

PREDICTION OF MISSING DATA USING REGRESSION MODELS: A PROGRAMMED APPROACH FOR LARGE SPSS SYSTEM FILES

RON BOBNER
YOUNGSTOWN STATE UNIVERSITY
HAL STALCUP
U.S. AIR FORCE, WRIGHT-PATTERSON AFB
ISADORE NEWMAN AND CAROLYN BENZ
THE UNIVERSITY OF AKRON

Abstract

The purpose of this article is to present computer programming capable of automatically predicting missing data in an SPSS system file using multiple regression techniques. Two versions are presented. The first is an interactive approach designed to run via VSPC under OS/VS MVS. The second is designed to be used in a card batch environment. In addition, an equation which predicts the amount of CPU time required is included.

Incomplete data are often a major concern in behavioral research. SPSS, SASS, and BIOMED, three of the most commonly available statistical packages, all have provisions for handling missing data. These subroutines generally take the form of some type of case deletion. This approach, however, decreases the N size and therefore the power of the analysis. In addition, if regression methods are used, a decrease in the N size will increase the upward bias of the R (Newman, 1972) and decrease the stability of the weights (Newman, 1973).

Two other alternatives for handling missing data are sometimes used. One method is the substitution of the mean

value calculated from the complete cases. However, this technique tends to decrease the variance and may therefore be inappropriate for analysis of variance procedures. The other approach is to predict the missing values through the use of a regression equation whose weights are calculated from the complete cases. This latter approach tends to bias any subsequent analysis to a lesser extent than the insertion of means.

While neither SPSS, SASS, or BIOMED makes direct provision for the insertion of means for missing values, it is possible by utilizing various calculation and data management subroutines within the packages. By comparison, the calculation and insertion of missing data based on the regression approach, utilizing subroutines existing within the packages, is so complicated that this procedure is impractical even when dealing with small amounts of missing data.

The following programming was designed to provide a solution to this problem. Two versions of the program are presented. The first was designed to run interactively via VSPC under OS/VS MVS where VS FORTRAN, the RUML subroutine of the IMSL (International Mathematical and Statistical Libraries, Inc., 1980), and SPSS Release 8.1 (McGraw-Hill, 1979) routines have been implemented. The second version was designed for use with a card deck in a batch mode and only requires the implementation of SPSS, the IMSL, and FORTRAN. The final section of this paper presents a method for predicting the amount of CPU time that will be required to predict and insert the missing data.

Interactive Programming

The interactive programming consists of two routines designed to run consecutively: INTRO and INTRO1. They will call on SPSS system file, sort the complete from the incomplete cases, and further sort the incomplete cases into two categories: $\leq 10\%$ and $> 10\%$ missing data. Next, the program will scan the first case in the $\leq 10\%$ missing data file and locate the first missing variable, build a regression equation predicting this variable utilizing only those predictors available from that particular case, and calculate the weights utilizing the complete case file. These weights, along with the values for the predictors from the first case, will be used to predict the missing value for the variable. An a priori decision was made to round the predicted value up or down as appropriate if the first three digits dropped are $> \pm .455$ of the last significant digit kept. If the first three digits dropped are in the range of $< \pm .455$ of the last significant digit kept, the value of the last significant digit is randomly rounded up or down. This value is then inserted into the data matrix and the program then continues its scan of the first case in the $\leq 10\%$ missing data file for other missing variables. If it should find another missing value the process is repeated. In no instance, however, are previously predicted variables entered into the prediction process of subsequent variables for that particular case.

Once the scan and prediction of missing variables for the first case is completed, the program repeats the process for the second and each subsequent case in the $\leq 10\%$ missing data file. After the last case has been completed, the program combines the complete case file with newly completed version of the $\leq 10\%$ missing data file. This new file is in BCD format and can be used as a raw data file for input into SPSS. Its name is USER FINISH CASES.

This new system file will contain complete data on all cases which initially were complete or had $\leq 10\%$ missing data. Those cases which had $> 10\%$ missing data will have been discarded. The choice of the 10% cutoff is based on the fact that $> 10\%$ missing data for any one case will not unduly bias the prediction of the missing data (Cohen & Cohen, 1975).

The following sections document the interactive versions of INTRO and INTROL.

