A Computer Program for Tests of Equality of Covariance Matrices

by

Regina Becker and K. C. S. Pillai Purdue University

Technical Report #85-8

Department of Statistics Purdue University

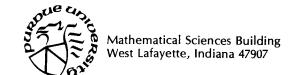
April 1985



Enclosed is a technical report containing documentation and a listing of a computer program for performing the computations described in the paper, "A Computer Program for Tests of Equality of Covariance Matrices," by Regina Becker and K. C. S. Pillai. Dr. Pillai passed away on June 5, 1985. He was a distinguished scholar and a leader in the field of multivariate analysis.

The research and programming for this project was supported by the Purdue University Department of Statistics. Accordingly, this work is freely distributed to anyone who requests it.

If, however, you have funds available to support this type of research, we would be very grateful to receive a charitable contribution of \$50.00 or more. The enclosed form can be used for this purpose. The money so collected will be used to start an annual memorial lecture on <u>Multivariate Analysis and Its Applications</u> here at Purdue.



Contribution to the Purdue Foundation*

Date _		 	·		
Amount		 	<u> </u>		
From	- n ·	 		 	

*Please make check payable to Purdue Foundation and put the $\underline{\text{Pillai Memorial Lecture Fund}}$ in the corner.

Mail to: Department of Statistics Purdue University West Lafayette, Indiana 47907

I. Introduction

This document is a supplement to the paper, "Percentage points of the largest characteristic root of the multivariate beta matrix", by K. C. S. Pillai and Bernhard N. Flury in Communications in Statistics, Theory and Methods (1984) 13(18), 2199-2237. A program is presented here which is an algorithm for the test of equality of the covariance matrices of two p-variate normal populations based on Roy's largest root test statistic, illustrated by a numerical example. The program runs on the Purdue University Computing Center's CDC 6000 system. Some minor modifications may be needed for other computers.

Seven subroutines from the IMSL package are needed to run the program, namely, 1) becovm, 2) linv2p, 3) ludecp, 4) vmulff, 5) vmulfs, 6) eigrs, and 7) algama. Usage of the subroutines is described in the sequel.

Users who do not have access to the IMSL package could substitute their own subroutines. It is possible that IMSL could allow a user to obtain copies of these seven routines only for use in this program. For information, contact IMSL, NBC Building, 7500 Bellaire Blvd., Houston, Texas 77036; (713) 772-1927.

The test mentioned above may be described in more detail. Let S_h (pxp) be the sum product (SP) matrix computed from a random sample of size n_h from $N_p(\mu_h, \Sigma_h)$, i.e. p-variate normal population with mean vector μ_h and covariance matrix Σ_h which is symmetric positive definite, h = 1,2. The two random samples are independently selected and $p < n_1, n_2$. Then a.e. the p characteristic roots of $S_{l \times 2}^{S_0-1}$ are all positive. Further, let $\gamma_l, \ldots, \gamma_p$ be the characteristic roots of $\Sigma_{l \times 2}^{S_0-1}$. Now the test of the hypothesis H_0 : $\gamma_i = l$, i = l,...,p, versus H_a : $\gamma_i \ge 1$, i = 1,...,p, $\sum_{i=1}^p \gamma_i > p$ may be carried out at level α as follows:

Reject H_O if the largest characteristic root of $S_1(S_1+S_2)^{-1}=\theta_{p,m,n}>\theta_{1-\alpha,p,m,n}$, otherwise do not reject. Here $m=\frac{1}{2}(n_1-p-2)$, $n=\frac{1}{2}(n_2-p-2)$ and the upper $(l-\alpha)$ percentile may be found from the reference above (Pillai and Flury or references therein).

If one prefers to compute the P-value the program gives the method to do so, however the P-value may be used with confidence only if it does not exceed ten percent in view of the fact that it is computed from an approximation to the cdf of the largest root at the upper end (see Pillai and Flury (1984)).

