

## Application of NORM to the Multiple Imputation for Multivariate Missing Data.

Hyun-Jeong Kim<sup>1)</sup>, Sung-Ho Moon<sup>2)</sup>, Jae-Kyoung Shin<sup>3)</sup>

### Abstract

The statistical analysis of incomplete data sometimes requires handling of incomplete observations. Towards this end, each case with some missing values generally should be deleted, namely, resulting in only use of non-missing cases. EM algorithm(Dempster et al., 1977) which involves prediction and estimation steps is a general method among others. In this article, we use the free software NORM developed for multiple imputation, which uses DA(Data Augmentation) algorithm in its imputation, and evaluate its efficiency through a numerical example.

**Keywords** : Hot-deck, EM-algorithm, NORM, Multiple Imputation

### 1.

가  
가 , i)  
( , ii)  
( 가 ), iii)  
iii)  
(MCAR, missing completely at random)  
가 ,  
(MAR, missing at random)

---

1 1- 1  
E-mail: semikim@silla.ac.kr  
2 55- 1  
3 9

가 (Little and Rubin, 1987).

Little and Rubin(1987) Johnson and Wichern(1992) iv) EM  
(Dempster et. al., 1977) (Orchard and Woodbury, 1972)  
, 가

EM ( , single imputation)  
가  
i) iii) iv)

(multiple imputation) i) iii)  
가 iv)

Markov Chain Monte Carlo(MCMC) Data Augmentation(DA)  
NORM(Schafer, 1997a, Chapter

5) NORM 가  
CAT(Schafer, 1997a, Chapter 7-8), (general location  
model) 가 MIX(Schafer, 1997a, Chapter 9)  
(multivariate linear mixed-effects model) 가  
PAN(Schafer, 1997b) , 가

2 NORM , 3  
1, 2  
NORM , .

**2. NORM**

NORM 1) , , ( )  
EM , 2) DA 가  
3) 가 ( ) , 4)  
DA , (plot), 5)  
(overall) Rubin(1987)

NORM  
 $m (>1)$  . (complete data : 1 )  
 $m$  .  $m$  (standard  
complete data : 6 ) .  $m$   
Rubin

( 5 10 )

NORM Markov Chain Monte Carlo (MCMC)

DA

2.1 NORM

NORM Data, EM algorithm, Data Augmentation, Impute from parameters

4

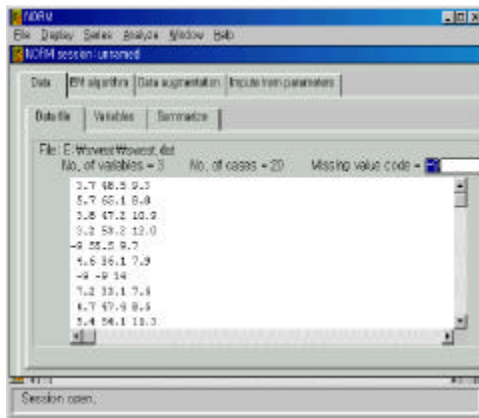
(1) Data

```

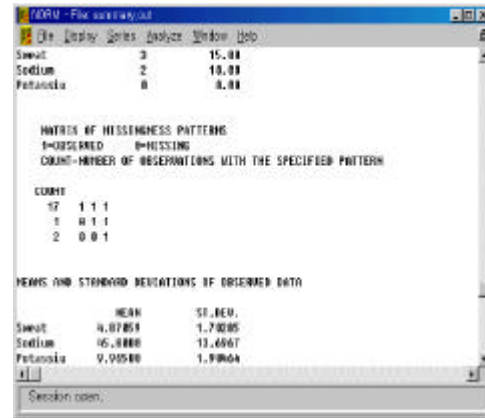
NORM file New session
import( ) ( ASCII) '*.dat'
Enter
. Data file
, 1 -9
. Variables
, Transformation none
*.imp (default *.imp
).

