

SUPPLEMENTARY MATERIALS FOR:

**'Jointness in Bayesian Variable Selection  
with Applications to Growth Regression'**

Eduardo Ley

*The World Bank, Washington DC, U.S.A.*

Mark F.J. Steel

*Department of Statistics, University of Warwick, U.K.*

**Version:** February 22, 2008

**Abstract.** This document contains some supplementary graphs, and computational notes to accompany our paper E. Ley and M.F.J. Steel (2007), "Jointness in Bayesian Variable Selection with Applications to Growth Regression," *Journal of Macroeconomics*, 29(3): 476–493.

**Address.** Eduardo Ley, The World Bank, Washington DC 20433. Email: eduley@gmail.com.  
Mark F.J. Steel, Department of Statistics, University of Warwick, Coventry CV4 7AL, U.K.  
Vox: +44(0)24 7652 3369. Email: M.F.Steel@stats.warwick.ac.uk

## 1. Overview

This document contains some supplementary materials for the paper E. Ley and M.F.J. Steel (2007), “Jointness in Bayesian Variable Selection with Applications to Growth Regression,” *Journal of Macroeconomics* (LS henceforth).

Files in \*.zip archive:

<b>File</b>	<b>Directory</b>	<b>Contents</b>	<b>Format</b>	<b>See</b>
Oreadme.pdf	1s5/	This file	PDF	-
1s5bma.f	1s5/code/	f77 source file	Text	Section on Fortran code
1s5bma.par	1s5/code/	Parameter file	Text	Section on Parameter file
grk41t72.dat	1s5/code/	FLS data file	Text	Section on Data
grk67t88.dat	1s5/code/	SDM data file	Text	Section on Data
k41i9_Std_.out	1s5/out/k41/std/	FLS standardized out file	Text	Section on Output Files
k41i9_Std_bj.dat	1s5/out/k41/std/	Bi-variate Jointness Measures	Text	Section on Output Files
k41i9_Std_tj.dat	1s5/out/k41/std/	Tri-variate Jointness Measures	Text	Section on Output Files
k41i9_NStd.out	1s5/out/k41/nstd/	FLS non-standardized out file	Text	Section on Output Files
k67i9_Std_.out	1s5/out/k67/std/	SDM standardized out file	Text	Section on Output Files
k67i9_Std_bj.dat	1s5/out/k67/std/	Bi-variate Jointness Measures	Text	Section on Output Files
k67i9_Std_tj.dat	1s5/out/k67/std/	Tri-variate Jointness Measures	Text	Section on Output Files
k67i9_NStd.out	1s5/out/k67/nstd/	SDM non-standardized out file	Text	Section on Output Files

## 2. Supplementary Figures and Tables for LS

Fig. 1 below displays the behavior of the measures of bivariate jointness for the SDM data, Fig. 2 from LS is also reproduced here for ease of comparison.

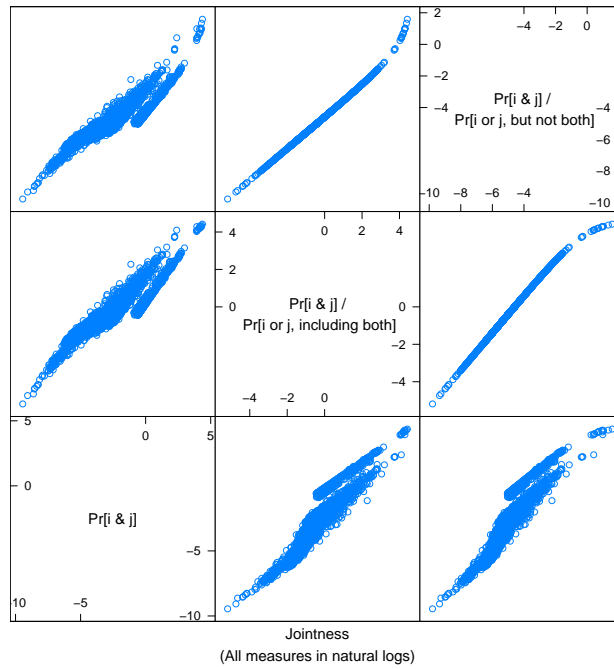


Fig 1. SDM data: Joint inclusion probabilities, bivariate jointness  $\mathcal{J}_{ij}^*$  and  $\mathcal{J}_{ij}$ .