While Section II describes the input, Section III gives the input for an example and Section IV the output for the example. The data consists of measurements on four variables (1) height (inches), (2) weight (pounds), (3) chest (inches) and (4) waist (inches) of male reserve officers in civilian status of the Armed Forces of the Philippines, hailing from two regions of the Philippine Islands but all within the age interval 29 to 31. The numbers of officers selected at random were 20 and 24 respectively from the first and second regions. (Normality assumption was found to be justified in view of earlier tests. Also for the original data see S. R. Ventura (1957). On the extreme roots of a matrix in multivariate analysis and associated tests. Unpublished thesis. The Statistical Center, University of the Philippines, Manila.)

Section V gives the program listing, including the usages of the seven subroutines listed above.

II. Input

Program: eqcov

Order of control cards:

- a. jobcard
- b. pfiles(get,datafile,x=tape2)
- c. pfiles(get,eqcov,id=dnt)
- d. mnf(i=eqcov,l=trash)
- e. #eor
- f. parameter card
- g. format card
- h. (input)
- j. #s

Preparation of cards specific to this program:

Format card

Columns 1-80 may be used to describe the data format. Use the usual fortran format statement omitting the word format.

Data or matrix input

- a. If data is input, enter each case on a new line. The format statement should describe one case. Enter n1+n2 cases, each case having nv variables. Enter cases for data set 1 followed by cases for data set 2.
- b. If matrices are input, enter the lower triangular portion or the complete matrix for data set 1 followed by the lower triangular portion or the complete matrix for data set 2.

III. Input for Example

4 20	24	5	2	
(4f10.0)				
61.9880				
157.8800	4043.80	000		
4.6200	295.70	000	73.5500	
1.0900	479.15	500	50.2250	81.6375
72.0163				
195.5875	3045.95	583		
31.3625	281.20	83	92.4583	
2.6062	388.68	3 7 5	36.8125	93.6563

IV. Output for Example

test for equality of covariance matrices based on the largest root test

number of variables is 4 sample size for data set 1 is 20 sample size for data set 2 is 24

data set 1 cross product matrix:

61.9880 157.8800 4043.8000 4.6200 295.7000 73.5500 1.0900 479.1500 50.2250 81.6375

data set 2 cross product matrix:

72.0163 195.5875 3045.9583 31.3625 281.2083 92.4583 2.6062 388.6875 36.8125 93.6563

eigenvalues:

.3736 .6583 1.2219 1.5195

the largest root is .6031

P-value is .6662

***note: the P-value is approximate.
if the P-value is .10 or below it can be used with confidence.
otherwise, the value of the largest root may be compared to the tables.

program main (input, output, tape5=input, tape6=output, tape2)

```
c*
C*
        This program is an algorithm for a test for equality
        of covariance matrices of two normal populations.
C*
С*
         It is based on the procedure described in:
C*
C*
        Pillai, K.C.S. and Flury, Bernhard N. (1984),
С*
        "Percentage Points of the Largest Characteristic
С*
        Root of the Multivariate Beta Matrix,"
C*
        Communications in Statistics- Theory and Methods,
c*
        13(18), 2199-2237.
С*
c*
                   Department of Statistics
c*
                     Purdue University
c*
                   West Lafayette, IN 47907
c*
c*
     programmer - Regina Becker
C*
     Department of Statistics
c*
     Purdue University
c*
     January, 1985
C*
С
С
      ***input information***
С
С
С
   variables:
              nv = number of variables (must be less
С
                   than or equal to 20)
С
С
              n1 = sample size of data set 1 (must be less than
С
                   or equal to 300)
С
С
              n2 = sample size of data set 2 (must be less than
С
                   or equal to 300)
С
C
              itape = the informational input source
С
                     5 if information is on input
С
                     2 if information is on tape
С
С
              ntype = the informational input code
C
                     0 if data is input
C
                     1 if covariance matrices input
C
                     2 if cross product matrices input
С
С
   order of input:
С
С
                        nv,n1,n2,itape,ntype
С
    1. parameter card
С
                        use (515) format
С
С
    2. format card
                        Columns 1-80 may be used to describe the
С
                        informational input format. Use the usual
С
                        fortran type format statement, omitting
```

the word format.