INTRO

10 DIMENSION H(72),NAME(80),IFORM(89),NAMSYS(8),IP(80)
20 + NAMSUR(8)
30 INTEGER YES,GO
40 DATA YES,NO,GO//‘Y’,‘N’,‘GO’ //
50 ITEST=9
60 CALL OFSYS(‘COMMAND’,‘ALL’,IRTN,IMESS,‘FILE JOROUT’)
70 CALL OFSYS(‘ALLOC’,‘JOBOUT’,9)
80 CALL CLEAR
90 WRITE(6,290)
100 200 FORMAT(6X,’ENTRY INTO THE DATA MANAGEMENT SOFTWARE.’/
110 +1X,’PURPOSE : DATA EDITING AND CASE FILE ADDITIONS.’/
120 +1X,’PROCEDURE:
130 +1X, 1. PLACE SPSS FILE ON DISK FILE. //
140 +1X, 2. SCAN CASES FOR COMPLETE CASES. //
150 +1X, 3. SPLIT CASES AND DELETE >10% MISSING VALUES. //
160 +1X, 4. USE COMPLETE CASES TO REPLACE MISSING VALUES. //
170 +1X, 5. REWRITE SPSS FILE WITH CORRECT CASES. ///)
180 CALL CLEAR
190 WRITE(6,201)
200 201 FORMAT(1X,’WELL, WELL RON; HOW THE HELL ARE YOU’)
210 READ(5,100) IANS
220 100 FORMAT(1A4)
230 CALL REPLY(IANS)
240 5 WRITE(6,202)
250 202 FORMAT(1X,’ARE YOU READY TO DUMP THE SPSS DATA FILE???’)
260 CALL ANSWER(IRN)
270 IF(IRN.EQ.0) GO TO 10
280 IF(ITEST.EQ.99) GO TO 99
290 WRITE(6,203)

LIT

300 203 FORMAT(6X, 'THEN WHY DID YOU START THIS PROGRAM ?????'/
310 +1X, 'I AM GOING TO ASSUME THAT YOU CAN NOT TYPE AND YOU NEED'/
320 +1X, 'ANOTHER CHANCE TO ANSWER THE QUESTION.'//
330 GO TO 5
340 10 WRITE(6, 204)
350 204 FORMAT(1X, 'VERY WELL THEN I AM READY TO TRANSFER THE DATA.')
360 WRITE(6, 205)
370 205 FORMAT(1X, 'I NEED SOME INFORMATION;')
380 WRITE(6, 206)
390 READ(5, *) NVARS
400 206 FORMAT(1X, 'FIRST, THE NUMBER OF VARIABLES IN THE FILE')
410 WRITE(6, 207)
420 READ(5, 101) (NAME(II), II=1, 55)
430 101 FORMAT(80A1)
440 207 FORMAT(1X, 'NEXT, THE NAMES OF THE VARIABLES IN THE FILE: ')
450 WRITE(6, 208)
460 READ(5, 101) (IFORM(II), II=1, 55)
470 208 FORMAT(1X, 'NEXT, THE FORMAT TO READ THE VARIABLES IN THE FILE: ')
480 WRITE(6, 209)
490 READ(5, 101) (NAMSYS(II), II=1, 8)
500 WRITE(6, 500)
510 READ(5, 101) (NAMSUB(II), II=1, 8)
520 209 FORMAT(1X, 'LAST, THE NAME OF THE SPSS SYSTEM FILE TO READ: ')
530 500 FORMAT(1X, 'AND THE NAME OF THE SUBFILE (??)')
540 WRITE(6, 210) (NAMSYS(II), II=1, 8), (NAMSYS(II), II=1, 8)

```

550 210 FORMAT(
560   ' //XXXXXX JOB 04423,STALCUP,CLASS=A',
570   '+//,JOPPARM SKIP=YES,FILE=2
580   '+// EXEC SPSS,DSN='&OUTCASE',UNITP=SYSDA,
590   '+// STATUS=NEW,DISP9=PASS,TIM9=600,'
600   '+// DCR9= '' (RECFM=FH,LRECL=400,BLKSIZE=4000) '' ,DENS=' ,8A1
610   '+GET FILE ' .8A1)
620   WRITE(9,591) (NAMSUR(II),II=1,9)
630 501 FORMAT(
640   'RUN SUFFLES ' ,8A1)
650   WRITE(9,211) (IFORM(II),II=1,55),(NAME(II),II=1,55)
660 211 FORMAT(
670   '+WRITE CASES ' ,55A1
680   '+'
690   WRITE(9,212) (IFORM(II),II=1,55)
700 212 FORMAT(
710   '+FINISH'
720   '+// EXEC FORT
730   '+// DIMENSION IDATA(200),DATA(200)'
740   '+// 200 FORMAT' ,55A1)
750   WRITE(9,213) INVARS
760 213 FORMAT(
770   '+IREC=0
780   '+INVARS=' ,I5)
790   WRITE(9,214)
800 214 FORMAT(
810   '+ 5 CONTINUE
820   '+// READ(11,200,END=99) (DATA(I),I=1,INVARS)'
830   '+// DO 10 I=1,INVARS
840   '+// IDATA(I)=DATA(I)
850   '+// 10 CONTINUE
860   '+// WRITE(12) (IDATA(I),I=1,INVARS)
870   '+// IREC=IREC+1
880   '+// GO TO 5