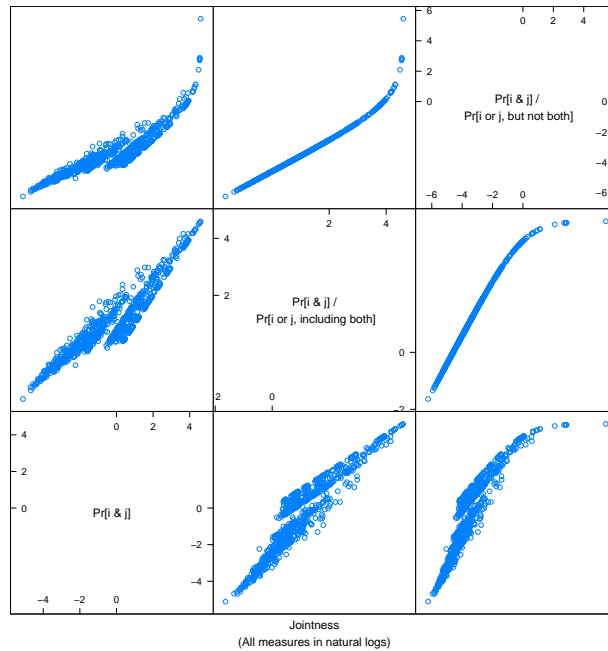


Fig 2. FLS data: Joint inclusion probabilities, bivariate jointness  $\mathcal{J}_{ij}^*$  and  $\mathcal{J}_{ij}$ .

The enclosed k67i9\_Std\_.out output file (described below) includes the complete Table 6 in LS for all 67 regressors. The \*NStd.out files also include the results for non-standardized regressors for both the SDM and FLS datasets.

### 3. Code and Data Files

Overview: In order to reproduce the results in the paper, you'll need to compile the f77 file `ls5bma.f` and generate an executable, say, `bma.exe`. You'll need the data (`*.dat`) file and the parameter (`ls5bma.par`) file in the same directory—this file controls some options explained below.

The successful execution will always produce an output file `*.out`, and, if the jointness option (`dojoint`) is set to true, also two additional files containing the jointness measures for all pairs and triplets (see below).

#### 3.1. Fortran code

The file `ls5bma.f` contains the **f77 source code** to produce the results in the paper. It is standard f77 code except for the time-date functions. (If you cannot link to the standard Unix (`-unixlib`) or VAX/VMS (`-vmslib`) libraries, either remove the calls to `wr_date()` and `wr_time()`, or empty the body of these subroutines so that they do nothing when they are called.)

There are several **parameters** that you may need to modify before compiling the code if you are going to make different runs than the ones included here:

`maxm` Specifies the maximum number of models expected to be visited, currently set to 300,000. If you run longer chains, or use different datasets, and find the chain visiting more than `maxm` models, you'll have to set this parameter to a larger value. The number of 'visits out' in the `*.out` file will be an upper bound on the increment needed to `maxm`, since some of these may be to repeated models.

`maxk` Maximum number of regressors in your dataset.

`maxn` Maximum number of observations in your dataset.

`maxnf` Maximum number of observations to be used for out-of-sample prediction, if `lpsloop` set to true in `ls5bma.par` file.

The program has been tested on different machines and with different compilers:

- **Absoft's Pro Fortran for MacOS X v9.0** on a G4 iMac. It is important to use static storage (`-s`). Other suggested options include: optimization (`-O3 -cpu:host`), fold to upper case (`-N109`), maximum internal handle (`-T 100000`), temporary string size (`-t102400`).
- **Intel Fortran Compiler 9.1 for MacOS** on an Intel Core 2 duo iMac, under MacOS X. You simply need to disable the compiler's inline expansion (`-Ob0`):

```
> ifort ls5bma.f -Ob0 -o bma.exe
```

- **Intel Fortran Compiler 10 for MacOS** on an Intel Core 2 duo iMac, under MacOS X; no directives required:

```
> ifort ls5bma.f -o bma.exe
```

A typical run on a 2.16GHz Intel Core 2 duo iMac takes less than 30 minutes.

- **Compaq Visual Fortran 6.6** on a dual Xeon workstation running Windows XP. Full optimization settings can be used.

### 3.2. Parameter File

The `ls5bma.par` file controls some of the execution time parameters.

The first line sets the name of the different output files (the bivariate and trivariate `*.dat` files will have ‘bj’ and ‘tj’ appended to this name).

