i.e, that the weights are frequencies. This will be an issue only with statistical inference. Otherwise, any multiple of the weights will produce the same results.

2.3.5 Active and Inactive Parameters

The **CML** global **_cml_Active** may be used to fix parameters to their start values. This allows estimation of different models without having to modify the function procedure. **_cml_Active** must be set to a vector of the same length as the vector of start values. Elements of **_cml_Active** set to zero will be fixed to their starting values, while nonzero elements will be estimated.

This feature may also be used for model testing. **__cml_NumObs** times the difference between the function values (the second return argument in the call to **CML**) is chi-squared distributed with degrees of freedom equal to the number of fixed parameters in **__cml_Active**.

2.4 Managing Optimization

The critical elements in optimization are scaling, starting point, and the condition of the model. When the data are scaled, the starting point is reasonably close to the solution, and the data and model go together well, the iterations converge quickly and without difficulty.

For best results therefore, you want to prepare the problem so that model is well-specified, the data scaled, and that a good starting point is available.

The tradeoff among algorithms and step length methods is between speed and demands on the starting point and condition of the model. The less demanding methods are generally time consuming and computationally intensive, whereas the quicker methods (either in terms of time or number of iterations to convergence) are more sensitive to conditioning and quality of starting point.

2.4.1 Scaling

For best performance, the diagonal elements of the Hessian matrix should be roughly equal. If some diagonal elements contain numbers that are very large and/or very small with respect to the others, **CML** has difficulty converging. How to scale the diagonal elements of the Hessian may not be obvious, but it may suffice to ensure that the constants (or "data") used in the model are about the same magnitude.

2.4.2 Condition

The specification of the model can be measured by the condition of the Hessian. The solution of the problem is found by searching for parameter values for which the gradient is zero. If, however, the Jacobian of the gradient (i.e., the Hessian) is very small for a particular parameter, then **CML** has difficulty determining the optimal values since a large region of the function appears virtually flat to **CML**. When the Hessian has very small elements, the inverse of the Hessian has very large elements and the search direction gets buried in the large numbers.

Poor condition can be caused by bad scaling. It can also be caused by a poor specification of the model or by bad data. Bad models and bad data are two sides of the same coin. If the problem is highly nonlinear, it is important that data be available to describe the features of the curve described by each of the parameters. For example, one of the parameters of the Weibull function describes the shape of the curve as it approaches the upper asymptote. If data are not available on that portion of the curve, then that parameter is poorly estimated. The gradient of the function with respect to that parameter is very flat, elements of the Hessian associated with that parameter is very small, and the inverse of the Hessian contains very large numbers. In this case it is necessary to respecify the model in a way that excludes that parameter.

2.4.3 Starting Point

When the model is not particularly well-defined, the starting point can be critical. When the optimization doesn't seem to be working, try different starting points. A closed form solution may exist for a simpler problem with the same parameters. For example, ordinary least squares estimates may be used for nonlinear least squares problems or nonlinear regressions like probit or logit. There are no general methods for computing start values and it may be necessary to attempt the estimation from a variety of starting points.

2.4.4 Diagnosis

When the optimization is not proceeding well, it is sometimes useful to examine the function, the gradient Ψ , the direction δ , the Hessian Σ , the parameters θ_t , or the step

length ρ , during the iterations. The current values of these matrices can be printed out or stored in the global **_cml_Diagnostic** by setting **_cml_Diagnostic** to a nonzero value. Setting it to 1 causes **CML** to print them to the screen or output file, 2 causes **CML** to store then in **_cml_Diagnostic**, and 3 does both.

When you have selected **_cml_Diagnostic** = 2 or 3, **CML** inserts the matrices into **_cml_Diagnostic** using the **VPUT** command. The matrices are extracted using the **VREAD** command. For example,

```
_cml_Diagnostic = 2;
call CMLPrt(CML("tobit",0,&lpr,x0));
h = vread(_cml_Diagnostic,"hessian");
d = vread(_cml_Diagnostic,"direct");
```

The following table contains the strings to be used to retrieve the various matrices in the **VREAD** command:

θ	"params"
δ	"direct"
Σ	"hessian"
Ψ	"gradient"
ρ	"step"

2.5 Constraints

There are two general types of constraints, nonlinear equality constraints and nonlinear inequality constraints. However, for computational convenience they are divided into five types: linear equality, linear inequality, nonlinear equality, nonlinear inequality, and bounds.

2.5.1 Linear Equality Constraints

Linear constraints are of the form:

$$A\theta = B$$

where A is an $m_1 \times k$ matrix of known constants, and B an $m_1 \times 1$ vector of known constants, and θ the vector of parameters.

The specification of linear equality constraints is done by assigning the A and B matrices to the **CML** globals, **_cml_A** and **_cml_B**, respectively. For example, to constrain the first of four parameters to be equal to the third,

```
_cml_A = { 1 0 -1 0 };
_cml_B = { 0 };
```

2.5.2 Linear Inequality Constraints

Linear constraints are of the form:

$$C\theta > D$$

where C is an $m_2 \times k$ matrix of known constants, and D an $m_2 \times 1$ vector of known constants, and θ the vector of parameters.

The specification of linear equality constraints is done by assigning the C and D matrices to the **CML** globals, **_cml_C** and **_cml_D**, respectively. For example, to constrain the first of four parameters to be greater than the third, and as well the second plus the fourth greater than 10:

2.5.3 Nonlinear Equality

Nonlinear equality constraints are of the form:

$$G(\theta) = 0$$

where θ is the vector of parameters, and $G(\theta)$ is an arbitrary, user-supplied function. Nonlinear equality constraints are specified by assigning the pointer to the user-supplied function to the **GAUSS** global, **_cml_EqProc**.

For example, suppose you wish to constrain the norm of the parameters to be equal to 1:

```
proc eqp(b);
    retp(b'b - 1);
endp;
_cml_EqProc = &eqp;
```

2.5.4 Nonlinear Inequality

Nonlinear inequality constraints are of the form:

$$H(\theta) \ge 0$$

where θ is the vector of parameters, and $H(\theta)$ is an arbitrary, user-supplied function. Nonlinear equality constraints are specified by assigning the pointer to the user-supplied function to the **GAUSS** global, **_cml_lneqProc**.

For example, suppose you wish to constrain a covariance matrix to be positive definite, the lower left nonredundant portion of which is stored in elements r:r+s of the parameter vector:

```
proc ineqp(b);
    local v;
    v = xpnd(b[r:r+s]); /* r and s defined elsewhere */
    retp(minc(eigh(v)) - 1e-5);
endp;
_cml_IneqProc = &ineqp;
```

This constrains the minimum eigenvalue of the covariance matrix to be greater than a small number (1e-5). This guarantees the covariance matrix to be positive definite.

2.5.5 Bounds

Bounds are a type of linear inequality constraint. For computational convenience they may be specified separately from the other inequality constraints. To specify bounds, the lower and upper bounds respectively are entered in the first and second columns of a matrix that has the same number of rows as the parameter vector. This matrix is assigned to the **CML** global, **__cml_Bounds**.

If the bounds are the same for all of the parameters, only the first row is necessary.

To bound four parameters:

Suppose all of the parameters are to be bounded between -50 and +50, then,

```
_cml_Bounds = { -50 50 };
```

is all that is necessary.

2.5.6 Example

The following example illustrates the estimation of a tobit model with linear equality constraints, nonlinearly inequality constraints, and bounds on the parameters. The nonlinear inequality constraint constraints the norm of the coefficients to be greater than three. The bounds are provided essentially to constrain the variance parameter to be greater than zero. The linear equality constraints constrain the first and second parameters to be equal.

```
library cml;
#include cml.ext;
cmlset;
proc lpr(x,z);
   local t,s,m,u;
   s = x[4];
   m = z[.,2:4]*x[1:3,.];
   u = z[.,1] ./= 0;
   t = z[.,1]-m;
   retp(u.*(-(t.*t)./(2*s)-.5*ln(2*s*pi)) + (1-u).*(ln(cdfnc(m/sqrt(s)))));
endp;
x0 = \{ 1, 1, 1, 1 \};
_{cml_A} = \{ 1 -1 0 0 \};
_cml_B = { 0 };
proc ineq(x);
   local b;
   b = x[1:3];
   retp(b'b - 3);
endp;
_cml_IneqProc = &ineq;
_{cml_Bounds} = \{ -10 10,
               -10 10,
               -10 10,
               .01 10 };
{ x,f,g,cov,ret } = CMLPrt(CML("tobit",0,&lpr,x0));
print "linear equality Lagrangeans";
print vread(_cml_Lagrange,"lineq");
print "nonlinear inequality Lagrangeans";
print vread(_cml_Lagrange,"nlinineq");
print;
print "bounds Lagangreans";
print vread(_cml_Lagrange,"bounds");
and the output looks like this:
_____
CML Version 1.0.0
                                               2/08/95 9:51 am
```

Data Set: tobit

return code = 0
normal convergence

Mean log-likelihood -1.34034

Number of cases 100

Covariance of the parameters computed by the following method: Inverse of computed Hessian

Parameters	Estimates	Std. err.	Gradient	
P01	-0.1832	0.0710	-0.2073	
P02	-0.1832	0.0710	0.1682	
P03	1.7126	0.0152	0.1825	
P04	1.0718	0.1589	-0.0000	

Number of iterations 8

Minutes to convergence 0.06683

linear equality Lagrangeans -0.1877

nonlinear inequality Lagrangeans 0.0533

bounds Lagangreans

.

The scalar missing value for the bounds Lagrangeans indicate that they are inactive. The linear equality and nonlinear inequality constraints are active.

At times the Lagrangeans will not be scalar missing values but yet will be equal to zero and thus inactive. This indicates that the constraints became active at some point during the iterations.

2.6 Gradients

2.6.1 Analytical Gradient

To increase accuracy and reduce time, you may supply a procedure for computing the gradient, $\Psi(\theta) = \partial L/\partial \theta$, analytically.

This procedure has two input arguments, a $K \times 1$ vector of parameters and an $N_i \times L$ submatrix of the input data set. The number of rows of the data set passed in the argument to the call of this procedure may be less than the total number of observations when the data are stored in a **GAUSS** data set and there was not enough space to store the data set in RAM in its entirety. In that case subsets of the data set are passed to the procedure in sequence. The gradient procedure must be written to return a gradient (or more accurately, a "Jacobian") with as many rows as the input submatrix of the data set. Thus the gradient procedure returns an $N_i \times K$ matrix of gradients of the N_i observations with respect to the K parameters. The **CML** global, **__cml__GradProc** is then set to the pointer to that procedure. For example,

```
library cml;
#include cml.ext;
cmlset;
                  /* Function - Poisson Regression */
proc lpsn(b,z);
   local m;
   m = z[.,2:4]*b;
   retp(z[.,1].*m-exp(m));
endp;
                  /* Gradient */
proc lgd(b,z);
   retp((z[.,1]-exp(z[.,2:4]*b)).*z[.,2:4]);
endp;
x0 = \{ .5, .5, .5 \};
_cml_GradProc = &lgd;
_cml_GradCheckTol = 1e-3;
{ x,f0,g,h,retcode } = CML("psn",0,\&lpsn,x0);
call CMLPrt(x,f0,g,h,retcode);
```

In practice, unfortunately, much of the time spent on writing the gradient procedure is devoted to debugging. To help in this debugging process, **CML** can be instructed to compute the numerical gradient along with your prospective analytical gradient for comparison purposes. In the example above this is accomplished by setting **__cml_GradCheckTol** to 1e-3.

2.6.2 User-Supplied Numerical Gradient

You may substitute your own numerical gradient procedure for the one used by **CML** by default. This is done by setting the **CML** global, **_cml_UserGrad** to a pointer to the procedure.

CML includes some numerical gradient functions in **gradient.src** which can be invoked using this global. One of these procedures, **GRADRE**, computes numerical gradients using the Richardson Extrapolation method. To use this method set

```
_cml_UserNumGrad = &gradre;
```

2.6.3 Analytical Hessian

You may provide a procedure for computing the Hessian, $\Sigma(\theta) = \partial^2 L/\partial\theta\partial\theta'$. This procedure has two arguments, the $K \times 1$ vector of parameters, an $N_i \times L$ submatrix of the input data set (where N_i may be less than N), and returns a $K \times K$ symmetric matrix of second derivatives of the objection function with respect to the parameters.

The pointer to this procedure is stored in the global variable **_cml_HessProc**.

In practice, unfortunately, much of the time spent on writing the Hessian procedure is devoted to debugging. To help in this debugging process, **CML** can be instructed to compute the numerical Hessian along with your prospective analytical Hessian for comparison purposes. To accomplish this **_cml_GradCheckTol** is set to a small nonzero value.

```
library cml;
    #include cml.ext;
   proc lnlk(b,z);
        local dev,s2;
        dev = z[.,1] - b[1] * exp(-b[2]*z[.,2]);
        s2 = dev'dev/rows(dev);
        retp(-0.5*(dev.*dev/s2 + ln(2*pi*s2)));
    endp;
   proc grdlk(b,z);
        local d,s2,dev,r;
        d = \exp(-b[2]*z[.,2]);
        dev = z[.,1] - b[1]*d;
        s2 = dev'dev/rows(dev);
        r = dev.*d/s2;
/*
        retp(r~(-b[1]*z[.,2].*r));
                                         correct gradient */
        retp(r~(z[.,2].*r));
                                     /* incorrect gradient */
    endp;
   proc hslk(b,z);
        local d,s2,dev,r, hss;
        d = \exp(-b[2]*z[.,2]);
```

```
dev = z[.,1] - b[1]*d;
   s2 = dev'dev/rows(dev);
   r = z[.,2].*d.*(b[1].*d - dev)/s2;
   hss = -d.*d/s2^r-b[1].*z[.,2].*r;
   retp(xpnd(sumc(hss)));
endp;
cmlset;
_cml_HessProc = &hslk;
_cml_GradProc = &grdlk;
_{cml}Bounds = { 0 10, 0 10 }; /* constrain parameters to */
                                /* be positive
_cml_GradCheckTol = 1e-3;
startv = { 2, 1 };
{ x,f0,g,cov,retcode } = CML("nlls",0,&lnlk,startv);
call CMLPrt(x,f0,g,cov,retcode);
```

The gradient is incorrectly computed, and **CML** responds with an error message. It is clear that the error is in the calculation of the gradient for the second parameter.

```
analytical and numerical gradients differ
```

```
analytical
  numerical
-0.015387035
                -0.015387035
0.031765317
                -0.015882659
```

analytical Hessian and analytical gradient

```
CML Version 1.0.0
                                                  2/08/95 10:10 am
```

```
Data Set: nlls
______
```

return code = function cannot be evaluated at initial parameter values

```
Mean log-likelihood
                            1.12119
Number of cases
                    150
```

The covariance of the parameters failed to invert

Parameters	Estimates	Gradient
P01	2.000000	-0.015387

```
P02 1.000000 -0.015883

Number of iterations .

Minutes to convergence .
```

2.6.4 User-Supplied Numerical Hessian

You may substitute your own numerical Hessian procedure for the one used by **CML** by default. This done by setting the **CML** global, **_cml_UserHess** to a pointer to the procedure. This procedure has three input arguments, a pointer to the log-likelihood function, a $K \times 1$ vector of parameters, and an $N_i \times K$ matrix containing the data. It must return a $K \times K$ matrix which is the estimated Hessian evaluated at the parameter vector.

2.6.5 Analytical Nonlinear Constraint Jacobians

When nonlinear equality or inequality constraints have been placed on the parameters, the convergence can be improved by providing a procedure for computing their Jacobians, i.e., $\dot{G}(\theta) = \partial G(\theta)/\partial \theta$ and $\dot{H}(\theta) = \partial H(\theta)/\partial \theta$.

These procedures have one argument, the $K \times 1$ vector of parameters, and return an $M_j \times K$ matrix, where M_j is the number of constraints computed in the corresponding constraint function. Then the **CML** globals, **_cml_EqJacobian** and **_cml_IneqJacobian** are set to pointers to the nonlinear equality and inequality Jacobians, respectively. For example,

```
library cml;
#include cml.ext;
cmlset;
proc lpr(c,x);
                   /* ordinary least squares model */
    local s,t;
    t = x[.,1] - x[.,2:4]*c[1:3];
    s = c[4];
    retp(-(t.*t)./(2*s)-.5*ln(2*s*pi)));
endp;
proc ineq(c);
                      /* constrain parameter norm to be > 1 */
    local b;
    b = c[1:3];
    retp(b'b - 1);
endp;
```

```
proc ineqj(c);
                      /* constraint Jacobian */
   local b:
   b = c[1:3];
   retp((2*b')~0);
endp;
                              /* bound residual to be larger */
_{cml}Bounds = \{ -1e256 \ 1e256, 
                               /* than a small number
                -1e256 1e256,
                -1e256 1e256,
                  1e-3 1e256 };
_cml_IneqProc = &ineq;
                                 /* set pointers */
_cml_IneqJacobian = &ineqj;
x0 = \{ 1, 1, 1, 1 \};
                             /* Y and X defined elsewhere */
{ x,f,g,cov,ret } = CMLPrt(CML(Y~X,0,&lpr,x0));
```

2.7 Inference

CML includes four broad classes of methods for analyzing the distributions of the estimated parameters:

- Taylor Series covariance matrix of the parameters. This includes two
 types: the inverted Hessian and the heteroskedastic- consistent
 covariance matrix computed from both the Hessian and the cross-product
 of the first derivatives.
- Confidence limits computed from the Taylor series covariance matrix of the parameters which take into account the constraints
- Bootstrap with additional procedures for kernel density plots, histograms, surface plots, and confidence limits
- Likelihood profile and profile t traces

CML computes a Taylor-series covariance matrix of the parameters that includes the sampling distributions of the Lagrangean coefficients. However, when the model includes inequality constraints, confidence limits computed from the usual t-statistics, i.e., by simply dividing the parameter estimates by their standard errors, are incorrect because they do not account for boundaries placed on the distributions of the parameters by the inequality constraints. **CML** includes a special procedure, **CMLClimits**, for computing confidence limits in the presence of inequality constraints.

The Taylor-series methods, however, assume that it is reasonable to truncate the Taylor-series approximation to the distribution of the parameters at the second order.

If this is not reasonable, bootstrapping is an alternative not requiring this assumption. **CML** includes the procedure **CMLBoot** which generates the mean vector and covariance matrix of the bootstrapped parameters, and **CMLHist** which produces histograms and surface plots. The likelihood profile and profile t traces explicated by Bates and Watts (1988) provide diagnostic material for evaluating parameter distributions. **CMLProfile** generates trace plots which are used for this evaluation.

2.7.1 Covariance Matrix of the Parameters

An argument based on a Taylor-series approximation to the likelihood function (e.g., Amemiya, 1985, page 111) shows that

$$\hat{\theta} \to N(\theta, A^{-1}BA^{-1})$$

where

$$A = E \left[\frac{\partial^{2} L}{\partial \theta \partial \theta'} \right]$$

$$B = E \left[\left(\frac{\partial L}{\partial \theta} \right)' \left(\frac{\partial L}{\partial \theta} \right) \right]$$

Estimates of A and B are

$$\hat{A} = \frac{1}{N} \sum_{i}^{N} \frac{\partial^{2} L_{i}}{\partial \theta \partial \theta'}$$

$$\hat{B} = \frac{1}{N} \sum_{i}^{N} \left(\frac{\partial L_{i}}{\partial \theta}\right)' \left(\frac{\partial L_{i}}{\partial \theta}\right)$$

Assuming the correct specification of the model plim(A) = plim(B) and thus

$$\hat{\theta} \to N(\theta, \hat{A}^{-1})$$

Without loss of generality we may consider two types of constraints, the nonlinear equality and the nonlinear inequality constraints (the linear constraints are included in nonlinear, and the bounds are regarded as a type of linear inequality). Furthermore, the inequality constraints may be treated as equality constraints with the introduction of "slack" parameters into the model:

$$H(\theta) \ge 0$$

is changed to

$$H(\theta) = \zeta^2$$

where ζ is a conformable vector of slack parameters.

Further distinguish active from inactive inequality constraints. Active inequality constraints have nonzero Lagrangeans, γ_j , and zero slack parameters, ζ_j , while the reverse is true for inactive inequality constraints. Keeping this in mind, define the diagonal matrix, Z, containing the slack parameters, ζ_j , for the inactive constraints, and another diagonal matrix, Γ , containing the Lagrangean coefficients. Also, define $H_{\oplus}(\theta)$ representing the active constraints, and $H_{\ominus}(\theta)$ the inactive.