```

```

860      +/' 99 WRITE(6,300) IREC
870      +/' 300 FORMAT(''1''//'
880      +/'      +25X,          PROCESSING ENDED WITH //:
890      +/'      +25X,          //:
900      +/'      +25X,          NORMAL RETURN CODE //:
910      +/'      +25X,          //:
920      +/'      +25X,          # RECORDS=====),I5)
930      +/'      STOP          ')
940      WRITE(9,215) NVARS
950      215 FORMAT(
960      +'      END
970      +/' // EXEC GOFORT
980      +/' //FT11F001 DD DSN=&&OUTCASE,DISP=(OLD,DELETE)
990      +/' //FT12F001 DD DSN=USER.CASE$(OUT),DISP=SHR
1000      +/' /* UNIT=SYSDA,
1001      +/' /* DCR=(RECFM=VBS,LRECL=800,RLKSIZE=804)
1010      +/' /* SPACE=(TRK,(200,5),RLSE),VOL=SER=ACAD01
1020      +/' // EXEC FORT
1030      +/' DIMENSION IDATA(200)
1040      +/' DATA INC1,INC2,INC3,INC4/-9,-8,-7,-6/
1050      +/' IREC=0
1060      +/' IREC1=0
1070      +/' IREC2=0
1080      +/' IFLAG=0
1090      +/' NREC=,I5
1100      +/' 5 READ(19,END=99) (IDATA(I),I=1,NREC))
1110      WRITE(7,216)
1120      216 FORMAT(
1130      +'      IREC=IREC+1
1140      +'      DO 10 I=1,NREC
1150      +'      IF(IDATA(I).EQ.INC1) IFLAG=IFLAG+1
1160      +'      IF(IDATA(I).EQ.INC2) IFLAG=IFLAG+1
1170      +'      IF(IDATA(I).EQ.INC3) IFLAG=IFLAG+1
1180      +'      IF(IDATA(I).EQ.INC4) IFLAG=IFLAG+1
1190      +'      10 CONTINUE

```

1200 +/- IF(IFLAG.EQ.0) GO TO 20
1210 +/- IF(IFLAG.LE.7) GO TO 25
1220 +/- IFLAG=0
1230 +/- GO TO 20
1240 +/- 20 WRITE(12) (IDATA(I),I=1,NREC)
1250 +/- IREC1=IREC1+1
1260 +/- WRITE(9,217)
1270 217 FORMAT(
1280 +/- IFLAG=0
1290 +/- GO TO 5
1300 +/- 25 WRITE(11) (IDATA(I),I=1,NREC)
1310 +/- IREC2=IREC2+1
1320 +/- IFLAG=0
1330 +/- GO TO 5
1340 +/- 99 CONTINUE
1350 +/- IREC3=IREC-(IREC1+IREC2)
1360 +/- WRITE(6,202) IREC,IREC1,IREC2,IREC3
1370 +/- 202 FORMAT(1X,I6,'' RECORDS PROCESSED'')
1380 +/- +1X,I6,'' RECORDS COMPLETE'')
1390 +/- +1X,I6,'' RECORDS INCOMPLETE'')
1400 +/- +1X,I6,'' RECORDS DISCARDED'')
1410 +/- WRITE(9,218)
1420 218 FORMAT(
1430 +/- STOP
1440 +/- END
1450 +/- // EXEC GOFORT
1460 +/- //FT10F001 DD DSN=USER.CASES(OUT).DISP=SHR