```

Summarize



1. NORM



2. Summary Out

Run summary.out  
default file 2

(2) EM algorithm

가 EM

(Dempster et al., 1977) . EM  
 (estimation) (maximization)

. EM  
 output file(\*.out) EM

parameter file(\*.prm) . EM output file  
 parameter file (Display)

. Computing option  
 (Maximum iterations) (Convergence criterion)

. 3 Maximum iterations 1000 , EM 1000  
 1000 MLE Posterior mode

, EM default MLE가  
 EM DA

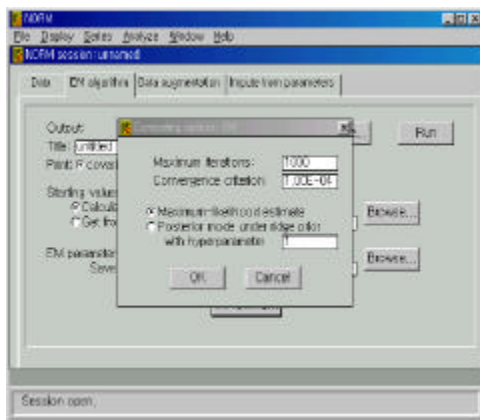
EM (Rate of missing information: Rubin, 1987)

(3) Data Augmentation (DA)  
 NORM DA

. DA  
 MCMC , DA EM

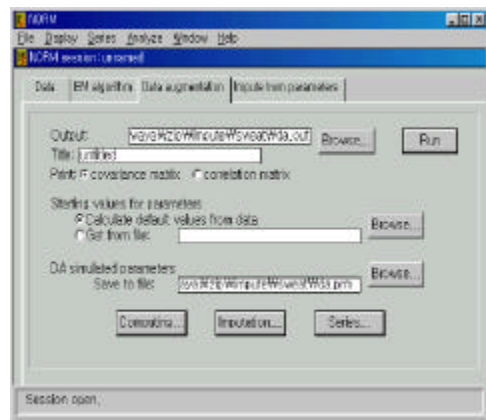
DA 4 Run Computing, Imputation, Series 3가  
 . Run button

(output file(\*.out)  
 parameter file(\*.prm) ). Output  
 da.out



3. EM

parameter file



4. Data Augmentation

. DA Computing options (Number

of iterations), (Random seed), (prior distribution)

. Imputation option  
 \*.imp file View imputation immediately  
 \*\_n.imp file

Series DA plotting parameter series(\*.prs) parameters

series(\*.prs) 가 . Save all parameters at every k-th cycle : ( , , )가 \*.prs . DA

k 가가 \*.prs file . Save only worst linear function of parameters : 가 ( ) \*.prs

series \*.prs Series plot Worst linear function Autocorrelation function . Series Save only worst linear function .

(4) Impute from parameter sheet

Impute from parameters NORM

. NORM parameter(\*.prm) NORM parameter series (\*.prs) output(\*.out) (\*.imp) . , , 가 . NORM DA . 5 Run parameter(\*.prm) file parameter series(\*.prs) file . Computing (seed) .

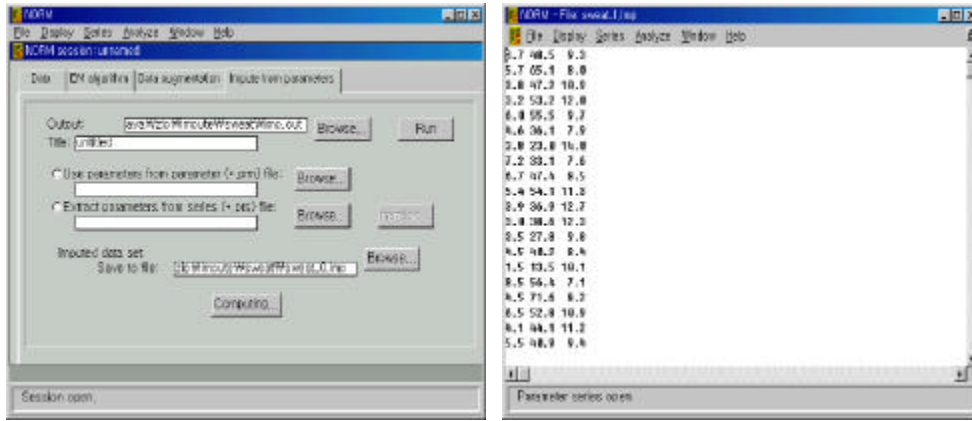
**2.2 NORM**

imp.out 6 1 display imputation default \*\_n.imp . n . NORM 3.2

NORM

NORM

URL <http://www.stat.psu.edu/~jls/misoftwa.html>



5. Impute from parameters sheet

6.