The likelihood function augmented by constraints is then

$$L_A = L + \lambda_1 g(\theta)_1 + \dots + \lambda_I g(\theta)^I + \gamma_1 h_{\oplus 1}(\theta) + \dots + \gamma_J h_{\oplus J}(\theta) + h_{\ominus 1}(\theta)_i - \zeta_1^2 + \dots + h_{\ominus K}(\theta) - \zeta_K^2$$

and the Hessian of the augmented likelihood is

$$E(rac{\partial^2 L_A}{\partial heta \partial heta'}) = \left[egin{array}{ccccccccc} \Sigma & 0 & 0 & G' & H'_\oplus & H'_\ominus \ 0 & 2\Gamma & 0 & 0 & 0 & 0 \ 0 & 0 & 0 & 0 & 0 & 2Z \ \dot{G} & 0 & 0 & 0 & 0 & 0 \ \dot{H}_\oplus & 0 & 0 & 0 & 0 & 0 \ \dot{H}_\ominus & 0 & 2Z & 0 & 0 & 0 \end{array}
ight]$$

where the dot represents the Jacobian with respect to θ , $L = \sum_{i=1}^{N} \log P(Y_i; \theta)$, and $\Sigma = \partial^2 L/\partial\theta\partial\theta'$. The covariance matrix of the parameters, Lagrangeans, and slack parameters is the Moore-Penrose inverse of this matrix. Usually, however, we are interested only in the covariance matrix of the parameters, as well as the covariance matrices of the Lagrange coefficients associated with the active inequality constraints and the equality constraints.

These matrices may be computed without requiring the storage and manipulation of the entire Hessian. Construct the partitioned array

$$\tilde{B} == \left[\begin{array}{c} \dot{G} \\ \dot{H}_{\oplus} \\ \dot{H}_{\ominus} \end{array} \right]$$

and denote the i-th row of \tilde{B} as \tilde{b}_i . Then the $k \times k$ upper left portion of the inverse, that is, that part associated with the estimated parameters, is calculated recursively. First, compute

$$\Omega_1 = \Sigma^{-1} - \frac{1}{\tilde{b}_1 \Sigma^{-1} \tilde{b}_1'} \Sigma^{-1} \tilde{b}_1' \tilde{b}_1 \Sigma^{-1}$$

then continue to compute for all rows of \tilde{B} :

$$\Omega_i = \Omega_{i-1} - \frac{1}{\tilde{b}_i \Omega_{i-1} \tilde{b}_i'} \Omega_{i-1} \tilde{b}_i' \tilde{b}_i \Omega_{i-1}$$

Rows associated with the inactive inequality constraints in \tilde{B} , i.e., with \dot{H}_{\ominus} , drop out and therefore they need not be considered.

Standard errors for some parameters associated with active inequality constraints may not be available, i.e., the rows and columns of Ω associated with those parameters may be all zeros.

2.7.2 Testing Constraints

Equality Constraints

When equality constraints are present in the model, their associated Lagrange coefficients may be tested to determine their reasonableness. An estimate of the covariance matrix of the joint distribution of the Lagrange coefficients associated with the equality constraints is $\dot{G}\Sigma^{-1}\dot{G}'$ and therefore

$$\hat{\lambda}' \dot{G} \Sigma^{-1} \dot{G}' \hat{\lambda}$$

is asymptotically $\chi^2(p)$ where p is the length of $\hat{\lambda}$. Individual constraints may be tested using their associated t-statistics.

When appropriate, CML inserts $\dot{G}\Sigma^{-1}\dot{G}'$ as "eqcov" in the global, **_cml_Lagrange**, using the GAUSS VPUT command.

Active Inequality Constraints

When inequality constraints are active, their associated Lagrange coefficients are nonzero. The expected value of their Lagrange coefficients is zero (assuming correct specification of the model), and they are active only in occasional samples. How many samples this occurs in depends on their covariance matrix, which is estimated by $\dot{H}_{\oplus}\Sigma^{-1}\dot{H}'_{\oplus}$.

When appropriate CML inserts $\dot{H}_{\oplus}\Sigma^{-1}\dot{H}_{\oplus}$ as "ineqcov" in the global, **_cml_Lagrange**, using the GAUSS VPUT command.

2.7.3 Heteroskedastic-consistent Covariance Matrix

When **_cml_CovPar** is set to 3, **CML** returns heteroskedastic-consistent covariance matrices of the parameters, and as well as the corresonding heteroskedastic-consistent covariance matrices of the Lagrange coefficients in the global, **_cml_Lagrange**.

Define

$$B = \left(\frac{\partial F}{\partial \theta}\right)' \left(\frac{\partial L}{\partial \theta}\right)$$

evaluated at the estimates. Then the covariance matrix of the parameters is $\Omega B\Omega$

2.7.4 Confidence Limits

When the model includes inequality constraints, confidence limits computed as the ratio of the parameter estimate to its standard error are not correct because they do not take into account that the distribution of the parameter is restricted by its boundaries.

A $1-\alpha$ joint confidence region for θ is the hyper-ellipsoid

$$JF(J, N - K; \alpha) = (\theta - \hat{\theta})'V^{-1}(\theta - \hat{\theta})$$
(2.1)

where V is the covariance matrix of the parameters, J is the number of parameters involved in the hypothesis, and $F(J, N-K; \alpha)$ is the upper α area of the F-distribution with J, N-K degrees of freedom.

If there are no constraints in the model, the $1-\alpha$ confidence interval for any selected parameter is

$$\hat{\theta} \pm \sqrt{\eta_k' V^{-1} \eta_k} \ t(N - K; \alpha/2)$$

where η_k is a vector of zeros with the k-th element corresponding to the parameter being tested set to one.

When there are constraints no such simple description of the interval is possible. Instead it is necessary to state the confidence limit problem as a parametric nonlinear programming problem.

The lower limit of the confidence limit is the solution to

$$\min\left\{\eta_k'\theta\mid (\theta-\hat{\theta})'V^{-1}(\theta-\hat{\theta})\geq JF(J,N-K;\alpha), G(\theta)=0, H(\theta)\geq 0)\right\}$$

where now η can be an arbitrary vector of constants and $J = \sum \eta_k \neq 0$, and where again we have assumed that the linear constraints and bounds have been folded in among nonlinear constraints. The upper limit is the maximum of the this same function.

In this form, the minimization is not convex and can't be solved by the usual methods. However, the problem can be re-stated as a parametric nonlinear programming problem (Rust and Burrus, 1972). Define the function

$$F(\phi) = \min((\theta - \hat{\theta})'V^{-1}(\theta - \hat{\theta}) \mid \eta_k'\theta = \phi, G(\theta) = 0, H(\theta) \ge 0)$$

The upper and lower limits of the $1-\alpha$ interval are the values of ϕ such that

$$F(\phi) = JF(J, N - K; \alpha)$$

To find this value it is necessary to iteratively refine ϕ by interpolation until 2.7.4 is satisfied (e.g., O'Leary and Rust, 1986). The **CML** procedure **CMLClimits** solves this problem.

2.7.5 Bootstrap

The bootstrap method is used to generate empirical distributions of the parameters, thus avoiding the difficulties with the usual methods of statistical inference described above.

CMLBoot

Rather than randomly sample with replacement from the data set, **CMLBoot** performs **_cml_NumSample** weighted maximum likelihood estimations where the weights are Poisson pseudo-random numbers with expected value equal to the the number of observations. This is asymptotically equivalent to simple random sampling with replacement. **_cml_NumSample** is set by the **CMLBoot** global variable. The default is 50 re-samplings. Efron and Tibshirani (1993:52) suggest that 100 is satisfactory, 50 is often enough to give a good estimate, and rarely are more than 200 needed.

The mean and covariance matrix of the bootstrapped parameters is returned by **CMLBoot**. In addition **CMLBoot** writes the bootstrapped parameter estimates to a **GAUSS** data set for use with **CMLHist**, which produces histograms and surface plots, **CMLDensity**, which produces kernel density plots, and **CMLBlimits**, which produces confidence limits based on the bootstrapped coefficients. The data set name can be specified by the user in the global **_cml_BootFname**. However, if not specified, **CMLBoot** selects the name BOOTxxxx, where xxxx starts at 0000 and increments by 1 until a name is found that is not already in use.

CMLDensity

CMLDensity is a procedure for computing kernel type density plots. The global, **__cml__Kernel** permits you to select from a variety of kernels, normal, Epanechnikov, biweight, triangular, rectangular, and truncated normal. For each selected parameter, a plot is generated of a smoothed density. The smoothing coefficients may be specified using the global, **__cml__Smoothing**, or **CMLDensity** will compute them.

CMLHist

CMLHist is a procedure for visually displaying the results of the bootstrapping in univariate histograms and bivariate surface plots for selected parameters. The univariate discrete distributions of the parameters used for the histograms are returned by **CMLHist** in a matrix.

Example

To bootstrap the example in Section 2.5.6, the only necessary alteration is the change the call to **CML** to a call to **CMLBoot**:

```
_cml_BootFname = "bootdata";
call CMLPrt(cmlboot("tobit",0,&lpr,x0));
call CMLDensity("bootdata",0);
call CMLHist("bootdata",0);
```

2.7.6 Profiling

The CML proc, CMLProfile generates profile t plots as well as plots of the likelihood profile traces for all of the parameters in the model in pairs. The profile t plots are used to assess the nonlinearity of the distributions of the individual parameters, and the likelihood profile traces are used to assess the bivariate distributions. The input and output arguments to CMLProfile are identical to those of CML. But in addition to providing the maximum likelihood estimates and covariance matrix of the parameters, a series of plots are printed to the screen using GAUSS' Publication Quality Graphics. A screen is printed for each possible pair of parameters. There are three plots, a profile t plot for each parameter, and a third plot containing the likelihood profile traces for the two parameters.

The discussion in this section is based on Bates and Watts (1988), pages 205-216, which is recommended reading for the interpretation and use of profile t plots and likelihood profile traces.

The Profile t Plot

Define

$$\tilde{\theta_k} = (\tilde{\theta}_1, \tilde{\theta}_2, ..., \tilde{\theta}_{k-1}, \theta_k, \tilde{\theta}_{k+1}, ..., \tilde{\theta}_K)$$

This is the vector of maximum likelihood estimates *conditional* on θ_k , i.e., where θ_k is fixed to some value. Further define the profile t function

$$\tau(\theta_k) = sign(\theta_k - \hat{\theta}_k)(N - K)\sqrt{2N\left[L(\tilde{\theta}_k) - L(\hat{\theta}_k)\right]}$$

For each parameter in the model, τ is computed over a range of values for θ_k . These plots provide exact likelihood intervals for the parameters, and reveal how nonlinear the estimation is. For a linear model, τ is a straight line through the origin with unit slope. For nonlinear models, the amount of curvature is diagnostic of the nonlinearity of the estimation. High curvature suggests that the usual statistical inference using the t-statistic is hazardous.

The Likelihood Profile Trace

The likelihood profile traces provide information about the bivariate likelihood surfaces. For nonlinear models the profile traces are curved, showing how the parameter estimates affect each other and how the projection of the likelihood contours onto the (θ_k, θ_ℓ) plane might look. For the (θ_k, θ_ℓ) plot, two lines are plotted, $L(\tilde{\theta}_k)$ against θ_k and $L(\tilde{\theta}_\ell)$ against θ_ℓ .

If the likelihood surface contours are long and thin, indicating the parameters to be collinear, the profile traces are close together. If the contours are fat, indicating the parameters to be more uncorrelated, the profile traces tend to be perpendicular. And if the contours are nearly elliptical, the profile traces are straight. The surface contours for a linear model would be elliptical and thus the profile traces would be straight and perpendicular to each other. Significant departures of the profile traces from straight, perpendicular lines, therefore, indicate difficulties with the usual statistical inference.

To generate profile t plots and likelihood profile traces from the example in Section 2.5.6, it is necessary only to change the call to **CML** to a call to **CMLProfile**:

```
call CMLPrt(cmlprofile("tobit",0,&lpr,x0));
```

CMLProfile produces the same output as **CML** which can be printed out using a call to **CMLPRT**.

For each pair of parameters a plot is generated containing an xy plot of the likelihood profile traces of the two parameters, and two profile t plots, one for each parameter.

The likelihood profile traces indicate that the distributions of parameters 1 and 2 are highly correlated. Ideally, the traces would be perpendicular and the trace in this example is far from ideal.

The profile t plots indicate that the parameter distributions are somewhat nonlinear. Ideally the profile t plots would be straight lines and this example exhibits significant nonlinearity. It is clear that any interpretations of the parameters of this model must be made quite carefully.

2.8 Run-Time Switches

If the user presses **Alt-H** during the iterations, a help table is printed to the screen which describes the run-time switches. By this method, important global variables may be modified during the iterations.

2. CONSTRAINED MAXIMUM LIKELIHOOD ESTIMATION

Alt-G	Toggle _cml_GradMethod
Alt-V	Revise _cml_DirTol
Alt-O	Toggle output
Alt-M	Maximum Tries
Alt-I	Compute Hessian
Alt-E	Edit Parameter Vector
Alt-C	Force Exit
Alt-A	Change Algorithm
Alt-J	Change Line Search Method
Alt-H	Help Table

The algorithm may be switched during the iterations either by pressing **Alt-A**, or by pressing one of the following:

Alt-1	Broyden-Fletcher-Goldfarb-Shanno (BFGS)
Alt-2	Davidon-Fletcher-Powell (DFP)
Alt-3	Newton-Raphson (NEWTON) or (NR)
Alt-4	Berndt, Hall, Hall & Hausman (BHHH)

The line search method may be switched during the iterations either by pressing **Alt-S**, or by pressing one of the following:

Shift-1	no search (1.0 or 1 or ONE)
Shift-2	cubic or quadratic method (STEPBT)
Shift-3	step halving method (HALF)
Shift-4	Brent's method (BRENT)
Shift-5	BHHH step method (BHHHSTEP)

2.9 Error Handling

2.9.1 Return Codes

The fourth argument in the return from **CML** contains a scalar number that contains information about the status of the iterations upon exiting **CML**. The following table describes their meanings:

2. CONSTRAINED MAXIMUM LIKELIHOOD ESTIMATION

- 0 normal convergence
- 1 forced exit
- 2 maximum iterations exceeded
- 3 function calculation failed
- 4 gradient calculation failed
- 5 Hessian calculation failed
- 6 line search failed
- 7 function cannot be evaluated at
 - initial parameter values
- 8 error with gradient
- 9 error with constraints
- 10 secant update failed
- 11 maximum time exceeded
- 12 error with weights
- 13 quadratic program failed
- 14 equality Jacobian failed
- 15 inequality Jacobian failed
- 20 Hessian failed to invert
- 34 data set could not be opened
- 99 termination condition unknown

2.9.2 Error Trapping

Setting the global **___output** = 0 turns off all printing to the screen. Error codes, however, still are printed to the screen unless error trapping is also turned on. Setting the trap flag to 4 causes **CML** to *not* send the messages to the screen:

trap 4;

Whatever the setting of the trap flag, **CML** discontinues computations and returns with an error code. The trap flag in this case only affects whether messages are printed to the screen or not. This is an issue when the **CML** function is embedded in a larger program, and you want the larger program to handle the errors.

2.10 References

Amemiya, Takeshi, 1985. Advanced Econometrics. Cambridge, MA: Harvard University Press.

Bates, Douglas M. and Watts, Donald G., 1988. *Nonlinear Regression Analysis and Its Applications*. New York: John Wiley & Sons.

2. CONSTRAINED MAXIMUM LIKELIHOOD ESTIMATION

- Berndt, E., Hall, B., Hall, R., and Hausman, J. 1974. "Estimation and inference in nonlinear structural models". *Annals of Economic and Social Measurement* 3:653-665.
- Brent, R.P., 1972. Algorithms for Minimization Without Derivatives. Englewood Cliffs, NJ: Prentice-Hall.
- Dennis, Jr., J.E., and Schnabel, R.B., 1983. Numerical Methods for Unconstrained Optimization and Nonlinear Equations. Englewood Cliffs, NJ: Prentice-Hall.
- Efron, Gradley, Robert J. Tibshirani, 1993. An Introduction to the Bootstrap. New York: Chapman & Hall.
- Fletcher, R., 1987. Practical Methods of Optimization. New York: Wiley.
- Gill, P. E. and Murray, W. 1972. "Quasi-Newton methods for unconstrained optimization." J. Inst. Math. Appl., 9, 91-108.
- Han, S.P., 1977. "A globally convergent method for nonlinear programming." *Journal of Optimization Theory and Applications*, 22:297-309.
- Hock, Willi and Schittkowski, Klaus, 1981. Lecture Notes in Economics and Mathematical Systems. New York: Springer-Verlag.
- Jamshidian, Mortaza and Bentler, P.M., 1993. "A modified Newton method for constrained estimation in covariance structure analysis." *Computational Statistics & Data Analysis*, 15:133-146.
- Judge, G.G., R.C. Hill, W.E. Griffiths, H. Lütkepohl and T.C. Lee. 1988. Introduction to the Theory and Practice of Econometrics. 2nd Edition. New York: Wiley.
- Judge, G.G., W.E. Griffiths, R.C. Hill, H. Lütkepohl and T.C. Lee. 1985. *The Theory and Practice of Econometrics*. 2nd Edition. New York: Wiley.
- O'Leary, Dianne P., and Rust, Bert W., 1986. "Confidence intervals for inequality-constrained least squares problems, with applications to ill-posed problems". American Journal for Scientific and Statistical Computing, 7(2):473-489.
- Rust, Bert W., and Burrus, Walter R., 1972. Mathematical Programming and the Numerical Solution of Linear Equations. New York: American Elsevier.
- White, H. 1981. "Consequences and detection of misspecified nonlinear regression models." *Journal of the American Statistical Association* 76:419-433.
- White, H. 1982. "Maximum likelihood estimation of misspecified models." *Econometrica* 50:1-25.

$2. \ CONSTRAINED \ MAXIMUM \ LIKELIHOOD \ ESTIMATION$

Chapter 3

Constrained Maximum Likelihood Reference

Purpose

Computes estimates of parameters of a constrained maximum likelihood function.

Library

cml

Format

```
\{x,f,g,cov,retcode\} = CML(dataset,vars,&fct,start)
```

Input

dataset string containing name of GAUSS data set

– or –

 $N \times NV$ matrix, data

vars $NV \times 1$ character vector, labels of variables selected for analysis

– or –

 $NV \times 1$ numeric vector, indices of variables selected for analysis.

If dataset is a matrix, vars may be a character vector containing either the standard labels created by **CML** (i.e., either V1, V2,..., or V01, V02,.... See discussion of the global variable **___vpad** below, or the

user-provided labels in **__altnam**).

&fct a pointer to a procedure that returns either the log-likelihood for one

observation or a vector of log-likelihoods for a matrix of observations (see discussion of the global variable **___row** in global variable section below).

start $K \times 1$ vector, start values.

Output

 $x K \times 1$ vector, estimated parameters

f scalar, function at minimum (the mean log-likelihood)

g $K \times 1$ vector, gradient evaluated at x

 $h K \times K$ matrix, covariance matrix of the parameters (see discussion of the

global variable **_cml_CovPar** below).

retcode scalar, return code. If normal convergence is achieved, then retcode = 0,

otherwise a positive integer is returned indicating the reason for the

abnormal termination:

- 0 normal convergence
- 1 forced exit.
- 2 maximum iterations exceeded.
- **3** function calculation failed.
- 4 gradient calculation failed.
- **5** Hessian calculation failed.
- 6 line search failed.
- 7 function cannot be evaluated at initial parameter values.
- 8 error with gradient
- 9 error with constraints
- 10 secant update failed
- 11 maximum time exceeded
- 12 error with weights
- 13 quadratic program failed
- 20 Hessian failed to invert
- 34 data set could not be opened.
- 99 termination condition unknown.