1370 //> F111F991 DD DSN=USER.LIN.CASES,DISP=SHR
1380 UNIT=SYSDA,
1390 DCB=(RECFM=VRS,LRECL=800,BLKSIZE=804),
1400 SPACE=(TRK,(100,5),RLSE),VOL=SER=ACAD01
1410 //> F112F991 DD DSN=USER.LIN.CASES,DISP=SHR
1420 UNIT=SYSDA,
1430 DCB=(RECFM=VBS,LRECL=800,BLKSIZE=804),
1440 SPACE=(TRK,(100,5),RLSE),VOL=SER=ACAD01
1450 ?? CONTINUE
1460 CALL CLEAR
1470 CALL OPSYS('COMMAND','ALL',IRTN,IMESS,PF 12 ''SUB JROUT'' C E)
1480 CALL OPSYS('COMMAND','ALL',IRTN,IMESS,PF 11 ''STATUS'' C E)
1490 CALL OPSYS('COMMAND','ALL',IRTN,IMESS,PF 10 ''RUN INTRO'' C E)
1500 WRITE(6,309)
1510 300 FORMAT(1X,'ALL RIGHT RON, NOW YOU NEED TO SUBMIT THE JOB: ',
1520 +/11X,'PRESS FF KEY #12'
1530 +/11X,'TO CHECK THE STATUS OF THE JOB'
1540 +/11X,'PRESS FF KEY #11'
1550 +/11X,'AND TO START THE NEXT PROGRAM'
1560 +/11X,'PRESS FF KEY #10.//')
1570 STOP
1580 END
1590 JROUTINE CLEAR
1600 WRITE(6,209)
1610 200 FORMAT(1X,'PLEASE CLEAR YOUR SCREEN AND PRESS ENTER')
1620 PEND(5,199)
1630 100 EDINIT(199)
1640 RETURN

1730
 1749
 1759
 1769
 1770
 1789
 1799
 1809
 1819
 1820
 1839
 1849
 1859
 1869
 1879
 1889
 1899
 1909
 1919
 1929
 1939
 1949
 1959
 1969
 1979
 1989
 1999
 2009
 2019
 2029
 2039
 2049
 2059
 2069
 2079
 2089

END
 SUBROUTINE REPLY(IANS)
 INTEGER Y/N, YES, NO, FINE, GOOD, NL, BAD
 DATA YES, NO, OK, FINE, GOOD, BL, BAD /'Y', 'N', 'OK'
 /'FINE', 'GOOD', 'BAD'/
 IF(IANS.EQ.FINE) GO TO 20
 IF(IANS.EQ.GOOD) GO TO 30
 IF(IANS.EQ.BAD) GO TO 40
 IF(IANS.EQ.OK) GO TO 50
 WRITE(6,500)
 GO TO 99
 20 WRITE(6,100)
 GO TO 99
 30 WRITE(6,200)
 GO TO 99
 40 WRITE(6,300)
 GO TO 99
 50 WRITE(6,400)
 99 CONTINUE
 100 FORMAT(1X, 'JUST FINE, HOW WHAT KIND OF ANSWER IS THAT ????')
 200 FORMAT(1X, 'WELL, HE WILL SEE ABOUT CHANGING THAT.')
 300 FORMAT(1X, 'SORRY TO HEAR THAT, PLEASE DON'T TAKE IT OUT ON ME.')
 400 FORMAT(1X, 'WELL DON'T COMMIT YOURSELF TO A POSITIVE ANSWER.')
 500 FORMAT(1X, 'NEXT TIME TRY HARDER, THINK OF A BETTER ANSWER.')
 RETURN
 END
 SUBROUTINE ANSFR(IRN)
 INTEGER YES
 DATA YES, NO /'Y', 'N'/
 WRITE(6,200)
 200 FORMAT(1X, '(Y/N):')
 10 READ(5,100) IANS
 100 FORMAT(1A1)
 IF(IANS.EQ.YES) GO TO 20
 IF(IANS.EQ.NO) GO TO 30
 WRITE(6,291)

2090 201 FURMATE(IX, 'PLEASE ANSWER Y OR N.....')
2100 GO TO 10
2110 20 IRN=9
2120 GO TO 99
2130 30 IRN=1
2140 99 CONTINUE
2150 RETURN
2160 END