3. data or matrix input

C

C C

С

C C

С

С

C C

С

С

c c

С

С

С

C C

C C

C C

С

С

CCC

С

C

С

C C

C C

С

c c

C C

С

С

C C

С

С

С

C C

С

C C

C

- a. If data is input, enter each case on a new line. The format statement should describe one case. Enter n1+n2 cases, each case having nv variables. Enter cases for data set 1 followed by cases for data set 2.
- b. If matrices are input, enter the lower triangular portion or the complete matrix for each data set. For example, if there are 3 variables the covariance matrix input should be:

s11 s21 s22 s31 s32 s33 s41 s42 s43 s44 r11 r21 r22 r31 r32 r33 r41 r42 r43 r44

note: the s matrix is covariance matrix of data set 1 the r matrix is covariance matrix of data set 2

The following subroutines are called by the program:

- 1. becovm from the IMSL library, this subroutine calculates the cross product matrix of each data set if data is input.
- 2. linv2p from the IMSL library, this subroutine finds the inverse of a matrix
- 3. ludecp from the IMSL library, this subroutine performs a triangular decomposition of a matrix x such that x = 1 * 1-transpose. l is a lower triangular matrix
- 4. vmulff from the IMSL library, this subroutine performs multiplication of two matrices stored in full storage mode
- 5. vmulfs from the IMSL library, this subroutine performs multiplication of two matrices a and b where a is stored in full storage mode and b is stored in symmetric storage mode
- 6. eigrs from the IMSL library, computes the eigenvalues of a real symmetric matrix
- 7. algama from the IMSL library, this function computes the natural log of the gamma function

