

Missing Data Using Stata

Paul Allison, Ph.D.

Upcoming Seminar:
August 15-16, 2017, Stockholm, Sweden

- 1 **Missing Data Using Stata**
- 2 **Basics**
- 3 **For Further Reading**
- 4 **Many Methods**
- 5 **Assumptions**
- 6 **Assumptions**
- 7 **Ignorability**
- 8 **Assumptions**
- 9 **Listwise Deletion (Complete Case)**
- 10 **Listwise Deletion (continued)**
- 11 **Listwise Deletion (continued)**
- 12 **Pairwise Deletion (Available Case)**
- 13 **Dummy Variable Adjustment**
- 14 **Imputation**
- 15 **Maximum Likelihood**
- 16 **Properties of Maximum Likelihood**
- 17 **ML with Ignorable Missing Data**
- 18 **ML for 2 x 2 Contingency Table**
- 19 **Maximizing the Likelihood with EM**
- 20 **ML for Quantitative Variables**
- 21 **EM Algorithm**
- 22 **EM for Multivariate Normal Data**
- 23 **EM for Multivariate Normal Data**
- 24 **College Example**
- 25 **Preliminary Analysis 1**
- 26 **Preliminary Analysis 2**
- 27 **Preliminary Analysis 3**
- 28 **EM in Stata**
- 29 **Convert Covariances to Correlations**

- 30 **EM As Input to regress**
- 31 **Direct ML**
- 32 **Direct ML**
- 33 **Direct ML (cont.)**
- 34 **SEM without Auxiliary Variable**
- 35 **SEM with Auxiliary Variable**
- 36 **SEM Output with Auxiliary Variable**
- 37 **Compare with Listwise Deletion**
- 38 **Regression with Mplus**
- 39 **Regression with Mplus**
- 40 **Logistic Regression with Mplus**
- 41 **Other Capabilities of Mplus**
- 42 **ML for Repeated Measures Data**
- 43 **Binary Example**
- 44 **Estimation in Stata**
- 45 **Figure 1**
- 46 **Limitations of Maximum Likelihood**
- 47 **Multiple Imputation**
- 48 **Regression Imputation**
- 49 **Adding a Random Component**
- 50 **Multiple, Random Imputations**
- 51 **Combining the Imputations**
- 52 **Formula for Standard Error**
- 53 **Random Variation in Parameters**
- 54 **Monotonic Missing Data**
- 55 **MI for Monotone Missing Data**
- 56 **Non-Monotone Missing Data**
- 57 **Two Iterative Solutions**
- 58 **MCMC**

- 59 **MCMC for Multivariate Normal**
- 60 **Software**
- 61 **Steps for MCMC in Stata**
- 62 **MCMC With Stata**
- 63 **Stata Output 1**
- 64 **Stata Output 2**
- 65 **Formulas**
- 66 **Imputation with the Dependent Variable**
- 67 **Should Missing Data on the Dependent Variable Be Imputed?**
- 68 **How Many Data Sets?**
- 69 **Options for mi impute mvn**
- 70 **Change the Number of Iterations**
- 71 **Change the Prior Distribution**
- 72 **Categorical Variables**
- 73 **Categorical Variables (cont.)**
- 74 **Some Things NOT to Do**
- 75 **Fully Conditional Specification**
- 76 **Logit Imputation of a Binary Variable**
- 77 **Predictive Mean Matching**
- 78 **Fill-In Phase of FCS**
- 79 **Imputation Phase of FCS**
- 80 **Downside of FCS**
- 81 **Software**
- 82 **FCS in Stata for NLSY Data**
- 83 **Impute Output**
- 84 **Estimate Output**
- 85 **Test Output**
- 86 **mi estimate with Other Commands**
- 87 **Multi-Parameter Inference**

- 88 **Restricted FMI Test**
- 89 **Unrestricted FMI Test**
- 90 **mi test command**
- 91 **Combining Chi-Squares**
- 92 **Stats Not Reported by mi estimate**
- 93 **mibeta for R-square & Standardized**
- 94 **mibeta Output**
- 95 **Interactions and Nonlinearities**
- 96 **Interaction Results**
- 97 **Imputation Model vs. Analysis Model**
- 98 **MI for Panel Data**
- 99 **Hip Fracture Example**
- 100 **Imputing Clustered Data in Stata**
- 101 **Imputation with Cluster Dummies**
- 102 **Imputation in Wide Form**
- 103 **Imputation Via Random Effects**
- 104 **Hip Fracture Example (cont.)**
- 105 **Why Didn't Imputation Do Better?**
- 106 **Nonignorable Missing Data**
- 107 **Nonignorable Missing Data**
- 108 **Heckman's Model for Selection Bias**
- 109 **Heckman's Model in Stata**
- 110 **Heckman's Model (cont.)**
- 111 **Pattern-Mixture Models with MI**
- 112 **MI for Pattern-Mixture Models**
- 113 **Summary and Review**
- 114 **Summary and Review**