The second line must contain the *exact* name of the data file.

The rest of the parameters are fairly self-explanatory. (Also, check the routine setup in `ls5bma.f` to see what the parameters do.) For a description of prediction issues and the G&M convergence estimate refer to FLS, and random  $\theta$  priors are discussed in LSb.

An example of a `ls5bma.par` file follows:

```
1 useanyname      OUT namefile, char10
2 grk41t72.dat    DAT namefile
3 -81665432       negative int to init random number generator
4 9               int 1-9, choice for g, see computefj and [FLS]
5 100000          warmup draws
6 500000          chain draws
7 T              standard --- standardise Xs?
8 F              LPSloop --- do prediction?
9 10             integer specifying nf
10 0.85d0         real taking care of sample split
11 T             wrpost --- write posterior results
12 T             dogm --- do G&M 2nd chain to assess convergence?
13 T             dojoint --- produce jointness results?
14 T             fixed theta --- F means random theta, see [LS]
15 20.5          prior expected model size
```

### 3.3. Data File

We enclose the two data files used in LS; FLS has  $k = 41, t = 72$ , and SDM has  $k = 67, t = 88$ . As noted, the exact name of the data file to use at execution must be specified in the second line in `ls5bma.par`.

The data file must have the following structure.

Line 1: The first line must have the number of observations (`ntot`).

Line 2: The second line the number of regressors (`kreg`).

Lines 3 to (`kreg + 2`): Then there must follow `kreg` lines with the corresponding variable names.

Lines (`kreg + 3`) to ( $2 \times \text{ntot} + 2$ ): Next, another `ntot` lines, each with the values of the dependent variable,  $y$ , and the `kreg` regressors,  $X$ , in free format—*i.e.*, each with a total of `kreg + 1` columns.

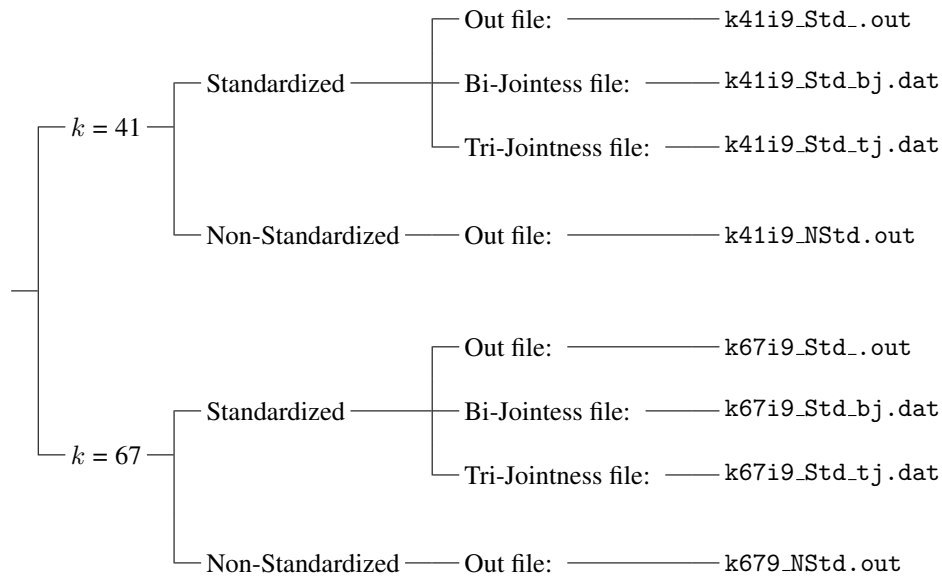
We show the listing of the first 48 lines of `grk41t72.dat` below. (Note that, to save space, lines from 44 onwards, each containing a  $(kreg + 1)$  array of data for each observation, are shown here truncated, without wrapup.)

```
1 72
2 41
3 GDPsh560
4 Confuncious
```



### 3.4. Output Files

For each run, we produce three output files. The \*.out file contains the output of the chain, while the two \*.dat files contain the different bi- or tri-variate jointness measures. (These two files are only generated when dojoint in the ls5bma.par file is set to true (T).)



### 4. Updating of the code to handle $k > 52$

A model is represented by an array of 0s and 1s:  $m = (m_1, m_2, \dots, m_k)$  with  $m_i \in \{0, 1\}$ . An index  $idx \in \{1, 2, \dots, 2^k\}$  tracks the  $2^k$  possible models.