Globals

The globals variables used by CML can be organized in the following categories according to which aspect of the optimization they affect:

Options _cml_Options

```
<u>Constraints</u> _cml_A, _cml_B, _cml_C, _cml_D, _cml_EqProc, _cml_IneqProc, _cml_EqJacobian, _cml_IneqJacobian, _cml_Bounds, _cml_Lagrange
```

```
<u>Descent and Line Search</u> _cml_Algorithm, _cml_Delta, _cml_LineSearch, _cml_Maxtry, _cml_Extrap, _cml_Interp, _cml_RandRadius, _cml_UserSearch
```

<u>Covariance Matrix of Parameters</u> _cml_CovPar, _cml_XprodCov, _cml_HessCov, _cml_FinalHess

```
<u>Gradient</u> _cml_GradMethod, _cml_GradProc, _cml_UserNumGrad, _cml_HessProc, _cml_UserNumHess, _cml_GradStep, _cml_GradCheckTol
```

Terminations Conditions _cml_DirTol, _cml_MaxIters, _cml_MaxTime

<u>Data</u> _cml_Lag, _cml_NumObs, __weight, __row, __rowfac

Parameters _cml_Active, _cml_ParNames

Miscellaneous __title, _cml_IterData, _cml_Diagnostic

The list below contains an alphabetical listing of each global with a complete description.

_cml_A $M_1 \times K$ matrix, linear equality constraint coefficient matrix **_cml_A** is used with **_cml_B** to specify linear equality constraints:

$$_{cml_A} * X = _{cml_B}$$

where X is the $K \times 1$ unknown parameter vector.

- **_cml_Active** vector, defines fixed/active coefficients. This global allows you to fix a parameter to its starting value. This is useful, for example, when you wish to try different models with different sets of parameters without having to re-edit the function. When it is to be used, it must be a vector of the same length as the starting vector. Set elements of **_cml_Active** to 1 for an active parameter, and to zero for a fixed one.
- **_cml_Algorithm** scalar, selects optimization method:
 - 1 BFGS Broyden, Fletcher, Goldfarb, Shanno method
 - 2 DFP Davidon, Fletcher, Powell method
 - 3 NEWTON Newton-Raphson method
 - 4 BHHH Berndt, Hall, Hall, Hausman method

Default = 3

- **_cml_Delta** scalar, floor for eigenvalues of Hessian in the NEWTON algorithm. When nonzero, the eigenvalues of the Hessian are augmented to this value.
- **__cml_B** $M_1 \times 1$ vector, linear equality constraint constant vector **__cml_B** is used with **__cml_A** to specify linear equality constraints:

$$_{cml_A} * X = _{cml_B}$$

where X is the $K \times 1$ unknown parameter vector.

- **_cml_Bounds** $K \times 2$ matrix, bounds on parameters. The first column contains the lower bounds, and the second column the upper bounds. If the bounds for all the coefficients are the same, a 1x2 matrix may be used. Default = $\{-1e256\ 1e256\ \}$.
- **__cml_C** $M_3 \times K$ matrix, linear inequality constraint coefficient matrix **__cml_C** is used with **__cml_D** to specify linear inequality constraints:

$$_{cml_C} * X = _{cml_D}$$

where X is the $K \times 1$ unknown parameter vector.

_cml_CovPar scalar, type of covariance matrix of parameters

- 0 not computed
- 1 computed from Hessian calculated after the iterations
- 2 heteroskedastic-consistent covariance matrix of the parameters

Default = 1;

_cml_D $M_3 \times 1$ vector, linear inequality constraint constant vector **_cml_D** is used with **_cml_C** to specify linear inequality constraints:

$$_{cml_C} * X = _{cml_D}$$

where X is the $K \times 1$ unknown parameter vector.

_cml_Diagnostic scalar.

- **0** nothing is stored or printed
- 1 current estimates, gradient, direction, function value, Hessian, and step length are printed to the screen
- the current quantities are stored in <u>_cml_Diagnostic</u> using the VPUT command. Use the following strings to extract from <u>_cml_Diagnostic</u> using VREAD:

function	"function"
estimates	"params"
direction	"direct"
Hessian	"hessian"
gradient	"gradient"
step	"step"

When **_cml_Diagnostic** is nonzer, **__output** is forced to 1.

- **_cml_DirTol** scalar, convergence tolerance for gradient of estimated coefficients. When this criterion has been satisfied CML exits the iterations. Default = 1e-5.
- **_cml_EqJacobian** scalar, pointer to a procedure that computes the Jacobian of the nonlinear equality constraints with respect to the parameters. The procedure has one input argument, the $K \times 1$ vector of parameters, and one output argument, the $M_2 \times 1$ vector of derivatives of the constraints with respect to the parameters. For example, if the nonlinear equality constraint procedure was,

```
proc eqproc(p);
    retp(p[1]*p[2]-p[3]);
endp;
the Jacobian procedure and assignment to the global would be,
    proc eqj(p);
     retp(p[2]~p[1]~-1);
endp;
_cml_EqJacobian = &eqj;
```

_cml_EqProc scalar, pointer to a procedure that computes the nonlinear equality constraints. For example, the statement:

```
_cml_EqProc = &eqproc;
```

tells **CML** that nonlinear equality constraints are to be placed on the parameters and where the procedure computing them is to be found. The procedure must have one input argument, the $K \times 1$ vector of parameters, and one output argument, the $M_2 \times K$ matrix of computed constraints that are to be equal to zero. For example, suppose that you wish to place the following constraint:

```
P[1] * P[2] = P[3]
The proc for this is:
    proc eqproc(p);
        retp(p[1]*[2]-p[3]);
    endp;
```

_cml_Extrap scalar, extrapolation constant in BRENT. Default = 2.

- **_cml_FinalHess** $K \times K$ matrix, the Hessian used to compute the covariance matrix of the parameters is stored in **_cml_FinalHess**. This is most useful if the inversion of the hessian fails, which is indicated when **CML** returns a missing value for the covariance matrix of the parameters. An analysis of the Hessian stored in **_cml_FinalHess** can then reveal the source of the linear dependency responsible for the singularity.
- _cml_GradCheckTol scalar. Tolerance for the deviation of numerical and analytical gradients when proc's exist for the computation of analytical gradients, Hessians, and/or Jacobians. If set to zero, the analytical gradients will not be compared to their numerical versions. When adding procedures for computing analytical gradients it is highly recommended that you perform the check. Set _cml_GradCheckTol to some small value, 1e-3, say when checking. It may have to be set larger if the numerical gradients are poorly computed to make sure that CML doesn't fail when the analytical gradients are being properly computed.

- **_cml_GradMethod** scalar, method for computing numerical gradient.
 - **0** central difference
 - 1 forward difference (default)
- **_cml_GradProc** scalar, pointer to a procedure that computes the gradient of the function with respect to the parameters. For example, the statement:

```
_cml_GradProc=&gradproc;
```

tells **CML** that a gradient procedure exists as well where to find it. The user-provided procedure has two input arguments, an $K \times 1$ vector of parameter values and an N×K matrix of data. The procedure returns a single output argument, an $N \times K$ matrix of gradients of the log-likelihood function with respect to the parameters evaluated at the vector of parameter values.

For example, suppose the log-likelihood function is for a Poisson regression, then the following would be added to the command file:

```
proc lgd(b,z);
    retp((z[.,1]-exp(z[.,2:4]*b)).*z[.,2:4]);
endp;
_cml_GradProc = &lgd;
```

Default = 0, i.e., no gradient procedure has been provided.

- **_cml_GradStep** scalar, increment size for computing gradient. When the numerical gradient is performing well, set to a larger value (1e-3, say). Default is the cube root of machine precision.
- **_cml_HessCov** $K \times K$ matrix. When **_cml_CovPar** is set to 3 the information matrix covariance matrix of the parameters, i.e., the inverse of the matrix of second order partial derivatives of the log-likelihood by observations, is returned in **_cml_HessCov**.
- **_cml_HessProc** scalar, pointer to a procedure that computes the hessian, i.e., the matrix of second order partial derivatives of the function with respect to the parameters. For example, the instruction:

```
_cml_HessProc = &hessproc;
```

tells **CML** that a procedure has been provided for the computation of the hessian and where to find it. The procedure that is provided by the user must have two input arguments, a $K \times 1$ vector of parameter values and an N×P data matrix. The procedure returns a single output argument, the $K \times K$ symmetric matrix of second order derivatives of the function evaluated at the parameter values.

_cml_IneqJacobian scalar, pointer to a procedure that computes the Jacobian of the nonlinear equality constraints with respect to the parameters. The procedure has one input argument, the $K \times 1$ vector of parameters, and one output argument, the $M_4 \times K$ matrix of derivatives of the constraints with respect to the parameters. For example, if the nonlinear equality constraint procedure was,

```
proc ineqproc(p);
    retp(p[1]*p[2]-p[3]);
endp;
```

the Jacobian procedure and assignment to the global would be,

```
proc ineqj(p);
    retp(p[2]~p[1]~-1);
endp;
_cml_IneqJacobian = &ineqj;
```

_cml_IneqProc scalar, pointer to a procedure that computes the nonlinear inequality constraints. For example the statement:

```
_cml_IneqProc = &ineqproc;
```

tells **CML** that nonlinear equality constraints are to be placed on the parameters and where the procedure computing them is to be found. The procedure must have one input argument, the $K \times 1$ vector of parameters, and one output argument, the $M_4 \times K$ matrix of computed constraints that are to be equal to zero. For example, suppose that you wish to place the following constraint:

```
P[1] * P[2] >= P[3]
The proc for this is:
    proc ineqproc(p);
        retp(p[1]*[2]-p[3]);
    endp;
```

- **_cml_Interp** scalar, interpolation constant in BRENT. Default = .25.
- **_cml_IterData** 3x1 vector, contains information about the iterations. The first element contains the # of iterations, the second element contains the elapsed time in minutes of the iterations, and the third element contains a character variable indicating the type of covariance matrix of the parameters.
- _cml_Lag scalar, if the function includes lagged values of the variables _cml_Lag may be set to the number of lags. When _cml_Lag is set to a nonzero value then __row is set to 1 (that is, the function must evaluated one observation at a time), and CML passes a matrix to the user-provided function and gradient procedures. The first row in this matrix is the (i -

_cml_Lag)-th observation and the last row is the i-th observation. The read loop begins with the ($_$ cml $_$ Lag+1)-th observation. Default = 0.

_cml_Lagrange vector, created using **VPUT**. Contains the Lagrangean coefficients for the constraints. They may be extracted with the **VREAD** command using the following strings:

"lineq"	linear equality constraints
"nlineq"	nonlinear equality constraints
"linineq"	linear inequality constraints
"nlinineq"	nonlinear inequality constraints
"bounds"	bounds
"eqcov"	covariance matrix of
	equality Lagrangeans
"ineqcov"	covariance matrix of
	inequality Lagrangeans
"boundcov"	covariance matrix of
	bounds' Lagrangeans

When an inequality or bounds constraint is active, its associated Lagrangean is nonzero. The linear Lagrangeans preced the nonlinear Lagrangeans in the covariance matrices.

_cml_LineSearch scalar, selects method for conducting line search. The result of the line search is a *step length*, i.e., a number which reduces the function value when multiplied times the direction..

- 1 step length = 1.
- **2** cubic or quadratic step length method (STEPBT)
- **3** step halving (HALF)
- 4 Brent's step length method (BRENT)
- **5** BHHH step length method (BHHHSTEP)

Default = 2.

Usually $_cml_LineSearch = 2$ is best. If the optimization bogs down, try setting $_cml_LineSearch = 1$, 4 or 5. $_cml_LineSearch = 3$ generates slower iterations but faster convergence and $_cml_LineSearch = 1$ generates faster iterations but slower convergence.

When any of these line search methods fails, **CML** attempts a random search of radius **_cml_RandRadius** times the truncated log to the base 10 of the gradient when **_cml_RandRadius** is set to a nonzero value. If **_cml_UserSearch** is set to 1, **CML** enters an interactive line search mode.

_cml_MaxIters scalar, maximum number of iterations.

- **_cml_MaxTime** scalar, maximum time in iterations in minutes. This global is most useful in bootstrapping. You might want 100 re-samples, but would be happy with anything more than 50 depending on the time it took. Set **_cml_NumSample** = 100, and **_cml_MaxTime** to maximum time you would be willing to wait for results. Default = 1e+5, about 10 weeks.
- **_cml_MaxTry** scalar, maximum number of tries to find step length that produces a descent.
- **_cml_NumObs** scalar, number of cases in the data set that was analyzed.
- **_cml_Options** character vector, specification of options. This global permits setting various **CML** options in a single global using identifiers. The following

```
_cml_Options = { newton stepbt forward screen };
```

the descent method to NEWTON, the line search method to STEPBT, the numerical gradient method to forward differences, and **__OUTPUT** = 2.

The following is a list of the identifiers:

Algorithms BFGS, DFP, NEWTON, BHHH

Line Search ONE, STEPBT, HALF, BRENT, BHHHSTEP

Covariance Matrix NOCOV, INFO, XPROD, HETCON

Gradient method CENTRAL, FORWARD

Output method NONE, FILE, SCREEN

- __output
- scalar, determines printing of intermediate results. Generally when **__output** is nonzero, i.e., where there some kind of printing during the iterations, the time of the iterations is degraded.
- **0** nothing is written
- 1 serial ASCII output format suitable for disk files or printers
- 2 output is suitable for screen only. ANSI.SYS must be active.
- \geq 5 same as **__output** = 1 except that information is printed only every **__output**-th iteration.

When **_cml_Diagnostic** is nonzero, **__output** is forced to 1.

- **_cml_ParNames** $K \times 1$ character vector, parameter labels.
- **_cml_UserNumGrad** scalar, pointer to user provided numerical gradient procedure.

 The instruction

_cml_UserNumGrad = &userproc;

tells **CML** that a procedure for computing the numerical gradients exists. The user-provided procedure has three input arguments, a pointer to a function that computes the log-likelihood function, a $K \times 1$ vector of parameter values, and an $K \times P$ matrix of data. The procedure returns a single output argument, an $N \times K$ matrix of gradients of each row of the input data matrix with respect to each parameter.

__row

scalar, specifies how many rows of the data set are read per iteration of the read loop. See the *REMARKS* Section for a more detailed discussion of how to set up your log-likelihood to handle more than one row of your data set. By default, the number of rows to be read is calculated by **CML**.

__rowfac

scalar, "row factor". If **CML** fails due to insufficient memory while attempting to read a **GAUSS** data set, then **__rowfac** may be set to some value between 0 and 1 to read a *proportion* of the original number of rows of the **GAUSS** data set. For example, setting

```
_{rowfac} = 0.8;
```

causes **GAUSS** to read in 80% of the rows of the **GAUSS** data set that were read when **CML** failed due to insufficient memory.

This global has an affect only when $__row = 0$. Default = 1.

__title

string title of run

_cml_UserNumHess scalar, pointer to user provided numerical Hessian procedure.

The instruction

```
_cml_UserHess = &hessproc;
```

tells **CML** that a procedure for computing the numerical Hessian exists. The user-provided procedure three input arguments, a pointer to a function that computes the log-likelihood function, a $K \times 1$ vector of parameter values, and an N×P matrix of data. The procedure returns a single output argument, a $K \times K$ Hessian matrix of the function with respect to the parameters.

- **_cml_UserSearch** scalar, if nonzero and if all other line search methods fail **CML** enters an interactive mode in which the user can select a line search parameter
- **__weight** vector, frequency of observations. By default all observations have a frequency of 1. zero frequencies are allowed. It is assumed that the elements of **__weight** sum to the number of observations.
- **_cml_XprodCov** $K \times K$ matrix. When **_cml_CovPar** is set to 3 the cross-product matrix covariance matrix of the parameters, i.e., the inverse of the cross-product of the first derivatives of the log-likelihood computed by observations, is is returned in **_cml_XprodCov**.

Remarks

Specifying Constraints.

There are five types of constraints: linear equality, linear inequality, nonlinear equality, nonlinear inequality, bounds Linear constraints are specified by initializing the appropriate <code>CML</code> globals to known matrices of constants. The linear equality constraint matrices are <code>_cml_A</code> and <code>_cml_B</code>, and they assume the following relationship with the parameter vector:

```
_{cml_A} * x = _{cml_B}
```

where x is the parameter vector.

Similarly, the linear *inequality* constraint matrices are **_cml_C** and **_cml_D**, and assume the following relationship with the parameter vector:

```
_{cml_C} * x >= _{cml_D}
```

The nonlinear constraints are specified by providing procedures and assigning their pointers to **CML** globals. These procedures take a single argument, the vector of parameters, and return a column vector of evaluations of the constraints at the parameters. Each element of the column vector is a separate constraint.

For example, suppose you wish to constrain the product of the first and third coefficients to be equal to 10, and the squared second and fourth coefficients to be equal to the squared fifth coefficient:

```
proc eqp(x);
   local c;
   c = zeros(2,1);
   c[1] = x[1] * x[3] - 10;
   c[2] = x[2] * x[2] + x[4] * x[4] - x[5] * x[5];
   retp(c);
endp;
_cml_EqProc = &eqp;
```

The nonlinear equality constraint procedure causes **CML** to find estimates for which its evaluation is equal to a conformable vector of zeros.

The nonlinear *inequality* constraint procedure is similar to the equality procedure. **CML** finds estimates for which the evaluation of the procedure is greater than or equal to zero. The nonlinear inequality constraint procedure is assigned to the global **_cml_IneqProc**. For example, suppose you wish to constrain the norm of the coefficients to be greater than one:

```
proc ineqp(x);
    retp(x'x-3);
endp;
_cml_IneqProc = &eqp;
```

Bounds are a type of linear inequality constraint. They are specified separately for computational and notational convenience. To declare bounds on the parameters assign a two column vector with rows equal to the number of parameters to the **CML** global, **__cml_Bounds**. The first column is the lower bounds and the second column the upper bounds. For example,

```
_cml_Bounds = { 0 10,
-10 0
-10 20 };
```

If the bounds are the same for all of the parameters, only the first row is required.

Writing the Log-likelihood Function

The user must provide a procedure for computing the log-likelihood for either one observation, or for a matrix of observations. The procedure must have two input arguments: first, a vector of parameter values, and second, one or more rows of the data matrix. The output argument is the log-likelihood for the observation or observations in the second argument evaluated at the parameter values in the first argument. Suppose that the function procedure has been named pfct, the following considerations apply:

The format of the procedure is:

```
 logprob = pfct(x,y); \\  where \\  x \qquad column vector of parameters of model \\  y \qquad one or more rows of the data set (if the data set has been transformed, or if <math>vars \neq 0, i.e., there is selection, then y is a transformed, selected observation)  if \begin{tabular}{l} --row = n, & then $n$ rows of the data set are read at a time \\  if \begin{tabular}{l} --row = 0, & the maximum number of rows that fit in memory is computed by $CML$.
```

The output from the procedure *pfct* is the log-likelihood for a single observation or a vector of log-likelihoods for a set of observations. If it is not possible to compute the log-likelihood for a set of observations, then either **___row** may be set to 1 to force

CML to send one observation at a time to pfct or the procedure computing the function may contain a loop. If possible, pfct should be written to compute a vector of log-likelihoods for a set of observations because this speeds up the computations significantly. If $_cml_Lag \ge 1$, then $__row$ is forced to 1.

Setting ___row= 0 causes CML to send the entire matrix to pfct if it is stored entirely in memory, or to compute the maximum number of rows if it is a GAUSS data set stored on disk (Note that even if the data starts out in a GAUSS data set, CML determines whether the data set fits in memory, and if it does, then it reads the data set into an array in memory). If you are getting insufficient memory messages, then set ___rowfac to a positive value less than 1.

Supplying an Analytical GRADIENT Procedure

To decrease the time of computation, the user may provide a procedure for the calculation of the gradient of the log-likelihood. The global variable **_cml_GradProc** must contain the pointer to this procedure. Suppose the name of this procedure is *gradproc*. Then,

```
g = gradproc(x, y);
```

where the input arguments are

- x vector of coefficients
- y one or more rows of data set.

and the output argument is

g row vector of gradients of log-likelihood with respect to coefficients, or a matrix of gradients (i.e., a Jacobian) if the data passed in y is a matrix (unless **_cml_Lag** ≥ 1 in which case the data passed in y is a matrix of lagged values but a row vector of gradients is passed back in g).

It is important to note that the gradient is row oriented. Thus if the function that computes the log-likelihood returns a scalar value ($__{row} = 1$), then a row vector of the first derivatives of the log-likelihood with respect to the coefficients must be returned, but if the procedure that computes the log-likelihood returns a column vector, then $_{cml}_{GradProc}$ must return a matrix of first derivatives in which rows are associated with observations and columns with coefficients.