Missing Data Using Stata

Paul D. Allison, Ph.D.
February 2016

www.StatisticalHorizons.com

1

Basics

Definition: Data are missing on some variables for some observations

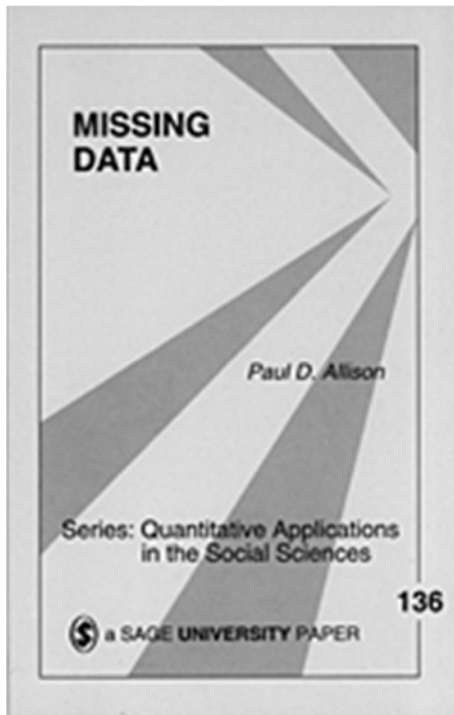
Problem: How to do statistical analysis when data are missing? Three goals:

- ▣ Minimize bias
- ▣ Maximize use of available information
- ▣ Get good estimates of uncertainty

NOT a goal: imputed values “close” to real values.

2

For Further Reading



Also:

Allison, Paul D. (2009) "Missing Data." Pp. 72-89 in The SAGE Handbook of Quantitative Methods in Psychology, edited by Roger E. Millsap and Alberto Maydeu-Olivares. Thousand Oaks, CA: Sage Publications Inc.
<http://statisticalhorizons.com/wp-content/uploads/2012/01/Milsap-Allison.pdf>

3

Many Methods

□ Conventional

- Listwise deletion (complete case analysis)
- Pairwise deletion (available case analysis)
- Dummy variable adjustment
- Imputation
 - Replacement with means
 - Regression
 - Hot deck

□ Novel

- Maximum likelihood
- Multiple imputation
- Inverse probability weighting (not discussed here)

4

Assumptions

Missing completely at random (MCAR)

Suppose some data are missing on Y . These data are said to be MCAR if the probability that Y is missing is unrelated to Y or other variables X (where X is a vector of observed variables).

$$\Pr(Y \text{ is missing} | X, Y) = \Pr(Y \text{ is missing})$$

- MCAR is the ideal situation.
- What variables must be in the X vector? Only variables in the model of interest.
- If data are MCAR, complete data subsample is a random sample from original target sample.
- MCAR allows for the possibility that *missingness* on one variable may be related to *missingness* on another
 - e.g., sets of variables may always be missing together

5

Assumptions

Missing at random (MAR)

Data on Y are missing at random if the probability that Y is missing does not depend on the value of Y , after controlling for observed variables

$$\Pr(Y \text{ is missing} | X, Y) = \Pr(Y \text{ is missing} | X)$$

E.g., the probability of missing income depends on marital status, but within each marital status, the probability of missing income does not depend on income.

- Considerably weaker assumption than MCAR
- Only X 's in the model must be considered. *But*, including other X 's (correlated with Y) can make MAR more plausible.
- *Can* test whether missingness on Y depends on X
- *Cannot* test whether missingness on Y depends on Y

6

Ignorability

The missing data mechanism is said to be ignorable if

- The data are missing at random *and*
 - Parameters that govern the missing data mechanism are distinct from parameters to be estimated (unlikely to be violated)
-

- In practice, "MAR" and "ignorable" are used interchangeably
- If MAR but not ignorable (parameters not distinct), methods assuming ignorability would still be good, just not optimal.
- If the missing data mechanism is ignorable, there is no need to model it.
- Any general purpose method for handling missing data must assume that the missing data mechanism is ignorable.

7

Assumptions

Not missing at random (NMAR)

If the MAR assumption is violated, the missing data mechanism must be modeled to get good parameter estimates.

Heckman's regression model for sample selection bias is a good example.

Effective estimation for NMAR missing data requires very good prior knowledge about missing data mechanism.