subroutine becovm

C C C

```
С
       usage: call becovm(x,ix,nbr,temp,xm,sp1,ier)
С
                   input matrix of dimension nbr(3) by nv for which cross
С
              x
С
                   product matrix is desired. nbr(3)=sample size
C
С
                   input, row dimension of x as specified in dimension
              ix
С
                   statement in calling program. ix = n1size
С
                   input vector of length 6.
С
              nbr
                   nbr(1) = nv
С
С
                   nbr(2) = number of observations per variable
C
                   nbr(3) = nbr(2)
С
                   nbr(4) = 1
                   nbr(5) = 1
С
С
                   nbr(6) = 1
С
С
              temp input vector of length nv
С
С
                   output vector of length nv containing variable means
              \mathbf{x}\mathbf{m}
С
С
              sp1
                   output matrix of dimension nv by nv stored in symmetric
                   storage mode requiring nv * (nv+1)/2 locations. spl
С
С
                   contains the cross products matrix
С
С
              ier
                   error parameter
С
С
       subroutine linv2p
С
С
С
       usage: call linv2p(sp2,nv,sp2inv,idgt,d1,d2,wkarea,ier)
С
C
                   nv by nv positive definite symmetric matrix to be
С
                   inverted.
                            sp2 is stored in symmetric storage mode
С
С
              nv
                   order of sp2
С
С
                      output vector of length nv(nv+1)/2 containing
С
                      the inverse of sp2. storage is symmetric mode.
С
С
              idgt the approximate number of digits in the answer which
С
                   were unchanged after improvement
С
С
              d1,
                   components of the determinant of sp2
              d2
С
                   determinant(sp2)=d1*2.**d2 (output)
С
С
              wkarea work area of dimension nsize (250)
С
C
              ier error parameter
С
C
       subroutine ludecp
С
С
С
       usage: call ludecp(sp2inv,ul,nv,d1,d2,ier)
С
С
                     input vector of length nv(nv+1)/2 containing the
С
                     nv by nv positive definite symmetric matrix stored
                     in symmetric storage mode
С
С
```

```
С
              ul
                   output vector of length nv(nv+1)/2 containing the
С
                   decomposed matrix 1 such that sp2inv=1*1-transpose.
С
                   1 is stored in symmetric storage mode. The diagonal
С
                   of 1 contains the reciprocals of the actual diagonal
С
                   elements.
С
С
              nv
                  order of sp2inv
С
                   components of the determinant of sp2inv.
С
              d1.
С
                   determinant(sp2inv) = d1*2.**d2 (output)
              d2
С
C
              ier
                  error parameter
С
С
C
       subroutine vmulff
С
С
       usage: call vmulff(c,l,r,c1,c2,ia,ib,lt,ic,ier)
С
С
                  r by cl matrix stored in full storage mode
              С
C
С
              1
                  cl by c2 matrix stored in full storage mode
С
С
              r
                  number of rows in c = nv
С
С
              c1
                  number of columns in c = nv
С
С
              c2
                  number of columns in l = nv
С
С
                  row dimension of matrix c as specified in the dimension
              ia
С
                  statement in the calling program = maxnv
С
                  row dimension of matrix l as specified in the dimension
С
              ib
С
                  statement in the calling program = maxnv
С
С
              lt
                  r by c2 matrix containing the product lt=c*l
С
С
                  row dimension of matrix lt as specified in the dimension
С
                  statement in the calling program = maxnv
С
С
              ier error parameter
С
С
С
       subroutine vmulfs
С
С
       usage: call vmulfs(lt,sp1,l,m,ia,c,ic)
С
                  1 by m matrix stored in full storage mode
С
              lt
С
С
                  m by m symmetric matrix stored in symmetric storage mode
С
                  number of rows in lt = nv
С
              1
С
C
             m
                  number of columns in lt = nv
С
С
              ia
                  row dimension of matrix lt as specified in dimension