Providing a procedure for the calculation of the first derivatives also has a significant effect on the calculation time of the Hessian. The calculation time for the numerical computation of the Hessian is a quadratic function of the size of the matrix. For large matrices, the calculation time can be very significant. This time can be reduced to a linear function of size if a procedure for the calculation of analytical first derivatives is

available. When such a procedure is available, **CML** automatically uses it to compute the numerical Hessian.

The major problem one encounters when writing procedures to compute gradients and Hessians is in making sure that the gradient is being properly computed. **CML** checks the gradients and Hessian when **_cml_GradCheckTol** is nonzero. **CML** generates both numerical and analytical gradients, and viewing the discrepancies between them can help in debugging the analytical gradient procedure.

Supplying an Analytical HESSIAN Procedure.

Selection of the NEWTON algorithm becomes feasible if the user supplies a procedure to compute the Hessian. If such a procedure is provided, the global variable **__cml_HessProc** must contain a pointer to this procedure. Suppose this procedure is called *hessproc*, the format is

```
h = hessproc(x,y);
```

The input arguments are

 $x K \times 1$ vector of coefficients

y one or more rows of data set

and the output argument is

 $K \times K$ matrix of second order partial derivatives evaluated at the coefficients in x.

To compare numerical and analytical Hessians set $_cml_GradCheckTol$ to a nonzero value.

Supplying Analytical Jacobians of the Nonlinear Constraints.

At each iteration the Jacobians of the nonlinear constraints, if they exist, are computed numerically. This is time-consuming and generates a loss of precision. For models with a large number of inequality constraints a significant speed-up can be achieved by providing analytical Jacobian procedures. The improved accuracy can also have a significant effect on convergence.

The Jacobian procedures take a single argument, the vector of parameters, and return a matrix of derivatives of each constraint with respect to each parameter. The rows are associated with the constraints and the columns with the parameters. The pointer to the nonlinear equality Jacobian procedure is assigned to **_cml_EqJacobian**. The pointer to the nonlinear *inequality* Jacobian procedure is assigned to **_cml_IneqJacobian**.

For example, suppose the following procedure computes the equality constraints:

```
proc eqp(x);
    local c;
    c = zeros(2,1);
    c[1] = x[1] * x[3] - 10;
    c[2] = x[2] * x[2] + x[4] * x[4] - x[5] * x[5];
    retp(c);
endp;
cml_EqProc = &eqp;
```

Then the Jacobian procedure would look like this:

```
proc eqJacob(x);
    local c;
    c = zeros(2,5);
    c[1,1] = x[3];
    c[1,3] = x[1];
    c[2,2] = 2*x[2];
    c[2,4] = 2*x[4];
    c[3,5] = -2*x[5];
    retp(c);
endp;
_cml_EqJacobian = &eqJacob;
```

The Jacobian procedure for the nonlinear inequality constraints is specified similarly, except that the associated global containing the pointer to the procedure is **_cml_IneqJacobian**.

Source

cml.src

Purpose

Produces kernel Density plots of bootstrapped parameters in GAUSS data set

Library

cml

Format

```
cl = CMLBlimits(dataset)
```

Input

```
\begin{array}{c} \textit{dataset} & \textit{string containing name of $\texttt{GAUSS}$ data set} \\ & -\textit{or} - \\ & \textit{N} \times \textit{K matrix, data} \end{array}
```

Output

cl $K \times 2$ matrix, lower (first column) and upper (second column) confidence limits of the selected parameters

Globals

```
_cml_Alpha (1 - _cml_Alpha)% confidence limits are computed. Default = .05
_cml_Select selection vector for selecting coefficients to be included in profiling, for example
_cml_Select = { 1, 3, 4 };
```

selects the 1st, 3rd, and 4th parameters for profiling.

Remarks

CMLBlimits sorts each column of the parameter data set and computes (1-_cml_Alpha)% confidence limits by measuring back _cml_Alpha/2 times the number of rows from each end of the columns. The confidence limits are the values in those elements. If amount to be measured back from each end of the columns doesn't fall exactly on an element of the column, the confidence limit is interpolated from the bordering elements.

Source

cmlblim.src

CMLClimits

3. CONSTRAINED MAXIMUM LIKELIHOOD REFERENCE

Purpose

Computes confidence limits for inequality constrained parameter estimates

Library

cml

Format

climits = CMLClimits(b, V)

Input

b $K \times 1$ vector, parameter estimates

 $V K \times K$ matrix, covariance matrix of parameters in b

Output

climits

 $K\times 2$ matrix, lower confidence limits in the first column and upper limits in second column

Globals

_cml_Alpha scalar, **CMLClimits** computes (1-**_cml_Alpha**)% confidence intervals, where **_cml_Alpha** varies between zero and one. Default = .05.

_cml_A $M_1 \times K$ matrix, linear equality constraint coefficient matrix **_cml_A** is used with **_cml_B** to specify linear equality constraints:

$$_{cml_A} * X = _{cml_B}$$

where X is the $K \times 1$ unknown parameter vector.

__cml_B $M_1 \times 1$ vector, linear equality constraint constant vector **__cml_B** is used with **__cml_A** to specify linear equality constraints:

$$_{cml_A} * X = _{cml_B}$$

where X is the $K \times 1$ unknown parameter vector.

_cml_Bounds $K \times 2$ matrix, bounds on parameters. The first column contains the lower bounds, and the second column the upper bounds. If the bounds for all the coefficients are the same, a 1x2 matrix may be used. Default = $\{-1e256\ 1e256\ \}$.

__cml_C $M_3 \times K$ matrix, linear inequality constraint coefficient matrix **__cml_C** is used with **__cml_D** to specify linear inequality constraints:

$$_{cml_C} * X = _{cml_D}$$

where X is the $K \times 1$ unknown parameter vector.

_cml_D $M_3 \times 1$ vector, linear inequality constraint constant vector **_cml_D** is used with **_cml_C** to specify linear inequality constraints:

$$_{cml_C} * X = _{cml_D}$$

where X is the $K \times 1$ unknown parameter vector.

_cml_EqJacobian scalar, pointer to a procedure that computes the Jacobian of the nonlinear equality constraints with respect to the parameters. The procedure has one input argument, the $K \times 1$ vector of parameters, and one output argument, the $M_2 \times K$ matrix of derivatives of the constraints with respect to the parameters. For example, if the nonlinear equality constraint procedure was,

```
proc eqproc(p);
    retp(p[1]*p[2]-p[3]);
endp;
```

the Jacobian procedure and assignment to the global would be,

```
proc eqj(p);
    retp(p[2]~p[1]~-1);
endp;
_cml_EqJacobian = &eqj;
```

_cml_EqProc scalar, pointer to a procedure that computes the nonlinear equality constraints. For example, the statement:

```
_cml_EqProc = &eqproc;
```

tells **CMLClimits** that nonlinear equality constraints are to be placed on the parameters and where the procedure computing them is to be found. The procedure must have one input argument, the $K \times 1$ vector of parameters, and one output argument, the $M_2 \times K$ matrix of computed constraints that are to be equal to zero. For example, suppose that you wish to place the following constraint:

```
P[1] * P[2] = P[3]
The proc for this is:
    proc eqproc(p);
        retp(p[1]*[2]-p[3]);
    endp;
```

_cml_IneqJacobian scalar, pointer to a procedure that computes the Jacobian of the nonlinear equality constraints with respect to the parameters. The procedure has one input argument, the $K \times 1$ vector of parameters, and one output argument, the $M_4 \times K$ matrix of derivatives of the constraints with respect to the parameters. For example, if the nonlinear equality constraint procedure was,

```
proc ineqproc(p);
    retp(p[1]*p[2]-p[3]);
endp;
```

the Jacobian procedure and assignment to the global would be,

```
proc ineqj(p);
    retp(p[2]~p[1]~-1);
endp;
_cml_IneqJacobian = &ineqj;
```

_cml_IneqProc scalar, pointer to a procedure that computes the nonlinear inequality constraints. For example the statement:

```
_cml_IneqProc = &ineqproc;
```

tells **CMLClimits** that nonlinear equality constraints are to be placed on the parameters and where the procedure computing them is to be found. The procedure must have one input argument, the $K \times 1$ vector of parameters, and one output argument, the $M_4 \times K$ matrix of computed constraints that are to be equal to zero. For example, suppose that you wish to place the following constraint:

```
P[1] * P[2] >= P[3]
The proc for this is:
```

```
proc ineqproc(p);
    retp(p[1]*[2]-p[3]);
endp;
```

_cml_MaxIters scalar, maximum number of iterations.

_cml_NumObs scalar, number of cases in the data set that was analyzed.

_cml_ParNames $K \times 1$ character vector, parameter labels.

_cml_Select selection vector for selecting coefficients to be included in the analysis, for example

```
_cml_Select = { 1, 3, 4 };
```

selects the 1st, 3rd, and 4th parameters for analysis.

Remarks

Confidence limits for inequality constrained models are the solutions to a parametric nonlinear programming problem. **CMLClimits** solves this problem given a covariance matrix, the vector of parameter estimates, and given the model constraints.

The calculation of confidence limits for large models can be time consuming. In that case it might be necessary to select parameters for analysis. This can be done using the **CML** global, **_cml_Select**.

The global **_cml_NumObs** must be set. If **CMLClimits** is called immediately after a call to **CML**, **_cml_NumObs** will be set by **CML**.

Source

cmlclim.src

Purpose

Computes confidence limits based on t-statistics

Library

cml

Format

cl = CMLTlimits(b, cov)

Input

 $b K \times 1$ vector, parameter estimates

cov $K \times K$ matrix, covariance matrix of parameter estimates

Output

cl $K \times 2$ matrix, lower (first column) and upper (second column) confidence limits of the selected parameters

Globals

_cml_Alpha (1-**_cml_Alpha**)% confidence limits are computed.

Default = .05

_cml_NumObs scalar, number of observations. Must be set.

_cml_Select selection vector for selecting coefficients to be included in profiling, for example

selects the 1st, 3rd, and 4th parameters for profiling.

Remarks

CMLTlimits returns $b[i] \pm t(_cml_NumObs - K; _cml_Alpha/2) \times \sqrt{cov[i, i]}$

_cml_NumObs must be set.

The global **_cml_NumObs** must be set. If **CMLTlimits** is called immediately after a call to **CML**, **_max_NumObs** will be set by **CML**.

Source

cml.src

Purpose

Computes bootstrapped estimates of parameters of a constrained maximum likelihood function.

Library

cml

Format

```
\{x,f,g,cov,retcode\} = CMLBoot(dataset,vars,&fct,start)
```

Input

datasetstring containing name of GAUSS data set

– or –

 $N\times NV$ matrix, data

 $NV \times 1$ character vector, labels of variables selected for analysis vars

– or –

 $NV \times 1$ numeric vector, indices of variables selected for analysis.

If dataset is a matrix, vars may be a character vector containing either the standard labels created by **CMLBoot** (i.e., either V1, V2,..., or V01, V02,..... See discussion of the global variable **___vpad** below, or the

user-provided labels in **__altnam**).

&fct a pointer to a procedure that returns either the log-likelihood for one

> observation or a vector of log-likelihoods for a matrix of observations (see discussion of the global variable **___row** in global variable section below).

start $K \times 1$ vector, start values.

Output

 $K \times 1$ vector, means of re-sampled parameters

scalar, mean re-sampled function at minimum (the mean log-likelihood) f

 $K \times 1$ vector, means of re-sampled gradients evaluated at the estimates

h $K \times K$ matrix, covariance matrix of the re-sampled parameters

retcodescalar, return code. If normal convergence is achieved, then retcode = 0, otherwise a positive integer is returned indicating the reason for the

abnormal termination:

3. CONSTRAINED MAXIMUM LIKELIHOOD REFERENCE

- 0 normal convergence
- 1 forced exit.
- 2 maximum iterations exceeded.
- **3** function calculation failed.
- 4 gradient calculation failed.
- **5** Hessian calculation failed.
- 6 line search failed.
- 7 function cannot be evaluated at initial parameter values.
- 8 error with gradient
- **9** error with constraints
- 10 secant update failed
- 11 maximum time exceeded
- **12** error with weights
- 13 quadratic program failed
- data set could not be opened.
- 99 termination condition unknown.

Globals

The **CML** procedure global variables are also applicable.

_cml_BootFname string, file name of **GAUSS** data set (do not include .DAT extension) containing bootstrapped parameter estimates. If not specified, **CMLBoot** selects a temporary filename.

_cml_MaxTime scalar, maximum amount of time spent in re-sampling. Default = 1e5 (about 10 weeks).

_cml_NumSample scalar, number of samples to be drawn. Default = 100.

Remarks

CMLBoot implements random sampling with replacement by computing **__cml_NumObs** pseudo-random Poisson variates and using them as weights in a call to **CML. CMLBoot** returns the mean vector of the estimates in the first argument and the covariance matrix of the estimates in the third argument.

A GAUSS data set is also generated containing the bootstrapped parameter estimates. The file name of the data set is either the name found in the global _cml_BootFname, or a temporary name. If CMLBoot selects a file name, it returns that file name in _cml_BootFname. The coefficients in this data set may be used as input to the CML procedures CMLHist and CMLDensity for further analysis.

Source

cmlboot.src

Purpose

Generates kernel density plots from GAUSS data sets

Library

cml, pgraph

Format

```
\{ px, py, smth \} = CMLDensity(dataset, vars)
```

Input

dataset string containing name of GAUSS data set

– or –

N×K matrix, data

vars $K \times 1$ character vector, labels of variables selected for analysis

- or -

 $K \times 1$ numeric vector, indices of variables selected for analysis.

If dataset is a matrix, vars may be a character vector containing either the standard labels created by **CMLDensity** (i.e., either V1, V2,..., or V01, V02,.... See discussion of the global variable **___vpad** below, or the

user-provided labels in __altnam).

Output

px __cml_NumPoints \times K matrix, abscissae of plotted points

py __cml_NumPoints \times K matrix, ordinates of plotted points

smth K \times 1 vector, smoothing coefficients

Globals

The **CML** procedure global variables are also applicable.

_cml_Kernel $K \times 1$ character vector, type of kernel:

NORMAL normal kernel
EPAN Epanechnikov kernel
BIWGT biweight kernel
TRIANG triangular kernel

CMLDensity

3. CONSTRAINED MAXIMUM LIKELIHOOD REFERENCE

RECTANG rectangular kernel

TNORMAL truncated normal kernel

If $_$ cml $_$ Kernel is scalar, the kernel is the same for all parameter densities. Default = NORMAL.

- **__cml__NumPoints** scalar, number of points to be computed for plots
- **_cml_EndPoints** $K \times 2$ matrix, lower (in first column) and upper (in second column) endpoints of density. Default is minimum and maximum, respectively, of the parameter values. If 1×2 matrix, endpoints are the same for all parameters.
- **_cml_Smoothing** $K \times 1$ vector, smoothing coefficients for each plot. If scalar, smoothing coefficient is the same for each plot. If zero, smoothing coefficient is computed by **CMLDensity**. Default = 0.
- **_cml_Truncate** $K \times 2$ matrix, lower (in first column) and upper (in second column) truncation limits for truncated normal kernel. If 1x2 matrix, truncations limits are the same for all plots. Default is minimum and maximum, respectively.
- **__output** If nonzero, K density plots are printed to the screen, otherwise no plots are generated.
- Source

cmldens.src

Purpose

Generates histograms and surface plots from GAUSS data sets

Library

cml, pgraph

Format

```
{ tab, cut } = CMLHist(dataset,vars)
```

Input

dataset string containing name of GAUSS data set

- or -

N×K matrix, data

vars $K \times 1$ character vector, labels of variables selected for analysis

– or –

 $K \times 1$ numeric vector, indices of variables selected for analysis.

If dataset is a matrix, vars may be a character vector containing either the standard labels created by **CMLHist** (i.e., either V1, V2,..., or V01, V02,..... See discussion of the global variable **___vpad** below, or the

user-provided labels in **__altnam**).

Output

tab __cml_NumCat imes K matrix, univariate distributions of bootstrapped

parameters

cut __cml_NumCat \times K matrix, cutting points

Globals

The **CML** procedure global variables are also applicable.

_cml_Center $K \times 1$ value of center category in histograms. Default is initial coefficient estimates.

_cml_CutPoint _cml_NumCat \times 1 vector, output, cutting points for histograms

_cml_Increment $K \times 1$ vector, increments for cutting points of the histograms. Default is $2 * _cml_Width * std dev / _cml_NumCat$.

_cml_NumCat scalar, number of categories in the histograms

CMLHist

3. CONSTRAINED MAXIMUM LIKELIHOOD REFERENCE

_cml_Width scalar, width of histograms, default = 2

__output If nonzero, K density plots are printed to the screen, otherwise no plots are generated.

Remarks

If **__output** is nonzero, K(K-1)/2 plots are printed to the screen displaying univariate histograms and bivariate surface plots of the bootstrapped parameter distributions in pairs.

The globals, _cml_Center, _cml_Width, and _cml_Increment may be used to establish cutting points (which is stored in _cml_Increment) for the tables of re-sampled coefficients in tab The numbers in _cml_Center fix the center categories, _cml_Width is a factor which when multiplied times the standard deviation of the estimate determines the increments between categories. Alternatively, the increments between categories can be fixed directly by supplying them in _cml_Increment.

Source

cmlhist.src

Library

cml, pgraph

Purpose

Computes profile t plots and likelihood profile traces for constrained maximum likelihood models

Format

```
\{x,f,g,cov,retcode\} = CMLProfile(dataset,vars,&fct,start)
```

Input

dataset string containing name of GAUSS data set

- or -

 $N\times NV$ matrix, data

vars $NV \times 1$ character vector, labels of variables selected for analysis

– or –

 $NV \times 1$ numeric vector, indices of variables selected for analysis.

If dataset is a matrix, vars may be a character vector containing either the standard labels created by **CMLProfile** (i.e., either V1, V2,..., or V01, V02,.... See discussion of the global variable **___vpad** below, or the

user-provided labels in **__altnam**).

&fct a pointer to a procedure that returns either the log-likelihood for one

observation or a vector of log-likelihoods for a matrix of observations (see discussion of the global variable **___row** in global variable section below).

start $K \times 1$ vector, start values.

Output

 $x K \times 1$ vector, parameter estimates

f scalar, log-likelihood at maximum

 $g K \times 1$ vector, gradients evaluated at the estimates

 $h K \times K$ matrix, covariance matrix of the parameters

retcode scalar, return code. If normal convergence is achieved, then retcode = 0,

otherwise a positive integer is returned indicating the reason for the

abnormal termination:

CMLProfile

3. CONSTRAINED MAXIMUM LIKELIHOOD REFERENCE

- **0** normal convergence
- 1 forced exit.
- 2 maximum iterations exceeded.
- **3** function calculation failed.
- 4 gradient calculation failed.
- **5** Hessian calculation failed.
- 6 line search failed.
- 7 function cannot be evaluated at initial parameter values.
- 8 error with gradient
- **9** error with constraints
- 10 secant update failed
- 11 maximum time exceeded
- 12 error with weights
- 13 quadratic program failed
- 34 data set could not be opened.
- 99 termination condition unknown.

Globals

The **CML** procedure global variables are also relevant.

- **_cml_NumCat** scalar, number of categories in profile table. Default = 16.
- **_cml_Increment** K \times 1 vector, increments for cutting points, default is 2 * **_cml_Width** * std dev / **_cml_NumCat**. If scalar zero, increments are computed by **CMLProfile**.
- **_cml_Center** $K \times 1$ vector, value of center category in profile table. Default values are coefficient estimates.
- **_cml_Select** selection vector for selecting coefficients to be included in profiling, for example

```
_cml_Select = { 1, 3, 4 };
```

selects the 1st, 3rd, and 4th parameters for profiling.

_cml_Width scalar, width of profile table in units of the standard deviations of the parameters. Default = 2.

Remarks

For each pair of the selected parameters, three plots are printed to the screen. Two of the are the profile t trace plots that describe the univariate profiles of the parameters, and one of them is the profile likelihood trace describing the joint distribution of the two parameters. Ideally distributed parameters would have univariate profile t traces that are straight lines, and bivariate likelihood profile traces that are two straight lines intersecting at right angles. This ideal is generally not met by nonlinear models, however, large deviations from the ideal indicate serious problems with the usual statistical inference.

Source

cmlprof.src

Purpose

Resets CONSTRAINED MAXIMUM LIKELIHOOD global variables to default values.