- Data contain no information about what models would be appropriate
- No way to test goodness of fit of missing data model
- Results often very sensitive to choice of model
- Listwise deletion able to handle one important kind of NMAR

8

Listwise Deletion (Complete Case)

Delete any unit with any missing data (only use complete cases)

Strengths

- Easy to implement
- Works for any kind of statistical analysis
- If data are MCAR, does not introduce any bias in parameter estimates
- Standard error estimates are appropriate

9

Listwise Deletion (continued)

Weaknesses

- May delete a large proportion of cases, resulting in loss of statistical power
- May introduce bias if MAR but not MCAR

Robust to NMAR for predictor variables in regression analysis

Let Y be the dependent variable in a regression (any kind) and X one of the predictors. Suppose

$$\Pr(X \text{ is missing} | X, Y) = \Pr(X \text{ is missing} | X)$$

Then listwise deletion will not introduce bias.

10

Listwise Deletion (continued)

Example: Estimate a regression with number of children as dependent variable and income as an independent variable.

- ❑ 30% of cases have missing data on income, persons with high income are less likely to report income
- ❑ But probability of missing income does not depend on number of children
- ❑ Then listwise deletion will not introduce any bias into estimates of regression coefficients

For logistic regression, listwise deletion is robust to NMAR on independent OR dependent variable (but not both)

Caveat: This property of listwise deletion presumes that regression coefficients are invariant across subgroups (no omitted interactions).

11

Pairwise Deletion (Available Case)

- ❑ For linear models, parameters are functions of means, variances and covariances (moments)
- ❑ Estimate each moment with all available nonmissing cases
- ❑ Plug moment estimates into formulas for parameters

Strengths:

- ❑ Approximately unbiased if MCAR
- ❑ Uses all available information

Weaknesses:

- ❑ Standard errors incorrect (no appropriate sample size)
- ❑ Biased if MAR but not MCAR
- ❑ May break down (correlation matrix not positive definite)

12

Dummy Variable Adjustment

A popular method for handling missing data on predictors in regression analysis (Cohen and Cohen 1985)

In a regression predicting Y , suppose there is missing data on a predictor X .

1. Create a new variable $D=1$ if X is missing and $D=0$ if X is present.
 2. When X is missing, set $X=c$ where c is some constant (e.g., the mean of X).
 3. Regress Y on both X and D (and any other variables)
- Produces biased coefficient estimates (Jones, JASA, 1996)
 - So does a related method: For categorical variables, create a separate missing data category
 - But may be appropriate for “doesn’t apply” missing data
 - May also be useful for predictive modeling with missing data.

13

Imputation

Any method that substitutes estimated values for missing values

- Replacement with means
- Regression imputation (replace with conditional means)

Problems

- Often leads to biased parameter estimates (e.g., variances)
- Usually leads to standard error estimates that are biased downward
 - Treats imputed data as real data, ignores inherent uncertainty in imputed values.

14

Maximum Likelihood

Choose as parameter estimates those values which, if true, would maximize the probability of observing what has, in fact, been observed.

Likelihood function: Expresses the probability of the data as a function of the data and the unknown parameter values.

Example: Let $p(y|\theta)$ be the probability density for y , given θ (a vector of parameters). For a sample of n independent observations, the likelihood function is

$$L(\theta) = \prod_{i=1}^n p(y_i | \theta)$$

15

Properties of Maximum Likelihood

To get ML estimates, we find the value of θ that maximizes the likelihood function.

Under usual conditions, ML estimates have the following properties:

- Consistent (implies approximately unbiased in large samples)
- Asymptotically efficient
- Asymptotically normal

16

ML with Ignorable Missing Data

Suppose we have 2 discrete variables X and Y , and there is ignorable missing data on X . Let $p(x,y|\theta)$ be the joint probability function.

For a single observation with X missing, the likelihood is

$$g(y|\theta) = \sum_x p(x,y|\theta) = E_x[p(y|x)]$$

The likelihood for the entire sample with m complete cases is

$$L(\theta) = \prod_{i=1}^m p(x_i, y_i | \theta) \prod_{i=m+1}^n g(y_i | \theta)$$

This likelihood may be maximized like any other.

17

ML for 2 x 2 Contingency Table

	<u>Vote</u>		
	Yes	No	
Male	36	37	Furthermore, voting was missing for 10 males and 15 females.
Female	22	52	

The parameters are $p_{11}, p_{12}, p_{21}, p_{22}$. If we exclude cases with missing data, the likelihood is

$$(p_{11})^{36}(p_{12})^{37}(p_{21})^{22}(p_{22})^{52}$$

If we allow for missing data, the likelihood is

$$(p_{11})^{36}(p_{12})^{37}(p_{21})^{22}(p_{22})^{52}(p_{11}+p_{12})^{10}(p_{21}+p_{22})^{15}$$

18

Maximizing the Likelihood with ℓ_{EM}

Freeware for Windows by Jeroen Vermunt:
<http://members.home.nl/jeroenvermunt/>

<u>Input</u>	<u>Output</u>
man 2	* P(sv) *
res 1	1 1 0.2380 (0.0339)
dim 2 2 2	1 2 0.2446 (0.0342)
lab r s v	2 1 0.1538 (0.0297)
sub sv s	2 2 0.3636 (0.0384)
mod sv	
dat [36 37 22 52 10 15]	

ℓ_{EM} fits a large class of models for categorical data, including log-linear, logit, latent class, and discrete time event history models.