When  $idx$  is a real\*8 variable, it can count only to  $2^{52}$ . This means that the loop below will print 53 — same result will be obtained in GAUSS and Matlab.

```

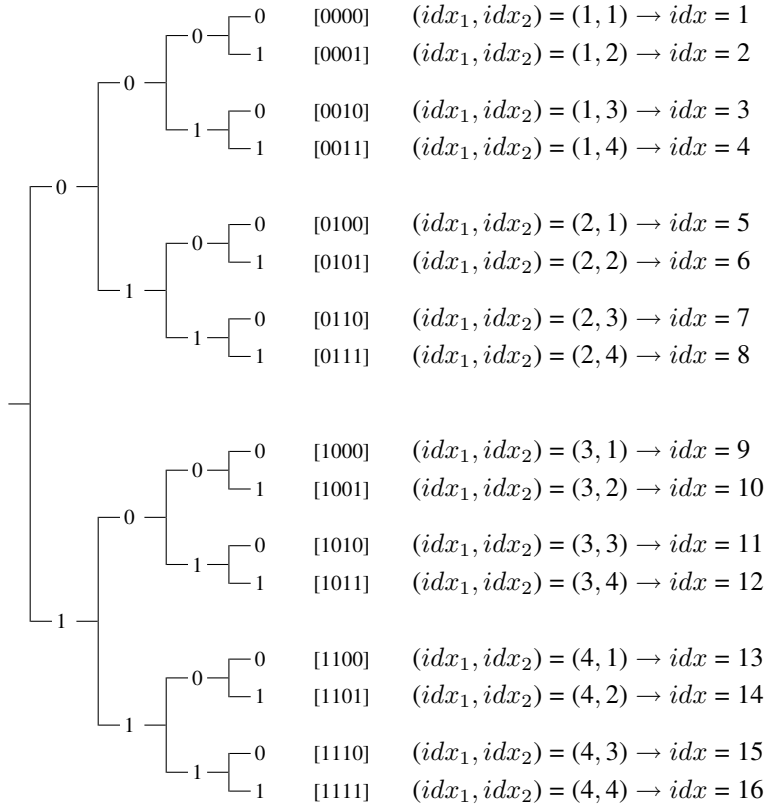
initiliaze (d = 1, i = 1)
while (d > 0) do
  i ← (i + 1)
  x1 ← 2i
  x2 ← (1 + 2i)
  d ← (x2 - x1)
enddo
print i

```

Thus, we cannot distinguish between  $2^{52}$  and  $1 + 2^{52}$

To get around this problem, we can use 2 indices to keep track of models —  $idx_1$  keeps track of the first 52 vars and  $idx_2$  keeps track of the remaining  $k - 52$  vars. Now we can handle up to 104 vars, or  $2^{104}$  models.

**Example:** 4 variables;  $2^4 = 16$  models. If we could only count to 4, we could just use 2 indexes,  $idx_1$  and  $idx_2$ , each in  $\{1, 2, 3, 4\}$  to account for the 16 models indexed by  $idx$ .



*Counting*—Notice that the first  $2^{k-1} = 2^3 = 8$  integers are associated with vectors that have with  $m_1 = 0$ . Start out with  $i = 1$ , then a 1 in  $m_1$  adds  $2^{k-1} = 8$  to  $idx$ , a 1 in  $m_2$  adds  $2^{k-2} = 4$  to  $idx$ , ..., a 1 in  $m_k$  simply adds  $2^0 = 1$  to  $idx$ .

The function `getmodidx(k,m)` starts with  $idx = 1$  and makes  $idx = idx + zm_i$  for  $z = 2^{k-1}$  down to 1. Now it will have to be replaced by a subroutine with 2 calls to the old FLS function, each call taking care of a portion of the binary array  $m$ .

```

1 c
2 c*****
3     real*8 function getmodidx(k,m)
4 c*****
5 c
6 c inputs:

```



```

7 c    k.le.52.....is the number of all possible regressors
8 c    m(k)...contains 0s or 1s in each coordinate
9 c output:
10 c   getmodidx..is the real*8 associated with the state vector
11 c
12     real*8 idx,z
13     integer  k,i,m(k)
14
15     z = 2.0d0**(k-1)
16     idx = 1
17     do 10 i=1,k
18         idx = idx + z*m(i)
19         z = z/2.0d0
20 10   continue
21     getmodidx = idx
22     return
23
24     end
25 c
26 c
27 c*****
28     subroutine get2modidx(k,m,idx1,idx2)
29 c*****
30 c
31 c inputs:
32 c    k>52.....is the number of all possible regressors
33 c    m(1:k)...contains 0s or 1s in each coordinate
34 c output:
35 c    idx1..is the real*8 associated with the 1st 52-var state vector
36 c    idx2..is the real*8 associated with the 2nd 52-var state vector
37 c
38     integer  i,k, m(k), m1(52), m2(52)
39     real*8 idx1,idx2,getmodidx
40
41     k1=min(52,k)
42     do 10 i=1,k1
43         m1(i)=m(i)
44 10   continue
45     idx1 = getmodidx(k1,m1)
46
47     k2 = k-52
48     if (k2.ge.1) then
49         do 15 i=1,k2
50             m2(i)=m(i+52)

```

```

51 15      continue
52          idx2 = getmodidx(k2,m2)
53      else
54          idx2=0.0d0
55      endif
56
57      return
58      end
59 c
60 c*****

```

The subroutine `gmodel(idx,k,m)` does the reverse—starting with `idx` constructs the binary array `m` representing the model. (If  $idx > 2^{k-j}$  then  $m_j = 1$ , otherwise  $m_j = 0$ , for  $j = 1, k - 1$ .) Now it will be replaced by another routine specifying 2 indexes: `idx1` and `idx2`.

```

1 c*****
2      subroutine gmodel(idx,k,m)
3 c*****
4 c
5 c given the model index idx and the number of all possible regressors
6 c k<53
7 c it returns a binary array, m, of dimension k. if the ith coordinate is 1
8 c the ith regressor is included in the model. otherwise it is excluded.
9 c
10 c inputs:
11 c     idx.....is the model index (real*8)
12 c     k.....is the number of all possible regressors k=1..52
13 c
14 c output:
15 c     m(k)...contains 0s or 1s in each coordinate
16 c
17      integer i, k, m(k)
18      real*8 idx,x,z
19
20      x = 2.0d0**(k-1)
21      z = idx
22
23      do 20 i=1,k
24          if (z.gt.x) then
25              m(i) = 1
26              z = z - x
27          else
28              m(i) = 0
29          endif
30      x = x / 2.0d0

```

```

31 20  continue
32      return
33      end
34 c
35 c*****
36      subroutine g2model(idx1,idx2,k,m)
37 c*****
38 c
39 c  given the model index idx and the number of all possible regressors k
40 c  it returns a binary array, m, of dimension k.  if the ith coordinate is 1
41 c  the ith regressor is included in the model.  otherwise it is excluded.
42 c
43 c inputs:
44 c   idx....is the model index (real*8)
45 c   k>52....is the number of all possible regressors
46 c
47 c output:
48 c   m(k)...contains 0s or 1s in each coordinate
49 c
50     integer  i,k,m1(52),m2(52),m(k)
51     real*8  idx1,idx2
52
53     k1 = min(k,52)
54     call gmodel(idx1,k1,m1)
55     do 10 i=1,k1
56         m(i)=m1(i)
57 10  continue
58
59     k2 = k-52
60     if (k2.ge.1) then
61         call gmodel(idx2,k2,m2)
62         do 15 i=1,k2
63             m(52+i)=m2(i)
64 15  continue
65     endif
66
67     return
68     end
69 c
70 c*****

```

## 5. References

- [FLS] Fernández, Carmen, Eduardo Ley and Mark F.J. Steel (2001) “Model Uncertainty in Cross-Country Growth Regressions,” *Journal of Applied Econometrics*, 16: 563–76.
- [LS5] Ley, Eduardo and Mark F.J. Steel (2007) “Jointness in Bayesian Variable Selection with Applications to Growth Regression,” *Journal of Macroeconomics* 29(3): 476–493.
- [LS6] Ley, Eduardo and Mark F.J. Steel (2008) “On the Effect of Prior Assumptions in Bayesian Model Averaging with Applications to Growth Regression,” *Journal of Applied Econometrics*, forthcoming.
- [SDM] Sala-i-Martin, Xavier X., Gernot Doppelhofer and Ronald I. Miller (2004) “Determinants of Long-Term Growth: A Bayesian averaging of classical estimates (BACE) approach.” *American Economic Review*, 94: 813–835.