С
                  statement in the calling program = maxnv
С
                  1 by m matrix containing the product c=lt*sp1
С
              С
```

```
С
С
                 row dimension of c as specified in the dimension
С
                 statement in the calling program = maxnv
С
С
       subroutine eigrs
С
С
С
       usage: call eigrs(ul,nv,jobn,d,c,ic,wkarea,ier)
С
             ul
                 input real symmetric matrix of order nv whose
С
                 eigenvalues are to be computed
С
С
                 order of the matrix ul
С
             nv
С
             jobn = 0 is an order to compute eigenvalues only
С
С
С
             d
                 output vector of length nv containing eigenvalues of ul
С
C
                 output nv by nv matrix of eigenvectors
С
С
             ic
                 input row dimension of c as specified in dimension
С
                 statement in calling program = nv
С
C
             wkarea work area of length at least nv
С
С
             ier error parameter
С
С
       subroutine algama
С
С
C
       usage: call algama(x)
С
С
             x
                    input argument
С
С
             algama
                    output value of the log base e of the absolute
                    value of gamma(x)
C
С
С
                    a parameter statement has been used to set maximum
С
      a special note:
                    sizes for these variables:
С
С
С
                    n1
                         maximum is nlsize = 300
С
С
                         maximum is n2size = 300
                    n2
С
С
                    nv
                         maximum is nvsize = 20
С
                     ind maximum is nvsize*(nvsize+1)/2 = 210
C
C
C
                    wkarea maximum is nsize. nsize must be at least
С
                          nv(nv+1)/2 + n
С
С
С
      These sizes may be increased if storage space on your computer allows
С
      for larger array storage.
                             In that sense, variable limits can be
С
      controlled by the programmer.
С
```

```
parameter (n1size=300, n2size=300, nvsize=20, ind=210, nsize=250)
      real pval, x (nlsize, nvsize), temp (nvsize), sp1 (ind), sp2 (ind),
      *wkarea(nsize), sp2inv(ind), lt(nvsize, nvsize), l(nvsize, nvsize),
     *t(nvsize),ul(ind),c(nvsize,nvsize),h(nvsize),m,n
      real xm (nvsize), d (nvsize)
      integer df1, df2, fmt (80), nbr (6)
      double precision sum, f (nvsize), ff
        read in nv-number of variables, n1-sample size of data set1
С
        n2=sample size of data set2, itape (5 if data is input,
С
С
        2 if data is on pfiles), ntype(0 if data, 1 if covariance matrix.
C
        2 if cross product matrix)
        read(5,10) nv,n1,n2,itape,ntype
10
        format (3i5, 2i5)
        maxnv=20
        write (6,410)
        format(lx,'***test for equality of covariance matrices based on',
410
       ' the largest root test***',//)
        write(6,412) nv,n1,n2
412
        format(1x, 'number of variables is ',i2/x, 'sample size for data'
     * 'set 1 is ',i3/x, 'sample size for data set 2 is ',i3/)
        check that nv le maxnv
С
        if (nv.qt.maxnv) then
        write(6,8)
        format(1x,'***number of variables must be less than 21'/)
8
        goto 900
        endif
        if (ntype.lt.0.or.ntype.gt.2) then
        write(6,12)
        format(1x,'***check input type specification'/
12
        1x, 'ntype=0 if data , 1 if covariance matrix, 2 if cross product
        matrix')
        goto 900
        endif
        if (itape.ne.5.and.itape.ne.2) then
        write(6,9)
        format(1x,'***incorrect input medium specified'/
9
        1x, 'use itape=5 if data is from input'/
        1x, 'use itape=2 if data is from outside file'/)
        goto 900
        endif
        df1=n1-1
        df2=n2-1
        min=df1
        if (df2.lt.df1) min=df2
        if (nv.gt.dfl.or.nv.gt.df2) then
        write(6,11) min
        format(1x, '***number of variables must be less than the'/
11
        1x, 'smaller degrees of freedom which is ', i4//)
        goto 900
        endif
```

```
pie=4.*atan(1.0)
         read data format and data
С
         read(5,20) (fmt(i), i=1,80)
20
         format (80a1)
         if (ntype.eq.0) then
         do 30 i=1,n1
         read(itape, fmt) (x(i,j), j=1,nv)
30
         continue
C