Library

cml

■ Format

CMLSet;

Input

None

Output

None

Remarks

Putting this instruction at the top of all command files that invoke **CML** is generally good practice. This prevents globals from being inappropriately defined when a command file is run several times or when a command file is run after another command file has executed that calls **CML**.

Source

cml.src

Formats and prints the output from a call to CML.

Library

cml

Format

```
\{x,f,g,h,retcode\} = CMLPrt(x,f,g,h,retcode);
```

Input

Output

The input arguments are returned unchanged.

Globals

__header string. This is used by the printing procedure to display information about the date, time, version of module, etc. The string can contain one or more of the following characters:

```
"t" print title (see ___title)

"l" bracket title with lines

"d" print date and time Example:

"v" print version number of program

"f" print file name being analyzed

__header = "tld";

Default = "tldvf".

__title string, message printed at the top of the screen and printed out by

CMLPrt. Default = "".
```

Remarks

The call to CML can be nested in the call to CMLPrt:

```
{ x,f,g,h,retcode } = CMLPrt(CML(dataset,vars,&fct,start));
```

Source

cml.src

3. CONSTRAINED MAXIMUM LIKELIHOOD REFERENCE

Purpose

Formats and prints the output from a call to CML along with confidence limits

Library

cml

Format

```
{ x,f,g,cl,retcode } = CMLCLPrt(x,f,g,cl,retcode);
```

Input

```
x 	 K \times 1 vector, parameter estimates
```

f scalar, value of function at minimum

g $K \times 1$ vector, gradient evaluated at x

cl $K \times 2$ matrix, lower and upper confidence limits

The lower limits are in the first column and the upper limits are in the second column.

retcode scalar, return code.

Output

The input arguments are returned unchanged.

Globals

__header string. This is used by the printing procedure to display information about the date, time, version of module, etc. The string can contain one or more of the following characters:

```
"t" print title (see ___title)
"l" bracket title with lines
"d" print date and time Example:
"v" print version number of program
"f" print file name being analyzed
__header = "tld";

Default = "tldvf".
```

___title string, message printed at the top of the screen and printed out by \mathbf{CMLPrt} . Default = "".

Remarks

Confidence limits computed by **CMLBlimits**, **CMLClimits**, or **CMLTlimits** may be passed in the fourth argument in the call to **CMLCLPrt**:

```
{ b,f,g,cov,ret } = CMLBoot("tobit",0,&lpr,x0);
cl = CMLBlimits(_cml_BootFname);
call CMLCLPrt(b,f,g,cl,ret);
```

CMLCLPrt

3. CONSTRAINED MAXIMUM LIKELIHOOD REFERENCE

Chapter 4

Constrained Event Count and Duration Regression

by

Gary King

Department of Government

Harvard University

This module contains procedures for estimating statistical models of event count or duration data with general nonlinear equality and inequality constraints on the parameters

The programs included in this module implement maximum likelihood estimators for parametric statistical models of events data. Data based on events come in two forms: event counts and durations between events. Event counts are dependent variables that take on only nonnegative integer values, such as the number of wars in a year, the number of medical consultations in a month, the number of patents per firm, or even the frequency in the cell of a contingency table. Dependent variables that are measured as durations between events measure time and may take on any nonnegative real number; examples include the duration of parliamentary coalitions or time between coups d'etat. Note that the same underlying phenomena may be represented as either event counts (e.g., number of wars) or durations (time between wars), and some of the programs included in the CONSTRAINED COUNT module enable you to estimate exactly the same parameters with either form of data.

A variety of statistical models have been proposed to analyze events data, and the programs here provide some that I have developed, along with others I have found particularly useful in my research. I list here the specific programs included in this module, the models each program can estimate, and citations to the work for which I wrote each program. More complete references to the literature on event count and duration models appear at the end of this document.

CMLPoisson	Poisson regression (King, 1988, 1987), truncated Pois-
	son regression (1989d: Section 5), and log-linear and log-
	proportion models for contingency tables (1989a: Chapter
	6).
CMLNegbin	Negative binomial regression (1989b), truncated negative
	binomial regression (1989d: Section 5), truncated or un-
	truncated variance function models (1989d: Section 5),
	overdispersed log-linear and log-proportion models for con-
	tingency tables (1989a: Chapter 6).
CMLHurdlep	Hurdle Poisson regression model (1989d: Section 4).
CMLSupreme	Seemingly unrelated Poisson regression model (1989c).
CMLSupreme2	Poisson regression model with unobserved dependent vari-
	ables (1989d: Section 6).
CMLExpon	Exponential duration model with or without censoring
•	(King, Alt, Burns, and Laver, 1989).
CMLExpgam	Exponential-Gamma duration model with or without cen-
	soring (King, Alt, Burns, and Laver, 1989).
CMLPareto	Pareto duration model with or without censoring (King,
	Alt, Burns, and Laver, 1989).

4.1 Getting Started

GAUSS 3.1.0+ is required to use these routines.

4.1.1 README Files

The file **README.ccn** contains any last minute information on this module. Please read it before using the procedures in this module.

4.1.2 Setup

In order to use the procedures in the *CONSTRAINED COUNT* Module, the **CMLCount** library must be active. This is done by including count in the **LIBRARY** statement at the top of your program or command file:

4. CONSTRAINED EVENT COUNT AND DURATION REGRESSION

```
library cmlcount,quantal,pgraph;
```

This enables **GAUSS** to find the *CONSTRAINED COUNT* and required *CONSTRAINED MAXIMUM LIKELIHOOD* procedures. If you plan to make any right hand references to the global variables (which are described in a later section), you also need the statement:

```
#include cmlcount.ext;
```

To reset global variables in succeeding executions of the command file, the following instruction can be used:

```
cmlcountset;
```

This could be included with the above statements without harm and would insure the proper definition of the global variables for all executions of the command file.

The version number of each module is stored in a global variable. For the *CONSTRAINED COUNT* Module, this global is:

__cmlc__version 3×1 matrix, the first element contains the major version number of the *CONSTRAINED COUNT* Module, the second element the minor version number, and the third element the revision number.

If you call for technical support, you may be asked for the version number of your copy of this module.

4.2 About the CONSTRAINED COUNT Procedures

The format of the programs included in this module are all very similar:

```
{ b,vc,llik } = CMLExpon(dataset,dep,ind);
{ b,vc,llik } = CMLExpgam(dataset,dep,ind);
{ b,vc,llik } = CMLPareto(dataset,dep,ind);
{ b,vc,llik } = CMLPoisson(dataset,dep,ind);
{ b,vc,llik } = CMLNegbin(dataset,dep,ind1,ind2);
{ b,vc,llik } = CMLHurdlep(dataset,dep,ind1,ind2);
{ b,vc,llik } = CMLSupreme(dataset,dep1,dep2,ind1,ind2);
{ b,vc,llik } = CMLSupreme2(dataset,dep1,dep2,ind1,ind2,ind3);
```

An example program file looks like this:

```
library cmlcount;

CMLCountSet;

dep = { wars };
ind = { age, party, unem };
dataset = "sample";

_cml_A = { 0 1 -1 0 };
_cml_B = { 0 };
_cml_Bounds = { 0 10 };

call CMLPoisson(dataset,dep,ind);
```

This run constrains the coefficients of age and party to be equal, and bounds the coefficients to be positive and less than 10.

You may run these lines, or ones like them, from the **GAUSS** editor or interactively in command mode.

4.2.1 Inputs

The variable *dataset* is always the first argument. This may either be a matrix or a string containing the name of a **GAUSS** data set.

The dependent variable (or variables) is specified in each program by naming a symbol or a column number. For example,

```
dep = { durat };
or
dep = 7;
```

The independent variable vector (or vectors) is also specified in each program with variable names or column numbers. For example,

```
ind = { age, sex, race, height, size, iq };
or
ind = { 2, 4, 5, 6, 7 };
```

For each procedure, the data set and dependent variables must be specified. However, since constant terms are automatically included as part of independent variable vectors, you may occasionally wish to include no additional independent variables. You may do this easily by setting the relevant vector to zero. For example, ind = 0. For another example, you may wish to run the negative binomial regression model with a scalar dispersion parameter rather than a variance function: ind2 = 0.

4.2.2 Outputs

Printed output is controlled by the global **—output**, described in the section below. This section describes the outputs b, vc, and llik on the left hand side of the expressions above.

b vector, the maximum likelihood estimates for all the parameters. The mean vector comes first; the variance function, other mean vectors, and scalar dispersion parameters, if any, come next.

vc matrix, the variance-covariance matrix evaluated at the maximum. The standard errors are SQRT(DIAG(vc)). If you choose the CML global __cml_CovPar = 3, vc contains heteroskedastic-consistent parameter estimates.. See Section 2.7 for more discussion of options for statistical inference in constrained maximum likelihood models.

llik scalar, the value of the log-likelihood function at the maximum.

4.2.3 Global Control Variables

__cmlc__Inference string, determines the type of statistical inference.

BOOT generates bootstrapped estimates and covariance matrix of estimates

CML generates maximum likelihood estimates

Setting **_cmlc_Inference** to BOOT generates a **GAUSS** data set containing the bootstrapped parameters. The file name of this data set is either the default BOOTx, where x is a four digit number starting with 1 and increasing until a unique name is found, or the name in the **CML** global variable, **_cml_BootFname**. This data set can be used with **CMLBlimits** for generating confidence limits, with **CMLDensity** for generating density estimates and plots of the bootstrapped parameters, or with **CMLHist** for generating histogram and surface plots.

__cmlc_Censor scalar, allows you to include a variable indicating which observations are censored. This is used by the exponential, exponential-gamma, and Pareto models of duration data. Alternatively, you may set it to a symbol **__cmlc_Censor** = "varname" if you are using a **GAUSS** data set, or a number (**__cmlc_Censor** = 11) if the data set is a matrix in memory. The censoring variable should be 0 for censored observations and 1 for others.

By default, no observations are censored.

_cmlc_Fix scalar, name of index number of logged variable among the regressors with coefficient fixed to 1.0. By default, no logged variables are included.

In some of the programs, you have the option of including the log of a variable and fixing its coefficient to 1.0. To include the variable (the program takes the log), set **_cmlc_Fix** to a variable name or number $(_cmlc_Fix = "totals" \text{ or } _cmlc_Fix = 12).$ The default $(_cmlc_Fix = 12)$. 0) includes no additional variable. In most event count data, the observation period is the same length for all i (a year, month, etc.). However, in others, the observation period varies. For example, suppose one observed the number of times a citizen was contacted by a candidate in the interval between two public opinion polls; since polls typically take some time to administer, the observation period would vary over the individuals. In still other situations, the observation period may be the same length but the population of potential events varies. For example, if one observed the number of suicides per state, one would need some way to include information on differing state sizes in the analysis. It turns out that both of these situations can be dealt with in the same way by including an additional variable in the stochastic portion of the model. But (as explained in King, 1989, Section 5.8), this procedure turns out to be mathematically equivalent to including the log of this additional variable in the regression component, and constraining its coefficient to 1.0. There is often little harm in just including the log of this variable and estimating its coefficient with all the others, but several of the programs allow one to make this constraint.

- _cmlc_Dispersion scalar, set this to a value to change the starting value for only the dispersion parameter in the negative binomial (CMLNegbin), generalized event count (CMLHurdlep), exponential-gamma (CMLExpgam), Pareto (CMLPareto), and seemingly-unrelated Poisson models (CMLSupreme, CMLSupreme2). By default, a special starting value is not used for the dispersion parameter.
- **_cmlc_Precision** scalar, the number of digits printed to the right of the decimal point on output. Default = 4.
- **__cmlc_Start** scalar, selects method of calculating starting values. Possible values are:
 - o calculates them by regressing ln(y + 0.5) on the explanatory variables.
 - uses a vector of user supplied start values stored in the global variable **_cmlc_StartValue**.
 - 2 uses a vector of zeros.
 - 3 uses random uniform numbers on the interval $\left[-\frac{1}{2}, \frac{1}{2}\right]$.

Default = 0.

4. CONSTRAINED EVENT COUNT AND DURATION REGRESSION

_cmlc_StartValue L×1 vector, start values if **_cmlc_Start** = 1.

__cmlc_ZeroTruncate scalar, specifies whether or not the model is a truncated model. For the Poisson and negative binomial models, **__cmlc_ZeroTruncate** = 0 estimates a truncated-at-zero version of the model. By default, the regular untruncated model is estimated.

—altnam K×1 vector, alternate names for variables when a matrix is passed to a CMLCount procedure. When a data matrix is passed to a CMLCount procedure and the user is selecting from that matrix, the global variable —altnam, if it is used, must contain names for the columns of the original matrix.

__output scalar, determines printing of intermediate results.

- **0** nothing is written.
- 1 serial ASCII output format suitable for disk files or printers.
- 2 (DOS only) output is suitable for screen only. ANSI.SYS must be active.

Default = 2.

___row scalar, specifies how many rows of the data set are read per iteration of the read loop. By default, the number of rows to be read is calculated automatically.

__rowfac scalar, row factor. If a CONSTRAINED COUNT procedure fails due to insufficient memory while attempting to read a GAUSS data set, then __rowfac may be set to some value between 0 and 1 to read a proportion of the original number of rows of the GAUSS data set. For example, setting

 $_{rowfac} = 0.8;$

causes **GAUSS** to read in 80% of the rows originally calculated.

This global has an affect only when $__row = 0$.

Default = 1.

__title string, message printed at the top of the screen and printed out by CMLCountPrt. Default = "".

__vpad scalar, if *dataset* is a matrix in memory, the variable names are automatically created by **CML**. Two types of names can be created:

Variable names automatically created by **CML** are not padded to give them equal length. For example, V1, V2,...V10, V11,....

Variable names created by the procedure are padded with zeros to give them an equal number of characters. For example, V01, V02, ..., V10, V11,.... This is useful if you want the variable names to sort properly.

Default = 1.

4.2.4 Adding Constraints

There are two general types of constraints, nonlinear equality constraints and nonlinear inequality constraints. However, for computational convenience they are divided into five types: linear equality, nonlinear equality, linear inequality, nonlinear inequality, and bounds. For a discussion of specifying constraints, see Section 2.5.

The specification of constraints requires knowledge of the order of the parameters. For all models, the first parameter is a constant term, then one parameter for each explanatory variable, and then a dispersion parameter. For **CMLHurdlep** and **CMLSupreme2** another constant term and set of explanatory parameters follows the dispersion parameter. For example, suppose there are four explanatory variables, and you wish to constrain the coefficients and the dispersion parameter to be positive:

```
_cml_Bounds = { 0 1e200 };
```

To constrain the coefficients of the first two explanatory variables to be equal:

```
_cml_A = { 0 1 -1 0 0 0 };
_cml_B = { 0 };
```

To constrain the norm of the coefficients of the explanatory variables to be greater than 2:

```
proc eqp(b);
    local c;
    c = b[2:4];
    retp(c'c - 2);
endp;
_cml_EqProc = &eqp;
```

4.2.5 Statistical Inference

CML statistical inference features may be accessed through the **COUNT** global, **_cmlc_Inference**. **_cmlc_Inference** has the following settings:

4. CONSTRAINED EVENT COUNT AND DURATION REGRESSION

CML	constrained maximum likelihood estimates (default)
BOOT	bootstrapped estimates

That is to generate bootstrapped estimates, set

```
_cmlc_Inference = "boot";
```

Confidence limits for inequality constrained parameters are generated by first leaving **_cmlc_Inference** at its default setting and then calling **CMLClimits** with the covariance matrix of the parameters and the parameter estimates as arguments.

Confidence Limits of Constrained Parameters

When inequality constraints are present, the considerations discussed in Section 2.7 are relevant. You may need to use **CMLClimits** for correct confidence limits. In this case, use the covariance matrix and estimates returned from the **CMLCount** procedure as input to **CMLClimits** as well as any globals and procedures used for the constraints. For example,

Bootstrapping

In addition to the usual standard errors, you may generate bootstrap standard errors. Setting **_cmlc_Inference** to BOOT causes **CMLCOUNT** to call **CMLBoot**. This generates bootstrapped estimates and covariance matrices of the estimates.

The bootstrapped parameters are also stored in a **GAUSS** data set. The name of the data set can be determined by setting **_cml_BootFname** to a file name, or by default it will be set to BOOTx where x is a four digit number incremented from 0001 until a name not in use is found. For further details about the bootstrap, see Section 2.7.5.

The data set thus generated can be used for computing confidence intervals of the coefficients using **CMLBlimits**. Also, density estimates and plots can be generated using **CMLDensity**, and histograms and surface plots of the coefficients can be produced using **CMLHist**. For further details about **CMLDensity**, see Section 2.7.5, and for further details about **CMLHist** see Section 2.7.5.

4.2.6 Problems with Convergence

All the programs use maximum likelihood estimation by numerically maximizing a different likelihood function. As with virtually all nonlinear iterative procedures, convergence works most of the time, but not every time. Problems to be aware of include the following:

- 1. The explanatory variables in each regression function should not be highly collinear among themselves.
- 2. The model should have more observations than parameters; indeed, the more observations, the better.
- 3. Starting values should not be too far from the optimal values.
- 4. The model specified should fit the data.
- 5. The Poisson hurdle model must have at least some observations with $y_i = 0$ and should take on at least two other values greater than zero.
- 6. The truncated models should have no observations with zeros (if inadvertently included, a message appears and the program stops).
- 7. The models with scalar dispersion parameters and variance functions should have maximum likelihood estimates that are bounded so that, for example, in the negative binomial model $\hat{\sigma}^2 > 1$

If you avoid the potential problems listed in the last paragraph, you should have little problem with convergence. Of course, avoiding these problems with difficult data sets is not always easy nor obvious. In these cases, problems may be indicated by the following situations:

1. iterations sending the parameters off in unreasonable directions or creating very large numbers.

4. CONSTRAINED EVENT COUNT AND DURATION REGRESSION

- 2. the program actually bombing out.
- 3. a single iteration taking an extraordinarily long time.
- 4. the program taking more than 40 or 50 iterations with no convergence.

If one of these problems occur, you have several options. First, look over the list in the last paragraph. To verify that the problem does indeed exist, you might try running your data on the Poisson regression model if you have event count data, or the exponential regression model if you have duration data. Both are known to be globally concave and tend to converge very easily. If this model works, but another does not, you probably do have a problem.

In the case of problems, you must consider iteration a participatory process. When **CML** is iterating, you can press **Alt-H** to receive a list of options that may be changed during iteration. See *CML REFERENCE* for a full explanation of each. I find that the following practices tend to work well:

- 1. If the program has produced many iterations without much progress, try pressing Alt-I every few iterations to force the program to calculate the information matrix or switch Newton-Raphson iterations. Either of these may not work if the iterations are not far enough along.
- 2. The number of zeros to the right of the decimal point on the relative gradients (printed on the screen while the program is iterating) is the approximate precision of your final estimates. If the program is having trouble converging, but the gradients are small enough (i.e., you have sufficient precision for your substantive problem), press Alt-C to force the program to converge.
- 3. If the program bombs out very quickly, changing the starting values are your best bet (with the global **_cmlc_Start**). The default starting values created with least squares, **_cmlc_Start** = 0, usually works best. If that does not work, you can also try creating them yourself, by thinking about what the answer is likely to be or by running a simpler model. For example, the exponential-gamma model is sometimes problematic; however, the exponential model often provides good starting values for the effect parameters. Thus if the other methods do not work, you might try the following:

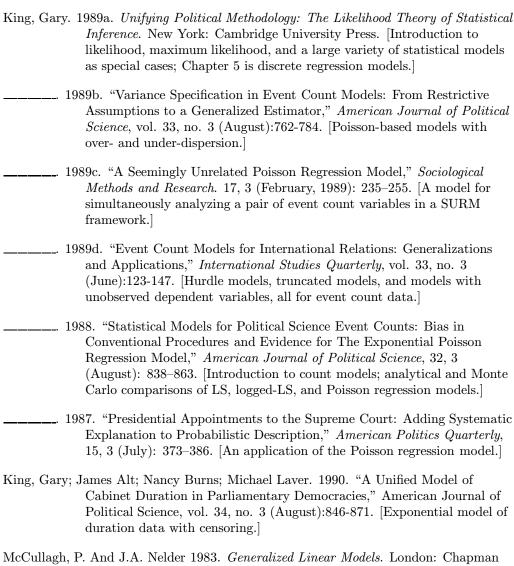
```
library cmlcount;
CMLCountSet;
dep = { durat };
ind = { unem, infl, age };
dataset = "datafile";
{ b,vc,llik } = CMLExpon(dataset,dep,ind);
_cmlc_StartValue = b;
_cmlc_Start = 1;
call CMLExpgam(dataset,dep,ind);
```

You can also choose one of the other methods of creating starting values by changing the **_cmlc_Start** global (described above).