(1997a, 1997b, 1999) Tanner and Wong(1987), Schafer and Olsen(1998)

3.

( $X_1$ : sweat rate,  $X_2$ : sodium content,  $X_3$ : potassium content) 20 ( $p = 3, n = 20$ )

3.1  
5, \*가 가 (Efron, 1994).

$$\theta = \text{maximum eigenvalue of } \hat{\Sigma}$$

$\hat{\Sigma}$

가 가

(hot-deck)

$$\hat{\theta} = 193.70,$$

$$\hat{\theta} = 173.84$$

(mean imputation)  $\hat{\theta} = 170.40$

가

$$\theta = 200.46$$

, (estimation) (maximization) EM  
 (Dempster et. al., 1977) 179.15 .  
 EM  
 NORM DA 2.1 .

3.1. Sweat Data

	$X_1$	$X_2$	$X_3$
Individual	Sweat rate	Sodium	Potassium
1	3.7	48.5	9.3
2	5.7	65.1	8.0
3	3.8	47.2	10.9
4	3.2	53.2	12.0
5	3.1*	55.5	9.7
6	4.6	36.1	7.9
7	2.4*	24.8*	14
8	7.2	33.1	7.6
9	6.7	47.4	8.5
10	5.4	54.1	11.3
11	3.9	36.9	12.7
12	4.5*	58.8*	12.3
13	3.5	27.8	9.8
14	4.5	40.2	8.4
15	1.5	13.5	10.1
16	8.5	56.4	7.1
17	4.5	71.6	8.2
18	6.5	52.8	10.9
19	4.1	44.1	11.2
20	5.5	40.9	9.4

Source : Johnson and Wichern(1992).

3.2

	$\hat{\theta}$
	$\theta = 200.46$
	193.70
hot - deck	173.84
	170.40
EM	179.15
NORM	196.58

NORM

1\_\_\_\_(Data). “-9” 가 1  
 . Variables VAR\_1  
 , Summarize 가  
 Run  
 2\_\_\_\_(EM). path default Computing  
 ( 3 ).  
 3\_\_\_\_(DA). Computing , Imputation Series  
 Run ( 4 ).  
 4\_\_\_\_(Impute). impute Run  
 Display impute file \*\_n.imp  
 6 n  
 가 .  
 NORM  
 196.58 . EM DA  
 가 .  
 (2001)  
 , 1 NORM S/W CAT,  
 MIX, PAN .

1. Dempster, A. P., Laird, N. M. and Rubin, D. B. (1977). Maximum-likelihood estimation from incomplete data via the EM algorithm (with discussion). *Journal of the Royal Statistical Society Series*, B39, 1-38.
2. Efron, B. (1994). Missing Data, Imputation, and Bootstrap, *Journal of the American Statistical Association*, Vol 89, No 426, 463-479.
3. Johnson, R. A. and Wichern, D. W. (1992). *Applied Multivariate Statistical Analysis 3rd Ed.*, Prentice Hall
4. Little, R. J. A. and Rubin, D. B. (1987). *Statistical Analysis with Missing Data*, John Wiley and Sons.
5. Orchard, T. and Woodbury, M. A. (1972). A missing information principle: Theory and application, *Proc. 6th Berkeley Symposium on Math. Statist. and Prob.* 1, 697-715.
6. Rubin, D. B. (1987). *Multiple Imputation for Nonresponse in Surveys*. New York: J. Wiley.
7. Schafer, J. L. (1997). *Analysis of Incomplete Multivariate Data*. London: Chapman & Hall, Chapter 3.
8. Schafer, J. L. and Olsen, M. K. (1998). Multiple imputation for multivariate missing-data problems: a data analyst and perspective. *Multivariate*



- Behavioral Research*, 33, 545-571.
9. Schafer, J. L. (1999). Multiple imputation: a primer. *Statistical Methods in Medical Research*, 8(1), 3-15.
  10. Tanner, M. A. and Wong, W. H. (1987). The calculation of posterior distributions by data augmentation (with discussion). *Journal of the American Statistical Association*, 82, 528-550.
  11. \_\_\_\_\_, \_\_\_\_\_ (2001). “ \_\_\_\_\_ ”, *\_\_\_\_\_* 2001 \_\_\_\_\_, 101-104.

[ 2002 8 \_\_\_\_\_, 2002 9 \_\_\_\_\_ ]