INTRO1

10 DIMENSION N(72),NAME(80),IFORM(80),NAMSYS(8),IP(80)
20 INTEGER YES,GO
30 DATA YES,NO GO/'Y','N','GO '/
31 WRITE(6,200)
32 200 FORMAT(IX,'READY? IF SO TYPE THE WORD GO.')
33 100 FORMAT(1A4)
1610 59 READ(5,100) IANS
1620 IF(IANS.EQ.GO) GO TO 55
1630 GO TO 50
1640 55 CONTINUE
1650 CALL OPSYS('COMMAND','ALL',IRTN,IMESS,'FILE JOROUT')
1660 CALL OPSYS('ALLOC','2078 REPLACE',10)
1661 CALL OPSYS('ALLOC','JROUT',9)
1670 REWIND 9
1671 DO 60 J=1,4
1672 110 FORMAT(89A1)
1680 READ(10,110) (IP(I),I=1,80)
1690 WRITE(9,110) (IP(I),I=1,80)

```

1700 60 CONTINUE
1701 WRITE(6,201)
1702 291 FORMAT(1X,'ENTER THE NUMBER OF VARIABLES: ')
1703 READ(5,*), NVARS
1704 WRITE(6,303) NVARS
1705 303 FORMAT(1X,' NVARS=',I5)
1706 65 CONTINUE
1707 READ(10,110,END=88),(IF(I),I=1,80)
1708 WRITE(9,110),(IF(I),I=1,80)
1709 GO TO 65
1710 88 CALL CLEAR
1711 WRITE(6,304)
1712 304 FORMAT(//1X,'NOW THE JOB FOR THE MLR IS READY.')
1713 +//1X,'PRESS PF KEY # 12 TO SUBMIT THE JOB'
1714 +//1X,'PRESS PF KEY # 11 TO CHECK THE STATUS OF THE JOB')
1715 WRITE(6,305)
1716 305 FORMAT(//1X,'CONTROL RETURNED TO YOU, OVER AND OUT.....')
1717 STOP
1718 END
1719 SUBROUTINE CLEAR
1720 WRITE(6,209)
1721 209 FORMAT(1X,'PLEASE CLEAR YOUR SCREEN AND PRESS ENTER.....')
1722 READ(5,100)
1723 100 FORMAT(1A1)
1724 RETURN
1725 END
1726 SUBROUTINE REPLY(IANS)
1727 INTEGER*4 YES,OK,FINE,GOOD,BL,BAD
1728 DATA YES,NO,OK,FINE,GOOD,BL,BAD/'Y','N','OK',
1729 +'FINE','GOOD','BAD',/
1730 IF(IANS.EQ.FINE) GO TO 20
1731 IF(IANS.EQ.GOOD) GO TO 30
1732 IF(IANS.EQ.BAD) GO TO 40
1733 IF(IANS.EQ.OK) GO TO 50
1734 WRITE(6,500)

```

```
210 GO TO 99
220 WRITE(6,100)
230 GO TO 99
240 WRITE(6,200)
250 GO TO 99
260 WRITE(6,300)
270 GO TO 99
280 WRITE(6,400)
290 CONTINUE
300 FORMAT(1X,'JUST FINE, NOW WHAT KIND OF ANSWER IS THAT ????')
310 FORMAT(1X,'WELL, WE WILL SEE ABOUT CHANGING THAT.')
320 FORMAT(1X,'SORRY TO HEAR THAT, PLEASE DON'T TAKE IT OUT ON ME')
330 FORMAT(1X,'WELL DON'T COMMIT YOURSELF TO A POSITIVE ANSWER.')
340 FORMAT(1X,'NEXT TIME TRY HARDER, THINK OF A BETTER ANSWER.')
350 RETURN
360 END
370 SUBROUTINE ANSWER(IRN)
380 INTEGER IES
390 DATA YES,NO/'Y','N'/
400 WRITE(6,200)
410 FORMAT(1X,'(Y/N):')
420 READ(5,100) IANS
430 FORMAT(1A1)
```

```
2449 IF (TRANS.EQ.YEE) GO TO 24  
2450 IF (TRANS.EQ.NO) GO TO 30  
2460 WRITE(6,261)  
2470 201 FORMAT(IX,'PLEASE ANSWER Y OR N.....')  
2480 GO TO 10  
2499 20 IYN=0  
2500 GO TO ??  
2510 30 IYN=1  
2520 99 CONTINUE  
2530 RETURN  
2540 END
```