4.3 Annotated Bibliography

- Allison, Paul. 1984. *Event History Analysis*. Beverly Hills: Sage. [A simple overview of event history methods for duration data.]
- Bishop, Yvonne M.M.; Stephen E. Fienberg; and Paul W. Holland. 1975. *Discrete Multivariate Analysis* Cambridge, Mass.: M.I.T. Press. [Models for contingency tables.]
- Cameron, A. Colin and Pravin K. Trivedi. 1986. "Econometric Models Based on Count Data: Comparisons and Applications of Some Estimators and Tests," Journal of Applied Econometrics 1, 29–53. [Review of the econometric literature on event counts.]
- Grogger, Jeffrey T. and Richard T. Carson. 1988. "Models for Counts from Choice Based Samples," Discussion Paper 88-9, Department of Economics, University of California, San Diego. [Truncated event count models.]
- Gourieroux, C.; A. Monfort; and A. Trognon. 1984. "Pseudo Maximum Likelihood Methods: Applications to Poisson Models," *Econometrica* 52: 701–720. [A three-stage robust estimation method for count data.]
- Hall, Bronwyn H.; Zvi Griliches; and Jerry A. Hausman. 1986. "Patents and R and D: Is there a Lag?" *International Economic Review.* 27, 2 (June): 265–83. [Nice example of a applying a variety of different estimators to single equation count models.]
- Hausman, Jerry; Bronwyn H. Hall; and Zvi Griliches. 1984. "Econometrics Models for Count Data with An Application to the Patents-R&D Relationship," *Econometrica.* 52, 4 (July): 909-938. [Count models for time-series cross sectional panels.]
- Holden, Robert T. 1987. "Time Series Analysis of a Contagious Process," *Journal of the American Statistical Association*. 82, 400 (December): 1019–1026. [A time series model of count data applied to airline hijack attempts.]
- Jorgenson, Dale W. 1961. "Multiple Regression Analysis of a Poisson Process," *Journal of the American Statistical Association* 56,294 (June): 235–45. [The Poisson regression model.]
- Kalbfleisch, J.D. and R.L. Prentice. 1980. The Statistical Analysis of Failure Time Data. New York: Wiley. [Summary of research on many models of duration data.]

4. CONSTRAINED EVENT COUNT AND DURATION REGRESSION



- and Hall. [A unified approach to specifying and estimating this class of models. Some count and duration models are covered.]
- Mullahy, John. 1986. "Specification and Testing of Some Modified Count Data Models," *Journal of Econometrics*. 33: 341–65. [Several hurdle-type models of event count data.]
- Tuma, Nancy Brandon and Michael T. Hannan. 1984. *Social Dynamics*. New York: Academic Press.

$4.\ CONSTRAINED\ EVENT\ COUNT\ AND\ DURATION\ REGRESSION$

Chapter 5

CMLCount Reference

Formats and prints the output from calls to CONSTRAINED COUNT procedures.

Library

cmlcount

Format

```
{ b,vc,llik } = CMLCountPrt(b,vc,llik);
```

Input

b (K+1)×1 vector, maximum likelihood estimates of the effect parameters stacked on top of the dispersion parameter.

vc (K+1)×(K+1) matrix, variance-covariance matrix

llik scalar, value of the log-likelihood function at the maximum.

Output

The input arguments are returned unchanged.

Remarks

The call to $CONSTRAINED\ COUNT$ procedures can be nested in the call to the **CMLCountPrt**:

```
{ b,vc,llik } = cmlcountprt(CMLExpgam(dataset,dep,ind));
```

Formats and prints the output from calls to $CONSTRAINED\ COUNT$ procedures with confidence limits

Library

cmlcount

Format

```
\{b,cl,llik\} = CMLCountCLPrt(b,cl,llik);
```

Input

```
b \hspace{1cm} (K+1)\times 1 \hspace{0.1cm} \text{vector, maximum likelihood estimates of the effect parameters} \\ vc \hspace{1cm} (K+1)\times 2 \hspace{0.1cm} \text{matrix, confidence limits} \\ llik \hspace{1cm} \text{scalar, value of the log-likelihood function at the maximum.} \\
```

Output

The input arguments are returned unchanged.

Remarks

Confidence limits computed by **CMLBlimits**, **CMLClimits**, or **CMLTlimits** may be passed in the fourth argument in the call to **CMLCountCLPrt**:

```
{ b,vc,llik } = CMLExpgam(dataset,dep,ind);
cl = CMLBlimits(_cml_BootFname);
call CMLCountCLPrt(b,cl,llik);
```

Source

ccount.src

Resets CONSTRAINED COUNT global variables to default values.

Library

cmlcount

Format

CMLCountSet;

Input

None

Output

None

Remarks

Putting this instruction at the top of all command files that invoke CONSTRAINED COUNT procedures is generally good practice. This prevents globals from being inappropriately defined when a command file is run several times or when a command file is run after another command file has executed that calls a CONSTRAINED COUNT procedure.

CMLCountSet calls CMLSet which calls GAUSSET.

Source

ccount.src

Estimates an exponential-gamma regression model, for the analysis of duration data, with maximum likelihood.

Library

cmlcount

■ Format

```
\{b,vc,llik\} = CMLExpgam(dataset,dep,ind);
```

Input

dataset string, name of **GAUSS** data set.

– or –

N×K matrix, data

dep string, the name of the dependent variable.

- or -

scalar, the index of the dependent variable.

ind K×1 character vector, names of the independent variables.

– or ·

K×1 numeric vector, indices of independent variables.

Set to 0 to include only a constant term.

If *dataset* is a matrix, *dep* or *ind* may be a string or character variable containing either the standard labels created by **CML** (V1, V2,..., or V01, V02,...., depending on the value of **__vpad**), or the user-provided labels in **__altnam**.

Output

b (K+1)×1 vector, maximum likelihood estimates of the effect parameters

stacked on top of the dispersion parameter.

vc (K+1)×(K+1) matrix, variance-covariance matrix of the estimated

parameters evaluated at the maximum. If you choose the CML global $_cml_CovPar = 3,\ vc$ contains heteroskedastic-consistent parameter

estimates.

llik scalar, value of the log-likelihood function at the maximum.

Globals

CML globals are also relevant, including constraint matrices and procedures.

__cmlc__Inference string, determines the type of statistical inference.

BOOT generates bootstrapped estimates and covariance matrix of estimates

CML generates maximum likelihood estimates

Setting **__cmlc__Inference** to BOOT generates a **GAUSS** data set containing the bootstrapped parameters. The file name of this data set is either the default BOOTx, where x is a four digit number starting with 1 and increasing until a unique name is found, or the name in the **CML** global variable, **__cml__BootFname**. This data set can be used with **CMLBlimits** for generating confidence limits, with **CMLDensity** for generating density estimates and plots of the boostrapped parameters, or with **CMLHist** for generating histogram and surface plots.

_cmlc_Censor string, the name of the censor variable from *dataset*.

– or –

scalar, the index of the censor variable from dataset.

By default, no censoring is used.

_cmlc_Start scalar, selects method of calculating starting values. Possible values are:

- ocalculates them by regressing ln(y + 0.5) on the explanatory variables.
- 1 uses a vector of user supplied start values stored in the global variable **_cmlc_StartValue**.
- **2** uses a vector of zeros.
- 3 uses random uniform numbers on the interval $\left[-\frac{1}{2}, \frac{1}{2}\right]$.

Default = 0.

_cmlc_StartValue $(K+1)\times 1$ vector, start values if **_cmlc_Start** = 1.

_cmlc_Precision scalar, number of decimal points to print on output. Default = 4.

__altnam K×1 vector, alternate names for variables when a matrix is passed to **CMLExpgam**. When a data matrix is passed to **CMLExpgam** and when the user is selecting from that matrix, the global variable **__altnam**, if it is used, must contain names for the columns of the original matrix.

__miss scalar, determines how missing data will be handled.

- **0** Missing values will not be checked for, and so the data set must not have any missings. This is the fastest option.
- 1 Listwise deletion. Removes from computation any observation with a missing value on any variable included in the analysis.

Default = 0.

__output scalar, determines printing of intermediate results.

- **0** nothing is written.
- 1 serial ASCII output format suitable for disk files or printers.
- 2 (DOS only) output is suitable for screen only. ANSI.SYS must be active.

Default = 2.

___row scalar, specifies how many rows of the data set will be read per iteration of the read loop. By default, the number of rows to be read will be calculated automatically.

___rowfac scalar, row factor. If **CMLExpgam** fails due to insufficient memory while attempting to read a **GAUSS** data set, then **___rowfac** may be set to some value between 0 and 1 to read a *proportion* of the original number of rows of the **GAUSS** data set. For example, setting

 $_{rowfac} = 0.8;$

will cause ${\sf GAUSS}$ to read in 80% of the rows originally calculated.

This global has an affect only when $__row = 0$.

Default = 1.

___title string, message printed at the top of the screen and printed out by CMLCountPrt. Default = "".

___vpad scalar, if *dataset* is a matrix in memory, the variable names are automatically created by **CML**. Two types of names can be created:

- Variable names automatically created by **CML** are not padded to give them equal length. For example, V1, V2,...V10, V11,....
- Variable names created by the procedure are padded with zeros to give them an equal number of characters. For example, V01, V02, ..., V10, V11,.... This is useful if you want the variable names to sort properly.

Default = 1.

Remarks

Let the *n* duration observations (nonnegative real numbers) for the dependent variable be denoted as y_1, \ldots, y_n . Assume that y_i follows a gamma distribution with expected

value μ_i and variance $\mu_i^2 \sigma^2$. Let the mean μ_i be an exponential-linear function of a vector of explanatory variables, x_i :

$$E(y_i) \equiv \mu_i = \exp(x_i \beta) \tag{5.1}$$

The program includes a constant term as the first column of x_i and allows one to include any number of explanatory variables. Note that μ_i from a duration model equals $1/\lambda_i$ from an event count model; thus, one need only change the sign of the effect parameters to get estimates of the same parameters from these different kinds of data.

The dispersion σ^2 is parametrized as follows:

$$\sigma_i^2 = \exp(\gamma) \tag{5.2}$$

EXPGAM reports estimates of β and γ .

For an introduction to the exponential gamma regression model see King, Alt, Burns, and Laver (1989) or Kalbfleisch and Prentice (1980).

Example

Constrained Exponential-Gamma Regression Model of Duration Data

A vector of effect parameters and a scalar dispersion parameter are estimated. The vector includes one element corresponding to each explanatory variable named in *ind* and a constant term. Five parameters are estimated in this example.

Constrained Censored Exponential-Gamma Regression Model of Duration Data

```
library cmlcount;
#include cmlcount.ext;
CMLCountset;
```

```
dataset = "wars";
dep = { wars };
ind = { unem, poverty, allianc };
_Censor = { v12 };
{ b,vc,llik } = CMLExpgam(dataset,dep,ind);
output file = cmlcount.out reset;
call CMLCountPrt(b,vc,llik);
output off;
```

A vector of effect parameters and a scalar dispersion parameter are estimated. The vector includes one element corresponding to each explanatory variable named in *ind* and a constant term. Five parameters are estimated in this example.

Source

cmlexpgm.src

Estimates a constrained exponential regression model or censored exponential regression model with maximum likelihood.

Library

cmlcount

Format

```
\{b,vc,llik\} = CMLExpon(dataset,dep,ind);
```

Input

dataset string, name of **GAUSS** data set.

– or –

N×K matrix, data

dep string, the name of the dependent variable

- or -

scalar, the index of the dependent variable

ind K×1 character vector, names of the independent variables

– or –

 $K \times 1$ numeric vector, indices of independent variables

Set to 0 to include only a constant term.

If dataset is a matrix, dep or ind may be a string or character variable containing either the standard labels created by **CML** (V1, V2,..., or V01, V02,..., depending on the value of **__vpad**), or the user-provided labels in **__altnam**.

Output

b K×1 vector, maximum likelihood estimates of the effect parameters.

vc K×K matrix, variance-covariance matrix of the estimated parameters

evaluated at the maximum. If the CML global _cml_CovPar is set to 3,

vc will contain heteroskedastic-consistent parameter estimates.

llik scalar, value of the log-likelihood function at the maximum.

Globals

CML globals are also relevant, including constraint matrices and procedures.

__cmlc__Inference string, determines the type of statistical inference.

BOOT generates bootstrapped estimates and covariance matrix of estimates

CML generates maximum likelihood estimates

Setting **__cmlc__Inference** to BOOT generates a **GAUSS** data set containing the bootstrapped parameters. The file name of this data set is either the default BOOTx, where x is a four digit number starting with 1 and increasing until a unique name is found, or the name in the **CML** global variable, **__cml__BootFname**. This data set can be used with **CMLBlimits** for generating confidence limits, with **CMLDensity** for generating density estimates and plots of the boostrapped parameters, or with **CMLHist** for generating histogram and surface plots.

_cmlc_Censor string, the name of the censor variable from dataset

scalar, the index of the censor variable from dataset

By default, no censoring is used.

_cmlc_Start scalar, selects method of calculating starting values. Possible values are:

- o calculates them by regressing $\ln(y+0.5)$ on the explanatory variables.
- will use a vector of user supplied start values stored in the global variable **_cmlc_StartValue**.
- **2** uses a vector of zeros.
- 3 uses random uniform numbers on the interval $\left[-\frac{1}{2}, \frac{1}{2}\right]$.

Default = 0.

_cmlc_StartValue $K \times 1$ vector, start values if **_cmlc_Start** = 1.

_cmlc_Precision scalar, number of decimal points to print on output. Default = 4.

__altnam K×1 vector, alternate names for variables when a matrix is passed to **CMLExpon**. When a data matrix is passed to **CMLExpon** and the user is selecting from that matrix, the global variable **__altnam**, if it is used, must contain names for the columns of the original matrix.

__miss scalar, determines how missing data will be handled.

- **0** Missing values will not be checked for, and so the data set must not have any missings. This is the fastest option.
- 1 Listwise deletion. Removes from computation any observation with a missing value on any variable included in the analysis.

Default = 0.

__output scalar, determines printing of intermediate results.

- **0** nothing is written.
- 1 serial ASCII output format suitable for disk files or printers.
- 2 (DOS only) output is suitable for screen only. ANSI.SYS must be active.

Default = 2.

___row scalar, specifies how many rows of the data set will be read per iteration of the read loop. By default, the number of rows to be read will be calculated automatically.

___rowfac scalar, row factor. If **EXPON** fails due to insufficient memory while attempting to read a **GAUSS** data set, then **___rowfac** may be set to some value between 0 and 1 to read a *proportion* of the original number of rows of the **GAUSS** data set. For example, setting

 $_{rowfac} = 0.8;$

will cause **GAUSS** to read in 80% of the rows originally calculated.

This global only has an affect when $__row = 0$.

Default = 1.

___title string, message printed at the top of the screen and printed out by **CMLCountPrt**. Default = "".

__vpad scalar, if *dataset* is a matrix in memory, the variable names are automatically created by **CML**. Two types of names can be created:

- Variable names automatically created by **CML** are not padded to give them equal length. For example, V1, V2,...V10, V11,....
- Variable names created by the procedure are padded with zeros to give them an equal number of characters. For example, V01, V02, ..., V10, V11,.... This is useful if you want the variable names to sort properly.

Default = 1.

Remarks

Let y_i $(i=1,\ldots,n)$ take on any non-negative real number representing a duration. Often y_i is only measured as an integer, such as the number of days or months. Even so, if your dependent variable is a measure of time, duration models, and not event count models, are called for. Let y_i be distributed exponentially with mean μ_i . Also let

 $E(y_i) \equiv \mu_i = \exp(x_i\beta)$. Note that μ_i from a duration model equals $1/\lambda_i$ from an event count model; thus, one need only change the sign of the effect parameters to get estimates of the same parameters from these different kinds of data.

For an introduction to the exponential regression model and the censored exponential regression model see Kalbfleisch and Prentice (1980) and King, Alt, Burns, and Laver (1989).

Example

Constrained Exponential Regression Model

A single vector of effect parameters are estimated. This vector includes one element corresponding to each explanatory variable named in *ind* and a constant term.

Constrained Censored Exponential Regression Model

```
library cmlcount;
#include cmlcount.ext;

CMLCountset;
dataset = "wars";
dep = { wars };
ind = { unem, poverty, allianc };
_cmlc_Censor = { notseen };
{ b,vc,llik } = CMLExpon(dataset,dep,ind);
output file = cmlcount.out reset;
call CMLCountPrt(b,vc,llik);
output off;
```

A single vector of effect parameters are estimated. This vector includes one element corresponding to each explanatory variable named in ind and a constant term.

Source

cmlexpon.src

Estimates a constrained hurdle Poisson regression model, for the analysis of event counts, with maximum likelihood.

Library

cmlcount

Format

```
\{b, vc, llik\} = CMLHurdlep(dataset, dep, ind);
```

Input

dataset string, name of GAUSS data set.

– or –

N×K matrix, data

dep string, the name of the dependent variable

– or -

scalar, the index of the dependent variable

ind1 K×1 character vector, names of first event independent variables

- or -

 $K \times 1$ numeric vector, indices of first event independent variables

ind2 K×1 character vector, names of second event independent variables

- or -

K×1 numeric vector, indices of second event independent variables

If dataset is a matrix, dep, ind1, or ind2 may be a string or character variable containing either the standard labels created by **CML** (V1, V2,..., or V01, V02,...., depending on the value of **__vpad**), or the user-provided labels in **__altnam**.

Output

b $(\mathrm{K}{+}\mathrm{L}){\times}1$ vector, maximum likelihood estimates of the effect parameters

stacked on top of the dispersion parameter.

vc (K+L)×(K+L) matrix, variance-covariance matrix of the estimated parameters evaluated at the maximum. If you choose the **CML** global **_cml_CovPar** = 3, vc will contain heteroskedastic-consistent parameter

estimates.

llik scalar, value of the log-likelihood function at the maximum.

CMLCount Reference

Globals

CML globals are also relevant, including constraint matrices and procedures.

_cmlc_Inference string, determines the type of statistical inference.

BOOT generates bootstrapped estimates and covariance matrix of estimates

CML generates maximum likelihood estimates

Setting **__cmlc__Inference** to BOOT generates a **GAUSS** data set containing the bootstrapped parameters. The file name of this data set is either the default BOOTx, where x is a four digit number starting with 1 and increasing until a unique name is found, or the name in the **CML** global variable, **__cml__BootFname**. This data set can be used with **CMLBlimits** for generating confidence limits, with **CMLDensity** for generating density estimates and plots of the boostrapped parameters, or with **CMLHist** for generating histogram and surface plots.

_cmlc_Start scalar, selects method of calculating starting values. Possible values are:

- ocalculates them by regressing ln(y + 0.5) on the explanatory variables.
- will use a vector of user supplied start values stored in the global variable **_cmlc_StartValue**.
- 2 uses a vector of zeros.
- 3 uses random uniform numbers on the interval $\left[-\frac{1}{2}, \frac{1}{2}\right]$.

Default = 0.

_cmlc_StartValue $(K+L)\times 1$ vector, start values if **_cmlc_Start** = 1.

_cmlc_Precision scalar, number of decimal points to print on output. Default = 4.

___altnam K×1 vector, alternate names for variables when a matrix is passed to **CMLHurdlep**. When a data matrix is passed to **CMLHurdlep** and the user is selecting from that matrix, the global variable **___altnam**, if it is used, must contain names for the columns of the original matrix.

__miss scalar, determines how missing data will be handled.

- **0** Missing values will not be checked for, and so the data set must not have any missings. This is the fastest option.
- 1 Listwise deletion. Removes from computation any observation with a missing value on any variable included in the analysis.

Default = 0.

__output scalar, determines printing of intermediate results.

- **0** nothing is written.
- 1 serial ASCII output format suitable for disk files or printers.
- 2 (DOS only) output is suitable for screen only. ANSI.SYS must be active.

Default = 2.

__row

scalar, specifies how many rows of the data set will be read per iteration of the read loop. By default, the number of rows to be read will be calculated automatically.