         call becovm to get variance-covariance matrix
        nbr(1) = nv
        nbr(2)=n1
        nbr(3)=n1
        nbr(4) = 1
        nbr(5) = 1
        nbr(6) = 1
        ix=n1size
        call becovm(x,ix,nbr,temp,xm,sp1,ier)
        spl has cross product matrtix since nbr(6)=1
С
C
        this is stored in symmetric storage mode
        do 40 i=1,n2
        read(itape, fmt) (x(i,j), j=1,nv)
40
        continue
        nbr(2)=n2
        nbr(3)=n2
        ix=n2size
        call becovm(x,ix,nbr,temp,xm,sp2,ier)
        else
        do 43 i=1,nv
        ii=(i-1)*i/2
        read(itape, fmt) (spl(ii+j), j=1,i)
43
        continue
        index=nv*(nv+1)/2
        do 45 i=1,nv
        ii = (i-1) * i/2
        read(itape, fmt) (sp2(ii+j), j=1,i)
45
        continue
        if (ntype.eq.2) go to 49
        do 47 i=1, index
        spl(i)=spl(i) * dfl
        sp2(i)=sp2(i) * df2
 47
        continue
49
        endif
        write(6,400)
400
        format(1x, 'data set 1 cross product matrix: '/)
        do 401 i=1,nv
        ii=(i-1)*i/2
        write (6,402) (spl(ii+j), j=1,i)
401
        continue
```

```
write (6,403)
403
        format(/x, 'data set 2 cross product matrix: '/)
        do 404 i=1,nv
        ii=(i-1)*i/2
        write (6,402) (sp2(ii+j),j=1,i)
402
        format(1x, 12f10.4)
404
        continue
С
        find sp2 inverse
        call linv2p(sp2,nv,sp2inv,idgt,d1,d2,wkarea,ier)
        sp2 inverse = sp2inv. decompose this to 1*1 transpose
С
        call ludecp(sp2inv,ul,nv,d1,d2,ier)
        ul has matrix 1 such that 1*1-transpose = sp2inv
С
С
        diagonal of 1 has reciprocals of actual diagonal elements
        ul(1) = 1./ul(1)
        jj=1
        do 50 i=2,nv
        jj=jj+i
        ul (jj) = 1./ul (jj)
50
        continue
        write 1 in full storage mode
С
        do 60 i=1,nv
        inc=(i-1) * i / 2
do 60 j=1,i
        l(i,j)=ul(inc+j)
        if(i.ne.j) 1(j,i)=0.0
60
        continue
        multiply 1-transpose * sp1 * 1
С
        do 70 i=1,nv
        do 70 j=1,nv
        lt(j,i)=l(i,j)
70
        continue
        call vmulfs(lt,sp1,nv,nv,maxnv,c,maxnv)
        call vmulff(c,1,nv,nv,nv,maxnv,maxnv,lt,maxnv,ier)
        this results in 1-transpose*sp1*1 into lt. need characteristic
С
С
        roots of lt. convert this to ssm, put in ul (not used now).
        call vcvtfs(lt,nv,maxnv,ul)
        get eigenvalues of ul
С
        call eigrs(ul,nv,0,d,c,nv,wkarea,ier)
        write(6,406)
406
        format(//x,'eigenvalues:'/)
```

```
write (6,666) (d(i),i=1,nv)
666
         format (1x, 12f10.4//)
С
         d contains eigenvalues, smallest to largest
С
         arrange theta(i) (called t(i)) largest to smallest
         do 80 i=1,nv
         t(nv-i+1)=d(i)/(1.+d(i))
80
         continue
        m=float(df1-nv-1)/2.
        n=float(df2-nv-1)/2.
        rnv=float(nv)
        nvt=nv/2
        if (nv.eq.2*nvt) then
        hh=1.0
        else
        xx=t(1)
        call mdbeta(xx,m+1.,n+1.,hh,ier)
        endif
С
        compute h(1) to h(nv)
        h(1)=1.0
        do 100 i=2,nv-1
        h(i) = float(nv-i+1) * (2.*m + float(nv-i+2)) * h(i-1) /
              (float(i-1)*(2.*m + 2.*n + float(2*nv-i+2)))
100
        continue
C
        compute f
        f(1) = (m+n) / (m+n+rnv)
        do 120 i=2,nv-1
        f(i) = ((m+n) *h(i) - (m+float(nv-i+1)) *f(i-1)) /
              (m+n+float(nv-i+1))
120
        continue
        sum=0.0
        do 130 i=1,nv-1
        sum=sum+(-1.0)**i*f(i)*t(1)**(nv-i)
130
        continue
        z1=.5*(2.*m+2.*n+2.*rnv+1.)
        z2=.5*(2.*m+rnv+1.)
        z3=.5*(2.*n+rnv+1.)
        z4=.5*rnv
        z5=.5*(2.*m+2.*n+rnv+2.)
        cst=(.5*alog(pie)+algama(m+n+rnv+1.)+algama(z1))-
             (alog(m+n) + algama(z2) + algama(z3) + algama(z4) + algama(z5))
        find antilog of cst
С
        cst=exp(cst)
        ff=hh + cst * (t(1) **m) * ((1.-t(1)) ** (n+1.)) *sum
        write(6,408) t(1)
408
        format(//x, 'the largest root is ', f10.4)
        p=1.-ff
        write(6,407) p
407
        format(//x, 'p-value is ', f6.4///)
```

```
write(6,409)
format(1x,'***note: the p-value is approximate.')
write(6,414)
format(4x,'if the p-value is .10 or below it can be used',

* 'with confidence.')
write(6,415)
format(4x,'otherwise, the value of the largest root may be',

* 'compared to the tables.'//)

900
stop
end
#eor
```