Batch Version

Both INTRO and INTROL are documented below in a form suitable for batch entry via cards. The cards which must be changed in each instance are noted.

```

//JOOPARM SKIP=YES,TIRE=2
// EXEC SPSS,DSN9='  OUTCASE',UNIT9=SYSDA,
// STATUS9=NEW,DISP9=PASS,TRK9=600,
// DCB9='(RECFM=FB,LRECL=400,BLKSIZE=4000)',DSN3=ABOD
GET FILE ABOD
RUN SUBFILES S10
WRITE CASES (67F3.0)
      SCHOOL,ITEM01,ITEM66,AGE,GRADE
FINISH
// EXEC FORT
      DIMENSION IDATA(200),DATA(200)
200 FORMAT(67F3.0)
      IREC=0
      INvars= 67
      S CONTINUE
      READ(11,200,END=99) (DATA(I),I=1,INvars)
      DO 10 I=1,INvars
      IDATA(I)=DATA(I)
10 CONTINUE
      WRITE(12) (IDATA(I),I=1,INvars)
      IREC=IREC+1
      GO TO S
99 WRITE(6,300) IREC
300 FORMAT('1'///
      +25X,' PROCESSING ZDED WITH   ',/
      +25X,'                                ',/
      +25X,' NORMAL RETURN CODE          ',/
      +25X,'                                ',/
      +25X,' RECORDS=====IS)           ',/
      STOP
      END

```

Computer Center

Initial SPSS file

Format of variables to be read from
INITIAL SPSS file

Number of variables to be read from
INITIAL SPSS file

```

// EXEC GOFOR7
//FT11F001 DD DSN= OUTCASE,DISP=(OLD,DELETE)
//FT12F001 DD DSN=USER.CASES(OUT),DISP=SHR
//3      UNIT=SYSDA,
//4      DCB=(RECFM=VBS,LRECL=800,BLKSIZE=804),
//5      SPACE=(TRK,(200,5),RLSE),VOL=SER=ACAD01
// EXEC FORT
      DIMENSION IDATA(200)
      DATA INC1,INC2,INC3,INC4/-9,-8,-7,-6/
      IREC=0
      IREC1=0
      IREC2=0
      IFLAG=0
      NREC= 67
      S READ(10,END=99) (IDATA(I),I=1,NREC)
      IREC=IREC+1
      DO 10 I=1,NREC
      IF(IDATA(I).EQ.INC1) IFLAG=IFLAG+1
      IF(IDATA(I).EQ.INC2) IFLAG=IFLAG+1
      C   IF(IDATA(I).EQ.INC3) IFLAG=IFLAG+1
      C   IF(IDATA(I).EQ.INC4) IFLAG=IFLAG+1
      10 CONTINUE
      IF(IFLAG.EQ.0) GO TO 20
      IF(IFLAG.LE.7) GO TO 25

```

Computer Center VCL

```
IFLAG=0
GO TO 5
20 WRITE(12) (IDATA(I),I=1,NREC)
IREC1=IREC1+1
IFLAG=0
GO TO 5
25 WRITE(11) (IDATA(I),I=1,NREC)
IREC2=IREC2+1
IFLAG=0
GO TO 5
99 CONTINUE
IREC3=IREC-(IREC1+IREC2)
WRITE(6,202) IREC,IREC1,IREC2,IREC3
202 FORMAT(1X,16,' RECORDS PROCESSED'
          +/1X,16,' RECORDS COMPLETE'
          +/1X,16,' RECORDS INCOMPLETE'
          +/1X,16,' RECORDS DISCARDED')
STOP
END
```

```
// EXEC CP0RT
//FT10F001 DD DSN=USER.CASES(OUT),DISP=SRR
//FT11F001 DD DSN=USER.INC.CASES,DISP=SRR
//*
    UNIT=SISDA,
//*
    DCB=(RECPN=VBS,LRECL=800,BLKSIZE=804),
//*
    SPACE=(TRK,(100,S),RLSE),VOL=SER=ACAD01
//FT12F001 DD DSN=USER.COM.CASES,DISP=SRR
//*
    UNIT=SISDA,
//*
    DCB=(RECPN=VBS,LRECL=800,BLKSIZE=804),
//*
    SPACE=(TRK,(100,S),RLSE),VOL=SER=ACAD01
```

Computer Center VCC

//||||||| JOB 04423,STALCUP,CLASS=C

/>JOBPARM SKIP=YES,TIME=15

// EXEC PORT

COMMON /COM1/ INCON(200)

NVARS= 67

Number of variables

I=0

NNNN=NVARS-1

10 READ(10,END=99) ISC,(INCON(J),J=1,NNNN)

WRITE(6,201) ISC,(INCON(J),J=1,NNNN)

CALL REP(NVARS,ISC)

CALL WRITE(NVARS,ISC)

I=I+1

GO TO 10

99 CONTINUE

WRITE(6,200) I

200 FORMAT('1',' JOB SUMMARY',

'//1E,'JOB FINISHED; COMPLETION CODE=0'

'/1E,'NUMBER OF CASES=',IS)

201 FORMAT(1X,'INCOMPLETE CASE LISTED BELOW:',

+50(/1E,30I4))

STOP

END

C.....