__rowfac

scalar, row factor. If **CMLHurdlep** fails due to insufficient memory while attempting to read a **GAUSS** data set, then **___rowfac** may be set to some value between 0 and 1 to read a *proportion* of the original number of rows of the **GAUSS** data set. For example, setting

```
_{rowfac} = 0.8;
```

will cause **GAUSS** to read in 80% of the rows originally calculated.

This global only has an affect when $__row = 0$.

Default = 1.

__title

string, message printed at the top of the screen and printed out by **CMLCountPrt**. Default = "".

__vpad

scalar, if *dataset* is a matrix in memory, the variable names are automatically created by **CML**. Two types of names can be created:

- Variable names automatically created by **CML** are not padded to give them equal length. For example, V1, V2,...V10, V11,....
- Variable names created by the procedure are padded with zeros to give them an equal number of characters. For example, V01, V02, ..., V10, V11,.... This is useful if you want the variable names to sort properly.

Default = 1.

Remarks

Let the n event count observations (nonnegative integers) for the dependent variable be denoted as y_1, \ldots, y_n . y_i is then a random dependent variable representing the number of events that have occurred during observation period i. Let λ_{0i} be the rate of the first event occurrence and λ_{+i} be the rate for all additional events after the first. If these are

the expected values of two separate Poisson processes, we have the hurdle Poisson regression model. These means are parametrized as usual:

$$\lambda_{0i} = \exp(x_i \beta) \tag{5.3}$$

and

$$\lambda_{+i} = \exp(z_i \gamma) \tag{5.4}$$

where x_i and z_i are (possibly) different vectors of explanatory variables. The program produces estimates of β and γ . If $\beta = \gamma$ and x = z, this model reduces to the Poisson.

For an introduction to the Hurdle Poisson regression model see Mullahy (1986) and King (1989d).

Example

Constrained Hurdle Poisson Regression Model:

Two vectors of effect parameters are estimated. Each includes one element corresponding to each explanatory variable plus a constant term (in the example, four parameters appear in the first regression function and seven in the second).

Source

cmlhurdl.src

Estimates a constrained negative binomial regression model or truncated-at-zero negative binomial regression model with maximum likelihood.

Library

cmlcount

Format

```
\{b,vc,llik\} = CMLNegbin(dataset,dep,ind1,ind2);
```

Input

dataset string, name of **GAUSS** data set.

- or -

N×K matrix, data

dep string, the name of the dependent variable

– or –

scalar, the index of the dependent variable

ind1 K×1 character vector, names of first event independent variables

– or –

 $K\times 1$ numeric vector, indices of first event independent variables

Set to 0 to include only a constant term.

ind2 K×1 character vector, names of second event independent variables

– or –

K×1 numeric vector, indices of second event independent variables

Set to 0 for a scalar dispersion parameter.

If dataset is a matrix, dep, ind1, or ind2 may be a string or character variable containing either the standard labels created by **CML** (V1, V2,..., or V01, V02,...., depending on the value of **___vpad**), or the user-provided labels in **___altnam**.

Output

b $(K+1)\times 1$ or $(K+L)\times 1$ vector, maximum likelihood estimates of the effect

parameters stacked on top of either the dispersion parameter or the

coefficients of the variance function.

vc $(K+1)\times(K+1)$ or $(K+L)\times(K+L)$ matrix, variance-covariance matrix of

the estimated parameters evaluated at the maximum. If you choose the

CML global $_$ cml $_$ CovPar = 3, vc will contain heteroskedastic-consistent parameter estimates.

llik scalar, value of the log-likelihood function at the maximum.

Globals

CML globals are also relevant, including constraint matrices and procedures.

__cmlc__Inference string, determines the type of statistical inference.

BOOT generates bootstrapped estimates and covariance matrix of estimates

CML generates maximum likelihood estimates

Setting **__cmlc__Inference** to BOOT generates a **GAUSS** data set containing the bootstrapped parameters. The file name of this data set is either the default BOOTx, where x is a four digit number starting with 1 and increasing until a unique name is found, or the name in the **CML** global variable, **__cml__BootFname**. This data set can be used with **CMLBlimits** for generating confidence limits, with **CMLDensity** for generating density estimates and plots of the boostrapped parameters, or with **CMLHist** for generating histogram and surface plots.

_cmlc_Fix scalar, name of index number of logged variable among the regressors with coefficient constrained to 1.0 By default, no logged variables are included.

_cmlc_Start scalar, selects method of calculating starting values. Possible values are:

- ocalculates them by regressing ln(y + 0.5) on the explanatory variables.
- will use a vector of user supplied start values stored in the global variable _cmlc_StartValue.
- **2** uses a vector of zeros.
- 3 uses random uniform numbers on the interval $\left[-\frac{1}{2}, \frac{1}{2}\right]$.

Default = 0.

- **_cmlc_StartValue** $(K+1)\times 1$ or $(K+L)\times 1$ vector, start values if **_cmlc_Start** = 1.
- **_cmlc_Dispersion** scalar, start value for scalar dispersion parameter. Default = 3.
- **_cmlc_Precision** scalar, number of decimal points to print on output. Default = 4.
- **_cmlc_ZeroTruncate** scalar, specifies which model is used:
 - 0 truncated-at-zero negative binomial model
 - 1 negative binomial model is used.

5. CMLCOUNT REFERENCE

__altnam

 $K\times 1$ vector, alternate names for variables when a matrix is passed to **CMLNegbin**. When a data matrix is passed to **CMLNegbin** and the user is selecting from that matrix, the global variable **__altnam**, if it is used, must contain names for the columns of the original matrix.

__miss

scalar, determines how missing data will be handled.

- **0** Missing values will not be checked for, and so the data set must not have any missings. This is the fastest option.
- 1 Listwise deletion. Removes from computation any observation with a missing value on any variable included in the analysis.

Default = 0.

__output

scalar, determines printing of intermediate results.

- **0** nothing is written.
- 1 serial ASCII output format suitable for disk files or printers.
- 2 (DOS only) output is suitable for screen only. ANSI.SYS must be active.

Default = 2.

__row

scalar, specifies how many rows of the data set will be read per iteration of the read loop. By default, the number of rows to be read will be calculated automatically.

__rowfac

scalar, row factor. If **CMLNegbin** fails due to insufficient memory while attempting to read a **GAUSS** data set, then **__rowfac** may be set to some value between 0 and 1 to read a *proportion* of the original number of rows of the **GAUSS** data set. For example, setting

```
_{-}rowfac = 0.8;
```

will cause **GAUSS** to read in 80% of the rows originally calculated.

This global only has an affect when $__row = 0$.

Default = 1.

__title

string, message printed at the top of the screen and printed out by $\mathbf{CMLCountPrt}$. Default = "".

__vpad

scalar, if dataset is a matrix in memory, the variable names are automatically created by ${\sf CML}$. Two types of names can be created:

Variable names automatically created by **CML** are not padded to give them equal length. For example, V1, V2,...V10, V11,....

Variable names created by the procedure are padded with zeros to give them an equal number of characters. For example, V01, V02, ..., V10, V11,.... This is useful if you want the variable names to sort properly.

Default = 1.

Remarks

Let y_i be a random dependent variable representing the number of events that have occurred during observation period i (i = 1, ..., n). Assume that y_i follows a negative binomial distribution with expected value λ_i and variance $\lambda_i \sigma^2$. Let the mean λ_i (the rate of event occurrence, which must be greater than zero) be an exponential-linear function of a vector of explanatory variables, x_i :

$$E(y_i) \equiv \lambda_i = \exp(x_i \beta) \tag{5.5}$$

The program includes a constant term as the first column of x_i and allows one to include any number of explanatory variables.

 σ^2 is parametrized as follows:

$$\sigma_i^2 = 1 + \exp(z_i \gamma) \tag{5.6}$$

where $z_i = 1$, if estimating a scalar dispersion parameter, or a vector of explanatory variables, if estimating a variance function. The program calculates estimates of β and γ .

For an introduction to the negative binomial regression model, see Hausman, Hall, and Griliches (1984) and King (1989b); for information on the truncated negative binomial model, see Grogger and Carson (1988), and on the variance function model with or without truncation see King (1989d: Section 5)

Example

Constrained Negative Binomial Regression Model

A single vector of effect parameters and one scalar dispersion parameter are estimated. The vector of effect parameters includes one element corresponding to each explanatory variable and a constant term. In the example, five parameters are estimated.

Constrained Negative Binomial Variance Function Regression Model

```
library cmlcount;
#include cmlcount.ext;

CMLCountset;
dataset = "wars";
dep1 = { wars };
ind1 = { unem, poverty, allianc };
ind2 = { partyid, x4 };
{ b,vc,llik } = CMLNegbin(dataset,dep,ind1,ind2);
output file = cmlcount.out reset;
call CMLCountPrt(b,vc,llik);
output off;
```

Two vectors of effect parameters are estimated, one for the mean ind1 and one for the variance function ind2. Each vector includes a constant term and one element corresponding to each explanatory variable. The example estimates seven parameters.

Constrained Truncated-at-zero Negative Binomial Regression Model

```
library cmlcount;
#include cmlcount.ext;

CMLCountset;
dataset = "wars";
dep1 = { wars };
ind1 = { unem, poverty, allianc };
_cmlc_ZeroTruncate = 0;
{ b,vc,llik } = CMLNegbin(dataset,dep,ind1,0);
output file = cmlcount.out reset;
call CMLCountPrt(b,vc,llik);
output off;
```

A single vector of effect parameters and one scalar dispersion parameter are estimated. The vector of effect parameters includes one element corresponding to each explanatory variable and a constant term. In the example, five parameters are estimated.

Constrained Truncated-at-zero Negative Binomial Variance Function Regression Model

```
library cmlcount;
#include cmlcount.ext;

CMLCountset;
dataset = "wars";
dep1 = { wars };
ind1 = { unem, poverty, allianc };
ind2 = { partyid, x4 };
_cmlc_ZeroTruncate = 0;
{ b,vc,llik } = CMLNegbin(dataset,dep,ind1,0);
output file = cmlcount.out reset;
call CMLCountPrt(b,vc,llik);
output off;
```

Two vectors of effect parameters are estimated, one for the mean and one for the variance function. Each vector includes a constant term and one element corresponding to each explanatory variable. In the example, the variables specified in ind1 pertain to the expected value and ind2 to the variance. Seven parameters are estimated.

Source

cmlnegbn.src

Purpose

Estimates a constrained Pareto regression model, for the analysis of duration data, with maximum likelihood.

Library

cmlcount

Format

```
\{b, vc, llik\} = CMLPareto(dataset, dep, ind);
```

Input

dataset string, name of GAUSS data set.

– or –

N×K matrix, data

dep string, the name of the dependent variable

– or –

scalar, the index of the dependent variable

ind K×1 character vector, names of the independent variables

– or -

K×1 numeric vector, indices of independent variables

Set to 0 to include only a constant term.

If dataset is a matrix, dep and ind may be a string or character variable containing either the standard labels created by **CML** (V1, V2,..., or V01, V02,...., depending on the value of **__vpad**), or the user-provided labels in **__altnam**.

Output

b (K+1)×1 vector, maximum likelihood estimates of the effect parameters

stacked on top of the dispersion parameter.

vc (K+1)×(K+1) matrix, variance-covariance matrix of the estimated

parameters evaluated at the maximum. If the **CML** global **_cml_CovPar** is set to 3, vc will contain heteroskedastic-consistent parameter estimates.

llik scalar, value of the log-likelihood function at the maximum.

Globals

CML globals are also relevant, including constraint matrices and procedures.

CMLCount Reference

__cmlc__Inference string, determines the type of statistical inference.

BOOT generates bootstrapped estimates and covariance matrix of estimates

CML generates maximum likelihood estimates

Setting **__cmlc__Inference** to BOOT generates a **GAUSS** data set containing the bootstrapped parameters. The file name of this data set is either the default BOOTx, where x is a four digit number starting with 1 and increasing until a unique name is found, or the name in the **CML** global variable, **__cml__BootFname**. This data set can be used with **CMLBlimits** for generating confidence limits, with **CMLDensity** for generating density estimates and plots of the boostrapped parameters, or with **CMLHist** for generating histogram and surface plots.

_cmlc_Censor string, the name of the censor variable from *dataset*

scalar, the index of the censor variable from dataset

Each element of censor variable is 0 if censored, or 1 if not.

By default, no censoring is used.

__cmlc__Start scalar, selects method of calculating starting values. Possible values are:

- ocalculates them by regressing ln(y + 0.5) on the explanatory variables.
- will use a vector of user supplied start values stored in the global variable _cmlc_StartValue.
- **2** uses a vector of zeros.
- 3 uses random uniform numbers on the interval $\left[-\frac{1}{2}, \frac{1}{2}\right]$.

Default = 0.

- **_cmlc_StartValue** $(K+1)\times 1$ vector, start values if **_cmlc_Start** = 1.
- **_cmlc_Dispersion** scalar, start value for scalar dispersion parameter. Default = 3.
- **_cmlc_Precision** scalar, number of decimal points to print on output. Default = 4.
- **__altnam** K×1 vector, alternate names for variables when a matrix is passed to **CMLPareto**. When a data matrix is passed to **CMLPareto** and the user is selecting from that matrix, the global variable **__altnam**, if it is used, must contain names for the columns of the original matrix.
- **__miss** scalar, determines how missing data will be handled.
 - **0** Missing values will not be checked for, and so the data set must not have any missings. This is the fastest option.

5. CMLCOUNT REFERENCE

1 Listwise deletion. Removes from computation any observation with a missing value on any variable included in the analysis.

Default = 0.

__output

scalar, determines printing of intermediate results.

- **0** nothing is written.
- 1 serial ASCII output format suitable for disk files or printers.
- 2 (DOS only) output is suitable for screen only. ANSI.SYS must be active.

Default = 2.

__row

scalar, specifies how many rows of the data set will be read per iteration of the read loop. By default, the number of rows to be read will be calculated automatically.

__rowfac

scalar, row factor. If **CMLPareto** fails due to insufficient memory while attempting to read a **GAUSS** data set, then **__rowfac** may be set to some value between 0 and 1 to read a *proportion* of the original number of rows of the **GAUSS** data set. For example, setting

$$_{rowfac} = 0.8;$$

will cause **GAUSS** to read in 80% of the rows originally calculated.

This global only has an affect when $__$ row = 0.

Default = 1.

__title

string, message printed at the top of the screen and printed out by **CMLCountPrt**. Default = "".

__vpad

scalar, if *dataset* is a matrix in memory, the variable names are automatically created by **CML**. Two types of names can be created:

- Variable names automatically created by **CML** are not padded to give them equal length. For example, V1, V2,...V10, V11,....
- Variable names created by the procedure are padded with zeros to give them an equal number of characters. For example, V01, V02, ..., V10, V11,.... This is useful if you want the variable names to sort properly.

Default = 1.

CMLCount Reference

Remarks

Let the *n* duration observations (non-negative real numbers) for the dependent variable be denoted as y_1, \ldots, y_n . Assume that y_i follows a Pareto distribution with expected value μ_i and variance $\mu_i \sigma^2 + \mu_i^2$. Let the mean μ_i be an exponential-linear function of a vector of explanatory variables, x_i :

$$E(y_i) \equiv \mu_i = \exp(x_i \beta) \tag{5.7}$$

The program includes a constant term as the first column of x_i and allows one to include any number of explanatory variables. Note that μ_i from a duration model equals $1/\lambda_i$ from an event count model; thus, one need only change the sign of the effect parameters to get estimates of the same parameters from these different kinds of data.

The dispersion σ^2 is parametrized as follows:

$$\sigma_i^2 = \exp(\gamma) \tag{5.8}$$

The program gives estimates of β and γ .

For an introduction to the Pareto regression model see Hannan and Tuma (1984) and King, Alt, Burns, and Laver (1989).

Example

Pareto Regression Model of Duration Data

A vector of effect parameters and a scalar dispersion parameter are estimated. The vector includes one element corresponding to each explanatory variable named in *ind* and a constant term. Five parameters are estimated in this example.

Constrained Censored Pareto Regression Model of Duration Data

CMLPareto

A vector of effect parameters and a scalar dispersion parameter are estimated. The vector includes one element corresponding to each explanatory variable named in ind and a constant term. Five parameters are estimated in this example.

Source

cmlparet.src

Purpose

Estimates a constrained Poisson regression model or truncated-at-zero Poisson regression model with maximum likelihood.

Library

cmlcount

Format

```
\{b,vc,llik\} = CMLPoisson(dataset,dep,ind);
```

Input

dataset string, name of **GAUSS** data set.

– or –

N×K matrix, data

dep string, the name of the dependent variable

– or –

scalar, the index of the dependent variable

ind K×1 character vector, names of the independent variables

– or –

K×1 numeric vector, indices of independent variables

Set to 0 to include only a constant term.

If dataset is a matrix, dep and ind may be a string or character variable containing either the standard labels created by **CML** (V1, V2,..., or V01, V02,..., depending on the value of **__vpad**), or the user-provided labels in **__altnam**.

Output

b K×1 vector, maximum likelihood estimates of the effect parameters.

vc K×K matrix, variance-covariance matrix of the estimated parameters

evaluated at the maximum. If you choose the **CML** global $_$ cml $_$ CovPar = 3, vc will contain heteroskedastic-consistent parameter estimates.

llik scalar, value of the log-likelihood function at the maximum.

Globals

CML globals are also relevant, including constraint matrices and procedures.

__cmlc__Inference string, determines the type of statistical inference.

BOOT generates bootstrapped estimates and covariance matrix of estimates

CML generates maximum likelihood estimates

Setting **__cmlc__Inference** to BOOT generates a **GAUSS** data set containing the bootstrapped parameters. The file name of this data set is either the default BOOTx, where x is a four digit number starting with 1 and increasing until a unique name is found, or the name in the **CML** global variable, **__cml__BootFname**. This data set can be used with **CMLBlimits** for generating confidence limits, with **CMLDensity** for generating density estimates and plots of the boostrapped parameters, or with **CMLHist** for generating histogram and surface plots.

__cmlc_Fix scalar, name of index number of logged variable among the regressors with coefficient constrained to 1.0 By default, no logged variables are included.

_cmlc_Start scalar, selects method of calculating starting values. Possible values are:

- ocalculates them by regressing ln(y + 0.5) on the explanatory variables.
- will use a vector of user supplied start values stored in the global variable _cmlc_StartValue.
- **2** uses a vector of zeros.
- 3 uses random uniform numbers on the interval $\left[-\frac{1}{2}, \frac{1}{2}\right]$.

Default = 0.

_cmlc_StartValue $K \times 1$ vector, start values if **_cmlc_Start** = 1.

_cmlc_Precision scalar, number of decimal points to print on output. Default = 4.

_cmlc_ZeroTruncate scalar, specifies which model is used:

- **0** truncated-at-zero negative binomial model
- 1 negative binomial model is used.

Default = 1.

__altnam K×1 vector, alternate names for variables when a matrix is passed to **CMLPoisson**. When a data matrix is passed to **CMLPoisson** and the user is selecting from that matrix, the global variable **__altnam**, if it is used, must contain names for the columns of the original matrix.

__miss scalar, determines how missing data will be handled.

- **0** Missing values will not be checked for, and so the data set must not have any missings. This is the fastest option.
- 1 Listwise deletion. Removes from computation any observation with a missing value on any variable included in the analysis.

Default = 0.

__output scalar, determines printing of intermediate results.

- **0** nothing is written.
- 1 serial ASCII output format suitable for disk files or printers.
- 2 (DOS only) output is suitable for screen only. ANSI.SYS must be active.

Default = 2.

___row scalar, specifies how many rows of the data set will be read per iteration of the read loop. By default, the number of rows to be read will be calculated automatically.

___rowfac scalar, row factor. If **POISSON** fails due to insufficient memory while attempting to read a **GAUSS** data set, then **___rowfac** may be set to some value between 0 and 1 to read a *proportion* of the original number of rows of the **GAUSS** data set. For example, setting

 $_{rowfac} = 0.8;$

will cause GAUSS to read in 80% of the rows originally calculated.

___title string, message printed at the top of the screen and printed out by $\mathbf{CMLCountPrt}$. Default = "".

__vpad scalar, if *dataset* is a matrix in memory, the variable names are automatically created by **CML**. Two types of names can be created:

- Variable names automatically created by **CML** are not padded to give them equal length. For example, V1, V2,...V10, V11,....
- Variable names created by the procedure are padded with zeros to give them an equal number of characters. For example, V01, V02, ..., V10, V11,.... This is useful if you want the variable names to sort properly.

Default = 1.

Remarks

Let the n event count observations (non-negative integers) for the dependent variable be denoted as y_1, \ldots, y_n . y_i is then a random dependent variable representing the number of events that have occurred during observation period i. By assuming that the events occurring within each period are independent and have constant rates of occurrence, y_i can be shown to follow a Poisson distribution:

$$f_p(y_i|\lambda_i) = \begin{cases} \frac{e^{-\lambda_i}(\lambda_i)^{y_i}}{y_i!} & \text{for } \lambda_i > 0 \text{ and } y_i = 0, 1, \dots \\ 0 & \text{otherwise} \end{cases}$$
 (5.9)

with expected value and variance λ_i . Under the Poisson regression model, λ_i (the rate of event occurrence, which must be greater than zero) is assumed to be an exponential-linear function of a vector of explanatory variables, x_i :

$$E(y_i) \equiv \lambda_i = \exp(x_i \beta) \tag{5.10}$$

The program includes a constant term as the first element of x_i and allows one to include any number of explanatory variables.

For an introduction to the Poisson regression model see King (1988); on the truncated model, see Grogger and Carson (1988) and King (1989d).

Example

Constrained Poisson Regression Model

Constrained Truncated-at-zero Poisson Regression Model

```
library cmlcount;
#include cmlcount.ext;
```

CMLPoisson

$5. \ \mathit{CMLCOUNT} \ \mathit{REFERENCE}$

```
CMLCountset;
dataset = "wars";
dep = { wars };
ind = { unem, poverty, allianc };
_cmlc_ZeroTruncate = 0;
{ b,vc,llik } = CMLPoisson(dataset,dep,ind);
output file = cmlcount.out reset;
call CMLCountPrt(b,vc,llik);
output off;
```

Source

cmlpoiss.src

Purpose

Estimates a constrained seemingly unrelated Poisson regression model, for the analysis of two event $CONSTRAINED\ COUNT$ variables, with maximum likelihood.

Library

cmlcount

Format

 $\{b,vc,llik\} = CMLSupreme(dataset,dep1,dep2,ind1,ind2);$

Input

datasetstring, name of GAUSS data set. – or – N×K matrix, data dep1string, name of the first dependent variable – or – scalar, index of the first dependent variable string, name of the second dependent variable dep2– or – scalar, index of the second dependent variable K×1 character vector, names of first event independent variables ind1K×1 numeric vector, indices of first event independent variables Set to 0 to include only a constant term. ind2K×1 character vector, names of second event independent variables – or – K×1 numeric vector, indices of second event independent variables Set to 0 to include only a constant term.

If dataset is a matrix, dep1, dep2, ind1 and ind2 may be a string or character variable containing either the standard labels created by **CML** (V1, V2,..., or V01, V02,...., depending on the value of **__vpad**), or the user-provided labels in **__altnam**.

Output

b (K+L+2)×1 vector, maximum likelihood estimates of the effect parameters of β and γ stacked on top of the covariance parameter ξ .

CMLCount Reference

vc $(K+L+2)\times(K+L+2)$ matrix, variance-covariance matrix of the estimated parameters evaluated at the maximum. If you choose the **CML** global $_cml_CovPar = 3$, vc will contain heteroskedastic-consistent parameter estimates.

llik scalar, value of the log-likelihood function at the maximum.

Globals

CML globals are also relevant, including constraint matrices and procedures.

__cmlc__Inference string, determines the type of statistical inference.

BOOT generates bootstrapped estimates and covariance matrix of estimates

CML generates maximum likelihood estimates

Setting **_cmlc_Inference** to BOOT generates a **GAUSS** data set containing the bootstrapped parameters. The file name of this data set is either the default BOOTx, where x is a four digit number starting with 1 and increasing until a unique name is found, or the name in the **CML** global variable, **_cml_BootFname**. This data set can be used with **CMLBlimits** for generating confidence limits, with **CMLDensity** for generating density estimates and plots of the boostrapped parameters, or with **CMLHist** for generating histogram and surface plots.

_cmlc_Start scalar, selects method of calculating starting values. Possible values are:

- ocalculates them by regressing ln(y + 0.5) on the explanatory variables.
- will use a vector of user supplied start values stored in the global variable _cmlc_StartValue.
- **2** uses a vector of zeros.
- 3 uses random uniform numbers on the interval $\left[-\frac{1}{2}, \frac{1}{2}\right]$.

Default = 0.

_cmlc_StartValue $(K+L+2)\times 1$ vector, start values if **_cmlc_Start** = 1.

_cmlc_Precision scalar, number of decimal points to print on output. Default = 4.

__altnam K×1 vector, alternate names for variables when a matrix is passed to **CMLSupreme**. When a data matrix is passed to **CMLSupreme** and the user is selecting from that matrix, the global variable **__altnam**, if it is used, must contain names for the columns of the original matrix.

__miss scalar, determines how missing data will be handled.

- **0** Missing values will not be checked for, and so the data set must not have any missings. This is the fastest option.
- 1 Listwise deletion. Removes from computation any observation with a missing value on any variable included in the analysis.

Default = 0.

__output scalar, determines printing of intermediate results.

- **0** nothing is written.
- 1 serial ASCII output format suitable for disk files or printers.
- 2 (DOS only) output is suitable for screen only. ANSI.SYS must be active.

Default = 2.

__row scalar, specifies how many rows of the data set will be read per iteration of the read loop. By default, the number of rows to be read will be calculated automatically.

__rowfac scalar, row factor. If **CMLSupreme** fails due to insufficient memory while attempting to read a **GAUSS** data set, then **__rowfac** may be set to some value between 0 and 1 to read a *proportion* of the original number of rows of the **GAUSS** data set. For example, setting

 $_{rowfac} = 0.8;$

will cause GAUSS to read in 80% of the rows originally calculated.

This global only has an affect when $__row = 0$.

Default = 1.

___title string, message printed at the top of the screen and printed out by **CMLCountPrt**. Default = "".

__vpad scalar, if *dataset* is a matrix in memory, the variable names are automatically created by **CML**. Two types of names can be created:

- Variable names automatically created by **CML** are not padded to give them equal length. For example, V1, V2,...V10, V11,....
- Variable names created by the procedure are padded with zeros to give them an equal number of characters. For example, V01, V02, ..., V10, V11,.... This is useful if you want the variable names to sort properly.

Default = 1.

Remarks

Suppose we observe two event count dependent variables y_{1i} and y_{2i} for n observations. Let these variables be distributed as a bivariate Poisson with $E(y_{1i}) = \lambda_{1i}$ and $E(y_{2i}) = \lambda_{2i}$. These means are parametrized as follows:

$$\lambda_{0i} = \exp(x_i \beta) \tag{5.11}$$

and

$$\lambda_{+i} = \exp(z_i \gamma) \tag{5.12}$$

where x_i and z_i are (possibly) different vectors of explanatory variables. The covariance parameter is ξ .

If you have convergence problems, you might try **CMLSupreme2** with argument ind3 = 0 instead.

For details about this model, see King (1989c).

Example

Constrained Seemingly Unrelated Poisson Regression Model (CMLSupreme)

Two vectors of effect parameters and one scalar covariance parameter are estimated. The vectors of effect parameters each include one element corresponding to each explanatory variable and a constant term. In the example, ten parameters are estimated.

Source

cmlsupr.src

Purpose

Estimates a constrained Poisson regression model with unobserved dependent variables, for the analysis of two observed (and three unobserved) event count variables, with maximum likelihood.

Library

cmlcount

Format

 $\{b,vc,llik\} = CMLSupreme2(dataset,dep1,dep2,ind1,ind2,ind3);$

Input

datasetstring, name of **GAUSS** data set. – or – N×K matrix, data string, name of the first dependent variable dep1scalar, index of the first dependent variable string, name of the second dependent variable dep2scalar, index of the second dependent variable ind1K×1 character vector, names of first event independent variables $K\times 1$ numeric vector, indices of first event independent variables Set to 0 to include only a constant term. ind2L×1 character vector, names of second event independent variables L×1 numeric vector, indices of second event independent variables Set to 0 to include only a constant term.

If dataset is a matrix, dep1, dep2, ind1, ind2, or ind3 may be a string or character variable containing either the standard labels created by **CML** (V1, V2,..., or V01, V02,..., depending on the value of **__vpad**), or the user-provided labels in **__altnam**.

Set to 0 to include only a constant term.

M×1 character vector, names of second event independent variables

M×1 numeric vector, indices of second event independent variables

Output

ind3

- b (K+L+M)×1 vector, maximum likelihood estimates of the effect parameters of β and γ stacked on top of the covariance parameter ξ .
- vc (K+L+M)×(K+L+M) matrix, variance-covariance matrix of the estimated parameters evaluated at the maximum. If you choose the **CML** global **_cml_CovPar** = 3, vc will contain heteroskedastic-consistent parameter estimates.
- llik scalar, value of the log-likelihood function at the maximum.

Globals

CML globals are also relevant, including constraint matrices and procedures.

_cmlc_Inference string, determines the type of statistical inference.

BOOT generates bootstrapped estimates and covariance matrix of estimates

CML generates maximum likelihood estimates

Setting **_cmlc_Inference** to BOOT generates a **GAUSS** data set containing the bootstrapped parameters. The file name of this data set is either the default BOOTx, where x is a four digit number starting with 1 and increasing until a unique name is found, or the name in the **CML** global variable, **_cml_BootFname**. This data set can be used with **CMLBlimits** for generating confidence limits, with **CMLDensity** for generating density estimates and plots of the boostrapped parameters, or with **CMLHist** for generating histogram and surface plots.

_cmlc_Start scalar, selects method of calculating starting values. Possible values are:

- o calculates them by regressing ln(y + 0.5) on the explanatory variables.
- will use a vector of user supplied start values stored in the global variable **_cmlc_StartValue**.
- 2 uses a vector of zeros.
- 3 uses random uniform numbers on the interval $\left[-\frac{1}{2}, \frac{1}{2}\right]$.

Default = 0.

__altnam

_cmlc_StartValue $(K+L+M)\times 1$ vector, start values if **_cmlc_Start** = 1.

_cmlc_Precision scalar, number of decimal points to print on output. Default = 4.

K×1 vector, alternate names for variables when a matrix is passed to **CMLSupreme2**. When a data matrix is passed to **CMLSupreme2** and the user is selecting from that matrix, the global variable **__altnam**, if it is used, must contain names for the columns of the original matrix.

__miss scalar, determines how missing data will be handled.

- **0** Missing values will not be checked for, and so the data set must not have any missings. This is the fastest option.
- 1 Listwise deletion. Removes from computation any observation with a missing value on any variable included in the analysis.

Default = 0.

__output scalar, determines printing of intermediate results.

- **0** nothing is written.
- 1 serial ASCII output format suitable for disk files or printers.
- 2 (DOS only) output is suitable for screen only. ANSI.SYS must be active.

Default = 2.

__row scalar, specifies how many rows of the data set will be read per iteration of the read loop. By default, the number of rows to be read will be calculated automatically.

___rowfac scalar, row factor. If **CMLSupreme2** fails due to insufficient memory while attempting to read a **GAUSS** data set, then **___rowfac** may be set to some value between 0 and 1 to read a *proportion* of the original number of rows of the **GAUSS** data set. For example, setting

$$_{rowfac} = 0.8;$$

will cause **GAUSS** to read in 80% of the rows originally calculated.

This global only has an affect when $__row = 0$.

Default = 1.

___title string, message printed at the top of the screen and printed out by $\mathbf{CMLCountPrt}$. Default = "".

___vpad scalar, if *dataset* is a matrix in memory, the variable names are automatically created by **CML**. Two types of names can be created:

- Variable names automatically created by **CML** are not padded to give them equal length. For example, V1, V2,...V10, V11,....
- Variable names created by the procedure are padded with zeros to give them an equal number of characters. For example, V01, V02, ..., V10, V11,.... This is useful if you want the variable names to sort properly.

Default = 1.

Remarks

This model assumes the existence of three independent unobserved variables, y_{1i}^* , y_{2i}^* , and y_{3i}^* , with means $E(y_{ji}^*) = \lambda_{ji}$, for j = 1, 2, 3. Although these are not observed, we do observe y_{1i} and y_{2i} , which are functions of these three variables:

```
y_{1i} = y_{1i}^* + y_{3i}^* 
 y_{2i} = y_{2i}^* + y_{3i}^*
```

The procedure estimates three separate regression functions, one for the expected value of each of the unobserved variables:

$$\lambda_{1i} = \exp(x_{1i}\beta_1)$$

$$\lambda_{2i} = \exp(x_{2i}\beta_2)$$

$$\lambda_{3i} = \exp(x_{3i}\beta_3)$$

$$(5.13)$$

where x_{1i} , x_{2i} and x_{3i} are (possibly) different sets of explanatory variables and β_1 , β_2 , and β_3 are separate parameter vectors. This option produces maximum likelihood estimates for these three parameter vectors.

Example

Poisson Regression Model with Unobserved Dependent Variables

Three vectors of effect parameters are estimated. Each includes one element corresponding to each explanatory variable plus a constant term. In the example, twelve parameters are estimated.

Source

cmlsupr2.src

 ${\bf CMLSupreme2}$

5. CMLCOUNT REFERENCE

Index

active parameters, 12	_cml_CutPoint, 61
algorithm, 31	_cml_D , 15, 37, 46
Alt-1 , 31	_cml_Delta , 37, 38
Alt-2 , 31	_cml_Diagnostic , 13, 38, 39, 44
Alt-3 , 31	_cml_DirTol , 31, 37, 39
Alt-4 , 31	_cml_EqJacobian , 22, 37, 39, 49, 53
Alt-A , 31	_cml_EqProc , 15, 37, 40, 53
Alt-H , 30	_cml_Extrap, 37, 40
altnam , 75, 77, 90, 95, 99, 104, 109,	_cml_FinalHess, 37, 40
114, 119, 123	_cml_GradCheckTol, 20, 37, 40
В	$_$ cml $_$ GradMethod $, 31, 37, 41$
D	_cml_GradProc , 37, 41, 48
BFGS, 10, 31, 38	$_$ cml $_$ GradStep, $37, 41$
BHHH, 31, 38	$_$ cml $_$ HessCov $,37,41$
BHHHSTEP, 12	_cml_HessProc , 20, 37, 41, 49
bootstrap, 24, 28, 75, 79	_cml_Increment, 61, 64
bounds, 16, 47, 78	_cml_IneqJacobian , 22, 37, 42, 49, 54
BRENT, 11	$_cml_EqProc, 15$
C	$_$ cml $_$ IneqProc $,\ 37,\ 42,\ 54$
C	_cml_Interp , 37, 42
ccount.src, 87, 88	$_$ cml $_$ lterData $,38,42$
CHGVAR, 6	_cml_Kernel , 28, 59
CML, 36	_cml_Lag , 38, 42
cml.src, 50, 56, 66, 67	$_$ cml $_$ Lagrange $, 26, 37, 43$
_cml_A , 14, 37, 38, 46, 52	$_$ cml $_$ LineSearch $,\ 37,\ 43$
_cml_Active, 12, 38	$_$ cml $_$ MaxIters, $37, 43, 54$
_cml_Algorithm, 37, 38	_cml_MaxTime , 37, 44, 58
_cml_Alpha, 51, 56	_cml_Maxtry, 37
_cml_B , 14, 37, 38, 46, 52	$_$ cml $_$ MaxTry $,31,44$
_cml_BootFname, 58	_cml_NumCat , 61, 64
_cml_Bounds, 16, 37, 47, 78	_cml_NumObs , 28, 38, 56, 58
_cml_C , 15, 37, 38, 46, 53	_cml_NumPoints, 59
_cml_Center, 61, 64	_cml_NumSample , 28, 44, 58, 64
_cml_CovPar, 25, 26, 36, 37, 39, 75	_cml_Options, 37, 44

cml ParNamas 29 44 54	CMI Poisson 119
_cml_ParNames, 38, 44, 54	CMLPoisson, 113
_cml_RandRadius, 11, 37, 43	cmlprof.src, 65 CMLProfile, 24, 63
_cml_Select, 51, 56, 64	
_cml_Smoothing, 28, 59	CMLPrt, 67 CMLCLPrt, 68
_cml_Truncate, 59	,
_cml_UserHess, 22	CMLSet, 66
_cml_UserNumGrad, 37, 44	cmlsupr.src, 121
_cml_UserNumHess, 37, 45	cmlsupr2.src, 125
_cml_UserSearch, 37, 45	CMLSupreme 78 122
_cml_Width , 61, 64	CMLTlimits 56
$_$ cml $_$ XprodCov $, 37, 45$	CMLTlimits, 56
${\tt cmlblim.src}, 51$	condition of Hessian, 13
CMLBlimits, 51, 80	confidence limits, 27
CMLBoot , 24, 28, 57	constraint Jacobians, 22
cmlboot.src, 58	constraints, 14, 22, 46, 78
_cmlc_Boot , 75	convergence, 43
_cmlc_Censor, 75	converting MAXLIK programs, 6
_cmlc_Dispersion, 75	covariance matrix, parameters, 23, 25, 26, 29, 39
_cmlc_Fix, 75	20, 29, 39 cubic step, 43
_cmlc_Inference, 75, 78	cubic step, 49
_cmlc_Precision, 75	D
_cmlc_Start, 75	
_cmlc_StartValue , 77, 90, 95, 99, 103,	derivatives, 9, 18, 39, 48
109,114,119,123	DFP, 10, 31, 38
_cmlc_ZeroTruncate, 75	diagnosis, 13
${\tt cmlclim.src}, 55$	E
CMLClimits , 23, 27, 52, 79	D
CMLCountCLPrt, 87	equality constraints, 14, 15, 38, 40, 46,
CMLCountPrt, 86	52, 53
CMLCountSet, 88	·
${\tt cmldens.src}, 60$	G
CMLDensity , 28, 59, 75, 80	
CMLExpgam, 89	global variables, 30
${\tt cmlexpgm.src},93$	gradient, 36, 57, 63, 67, 68
CMLExpon, 94	gradient procedure, 18, 41, 48
cmlexpon.src, 97	TT
CMLHist , 24, 28, 61, 75, 80	H
cmlhist.src, 62	IIAID 11
cmlhurdl.src, 101	HALF, 11
CMLHurdlep, 78, 98	Hessian, 10, 13, 25, 31
CMLNegbin, 102	Hessian procedure, 20, 22, 49
cmlnegbn.src, 107	heteroskedastic-consistent covariance
cmlparet.src, 112	matrix, 26, 39
CMLPareto, 108	I
cmlpoiss.src, 117	

INDEX

inactive parameters, 12	regression, Hurdle Poisson, 72
inequality constraints, 15, 38, 39, 42,	regression, negative binomial, 72
46, 53, 54 Installation, 1	regression, seemingly unrelated Poisson, 72
J	regression, truncated negative binomial, 72
Jacobian, 22	regression, truncated Poisson, 72 resampling, 28
L	row, 8, 36, 38, 42, 45, 47, 48, 57, 63
	rowfac, 38, 45, 77, 91, 96, 100, 104,
Lagrange coefficients, 26, 43 likelihood profile trace, 29, 30	110, 115, 120, 124 run-time switches, 30
line search, 10, 31	S
linear constraints, 14, 15, 38, 39, 46, 52,	S
53 log-likelihood function, 7, 8, 22, 36, 47,	scaling, 13
48, 57, 63	Shift-1, 31
log-linear, 72	Shift-2, 31 Shift-4, 31
M	Shift-3, 31
	Shift-5 , 31
maximum likelihood, 7, 36, 57, 71	starting point, 13
MAXLIK programs, converting, 6	statistical inference, 23, 78
miss , 90, 95, 99, 104, 109, 114, 119,	step length, 10, 31, 43
124	STEPBT, 11
N	T
NEWTON, 10, 31, 38, 49	title , 38, 45
nonlinear constraints, 15, 40, 42, 46, 53,	
54	U
NR, 31	UNIX, 1, 3
O	V
	V
output, 31, 44, 59, 61, 75, 77, 91, 96, 100, 104, 110, 115, 120, 124	VPUT , 14
	VREAD, 14
P	W
profile t plot, 29	
Ω	weight, 12, 38, 45 weighted maximum likelihood, 12
	Windows/NT/2000, 3
quadratic step, 43	
quasi-Newton, 10	
R	