SUBROUTINE REP(BRREC, ISC)
DIMENSION X:(1350,200),TEMP(200).

* I(200,7) , VARS(200) ,IREP(200) ,BRR(6) .
+ INDEX(200)
COMMON /COM1/INDATA(200)
COMMON /COM2/ISUB(1350,200)
COMMON /COM3/IFILE(1350,200)
COMMON /COM4/INFILE(200)
ICOMP=1
IT=ISC
BRREC=BRREC-1
REWIND 11
10 READ(11,END=99) ISC, (INFILE(I), I=1, BRREC)
C WRITE(6,777) ISC, (INFILE(I), I=1, BRREC)
777 FORMAT(1X,30I3/)
IF(IT.EQ.ISC) GO TO 20
GO TO 10
20 CONTINUE
J=0
L=0
C..... J = COUNTER FOR INCOMPLETE VARS
C..... L = COUNTER FOR COMPLETE VARS
C..... IREP = INDEX OF LOCATION FOR COMPLETE VARS
C..... INDEX = INDEX OF LOCATION FOR INCOMPLETE VARS
C..... IFILE = VALUES OF THE DEPENDENT VARS (COMPLETE)
C..... ISUB = VALUES OF THE INDEPENDENT VARS (INCOMPLETE)

```
DO 30 I=1,NNREC
IF(INDATA(I).EQ.-9.0R.INDATA(I).EQ.-8) GO TO 40
GO TO 50
40 J=J+1
INDEX(J)=I
GO TO 30
50 CONTINUE
L=L+1
IREP(L)=I
30 CONTINUE
N=0
35 N=N+1
DO 60 K=1,J
IFILE(N,K)=INFILE(INDEX(K))
60 CONTINUE
DO 70 K=1,L
ISUB(N,K)=INFILE(IREP(K))
70 CONTINUE
READ(11,ERR=88) ISC,(INFILE(LL),LL=1,NNREC)
IF(IT.EQ.ISC) GO TO 35
LL=L+J
L1=L+1
N1=L+1
RTN=N1-N1
C WRITE(6,888)
CALL SUB(J,L,N,NVARS,IREP,INDEX,IT,8,LL,L1,N1)
RETURN
88 CONTINUE
C WRITE(6,666)
```

```

555 FORMAT(1X,'INCRC= ',IS)
666 FORMAT(1X,'CALL AT END')
888 FORMAT(1X,'CALL AT RID')
LL=L+J
L1=L+1
R1=L+1
BTE=B1+R1
CALL SUB(J,L,N,IVARS,IREP,INDEX,IT,B,LL,L1,R1)
99 CONTINUE
RETURN
END

C.....SUBROUTINE SUB(J,L,N,IVARS,IREP,INDEX,IT,B,LL,L1,R1)
      DIMENSION IT(N,R1),TBEP(200),
     * ITBAR(200),A(5400),ABOVA(14),B(R1,7),VARB(2100),IREP(200),
     * RER(6),INDEX(200)
      COMMON /COM1/IBDATA(200)
      COMMON /COM2/ISUB(1350,200)
      COMMON /COM3/FILE(1350,200)
      COMMON /COM4/IFILE(200)

C.....      J = COUNTER FOR INCOMPLETE VARS
C.....      L = COUNTER FOR COMPLETE VARS
C.....      N = COUNTER FOR NUMBER OF CASES
C.....      IREP = INDEX OF LOCATION FOR COMPLETE VARS
C.....      INDEX = INDEX OF LOCATION FOR INCOMPLETE VARS
C.....      FILE = VALUES OF THE DEPENDENT VARS (COMPLETE)
C.....      ISUB = VALUES OF THE INDEPENDENT VARS (INCOMPLETE)
C.....      WRITE(6,666) B,L,J

```



```

      BBB(5) = 1
      BBB(6) = 1
      ALFA = 0.10
      IB = 01
      WRITE(6,555)
C      555 FORMAT(1X,'JUST BEFORE THE CALL',I10)
      CALL BECOVR (IT,IX,BBB,TEMP,XIBAB,I,IEB)

      CALL RLEBL (A,XIBAB,BB,BB,ALFA,ANOVA,B,IB,VABD,IEB)
C      WRITE(6,666) B(1,1),B(BB+1,1)
      SUM=B(BB+1,1)
      DO 30 II=1,BB
      SUM=SUM+B(II,1)+IBDATA(IREP(II))
30  CONTINUE
      CALL ROUND (SUM,IS)
      IBDATA(IINDEX(ICOB))=IS
      GO TO 5
999  CONTINUE
      RETURN
      END
C.....SUBROUTINE ROUND(I,I)
      I=X
      IX=I*100
      IX=IX*100
      Y=IX-IX
      IF(Y.GE.45.0.AND.Y.LE.55.0) GO TO 10
      IF(Y.LT.45.0) RETURN
      IF(Y.GT.55.0) I=I+1
      RETURN

```

```
200 FORMAT(1X,'FINISHED CASE LISTED BELOW:',  
+50(/1X,30I4))  
RETURN  
END  
// EXEC GOFORT,GOSIZE=3500K  
//FT10F001 DD DSN=USER.INC.CASES,DISP=SHR  
//FT11F001 DD DSN=USER.COM.CASES,DISP=SHR  
//FT16F001 DD DSN=USER.NAL01,DISP=SHR  
// EXEC FORT  
DIMENSION N(80)  
K=0  
1 READ(12,END=88) (N(J),J=1,65)  
WRITE(16,100) (N(I),I=1,65) Number of variables  
K=K+1  
GO TO 1  
88 CONTINUE  
5 READ(13,100,END=99) (N(I),I=1,65) Number of variables  
WRITE(16,100) (N(I),I=1,65)  
K=K+1  
GO TO 5  
99 CONTINUE  
WRITE(6,201) K  
100 FORMAT(80I3)  
200 FORMAT(/1X,3(1X,30I3/))  
201 FORMAT(///1X,' OF CASES ADDED=',I6)  
STOP  
END  
// EXEC GOFORT  
//FT12F001 DD DSN=USER.COM.CASES,DISP=SHR  
//FT13F001 DD DSN=USER.NAL01,DISP=SHR  
//FT16F001 DD DSN=USER.FINISH.CASES,DISP=(MOD,KEEP)  
//
```

Prediction of Time Required

The programming requires fairly large amounts of CPU time when dealing with large files. The total time required is a function of the number of variables, the number of cases, the number of cases with missing data, and the interaction between the number of missing cases and the total number of cases.

To date, run time data have been gathered on 23 files which ranged in size from an $N = 122$ to $N = 1,081$. The number of missing data cases ranged from $N = 12$ to $N = 191$. All the files had 67 variables. Required CPU times ranged from 2.145 minutes to 203.072 minutes.

This information was used to build to test a regression equation which could be used to predict the CPU time required. The $R^2 = .9960$ had a $p < .00000$, but the equation is limited by the fact that all the test files had 67 variables. Even considering this limitation, the following regression equation may be useful in predicting CPU time.

$$\begin{aligned} \text{CPUT} = & (.22248840)u + (.00779903)\text{NC} + (.08595389)\text{NM} \\ & + (.00100721)\text{NC} * \text{NM} + E \end{aligned}$$

Where:

- CPUT = CPU time required
- U = Unit vector (+)
- NC = Number of cases
- NM = Number of cases with 10% missing data
- E = Error vector (ignored when predicting CPUT)

References

Cohen, J., & Cohen, P. Applied multiple regression/correlation analysis for the behavioral sciences. New York: John Wiley & Sons, 1975.

International Mathematical and Statistical Libraries, Inc.

The IMSL library, "RLMUL" (vol. 3). Houston, Texas:

International Mathematical and Statistical Libraries, Inc.,

June 1980.

Newman, I. Variations between shrinkage estimation formulas and the appropriateness of their interpretation. Multiple Linear Regression Viewpoints, 1973, 4, 45-48.

Newman, I., & Fry, J. A response to "A note on multiple comparisons" and comment on shrinkage. Multiple Linear Regression Viewpoints, 1972, 2, 36-39.