19

ML for Quantitative Variables

Assume multivariate normality, which implies

- All variables are normally distributed
- All conditional expectation functions are linear
- All conditional variance functions are homoscedastic

A strong assumption but widely invoked as the basis for multivariate analysis

Several ways to get ML estimates with missing data, based on this assumption

- Factoring the likelihood for monotone missing data patterns
- EM algorithm
- Direct maximization of the likelihood

20

EM Algorithm

A general approach to getting ML estimates with missing data

Two-step procedure

1. Expectation (E): Find the expected value of the log-likelihood for the observed data, based on current parameter values.
2. Maximization (M): Maximize the expected log-likelihood to get new parameter estimates.

Repeat until convergence.

For multivariate normal data, parameters are means, variances, and covariances.

21

EM for Multivariate Normal Data

1. Choose starting values for means and covariance matrix.
2. If data are missing on x , use current values of parameters to calculate the linear regression of x on all variables present for each case.
3. Use linear regressions to impute values of x . (E-step)
4. After all data have been imputed, recalculate means and covariance matrix, with corrections for variances and covariances (*see next slide*). (M-step)
5. Repeat steps 2-4 until convergence.

22

EM for Multivariate Normal Data

Correction: Suppose X was imputed using variables W and Z .

Let $S^2_{x.wz}$ be the residual variance from that regression. Then, in calculating the variance for X , wherever you would use x^2_i , substitute $x^2_i + S^2_{x.wz}$

For covariances between two variables with missing values, there's a similar correction in which you add the residual covariance.

EM algorithm for multivariate normal data is available in many commercial software packages: SPSS, Systat, SAS, Splus, Stata

23

College Example

1994 U.S. News Guide to Best Colleges

- 1302 four-year colleges in U.S.
- Goal: estimate a regression model predicting graduation rate ($\#$ graduating/ $\#$ enrolled 4 years earlier \times 100)
- 98 colleges have missing data on graduation rate

Independent variables:

- 1st year enrollment (logged, 5 cases missing)
- Room & Board Fees (40% missing)
- Student/Faculty Ratio (2 cases missing)
- Private=1, Public=0
- Mean Combined SAT Score (40% missing)

- Auxiliary variable: Mean ACT scores (45% missing)

24

Preliminary Analysis 1

```
use c:\data\college.dta, clear
mi set wide
```

This declares the data to be a missing data set. It also specifies that imputed data are to be stored in the wide format. There are four different storage formats. But how it's stored usually doesn't matter, and we're not imputing yet anyway.

```
mi misstable summarize
```

This requests basic descriptive statistics.

25

Preliminary Analysis 2

Variable	Missing		Not missing		Obs<.	
	Obs=.	Obs>.	Obs<.	Unique values	Min	Max
gradrat	98		1,204	89	8	118
lenroll	5		1,297	>500	2.890372	8.912608
rmbrd	519		783	>500	1.26	8.7
stufac	2		1,300	208	2.3	91.8
csat	523		779	339	600	1410
act	588		714	17	11	31

26

Preliminary Analysis 3

mi misstab patterns

Missing-value patterns (1 means complete)						
Percent	Pattern					
	1	2	3	4	5	6
23%	1	1	1	1	1	1
12	1	1	1	0	1	1
12	1	1	1	1	1	0
12	1	1	1	1	0	0
9	1	1	1	1	0	1
9	1	1	1	0	1	0
8	1	1	1	0	0	0
6	1	1	1	0	0	1
1	1	1	0	0	1	1
1	1	1	0	1	0	0
1	1	1	0	1	1	1

<1	1	1	0	0	0	0
<1	1	1	0	1	0	1
<1	1	1	0	0	0	1
<1	1	1	0	0	1	0
<1	1	1	0	1	1	0
<1	0	0	0	0	1	1
<1	0	1	0	0	0	1
<1	1	0	0	0	0	0
<1	1	0	0	0	1	0
<1	1	0	1	0	0	1
<1	1	0	1	1	0	0

-----+-----
100%

Variables are (1) stufac (2) lenroll (3) gradrat (4) rmbrd (5) csat (6) act

27

EM in Stata

```
mi register impute gradrat lenroll rmbrd stufac csat
act private
mi impute mvn gradrat lenroll rmbrd stufac csat act
private, emonly
```

	gradrat	lenroll	rmbrd	stufac	csat	act	private
_cons	59.8618	6.169419	4.072555	14.86372	957.8762	22.2198	.6390169
Sigma							
gradrat	355.7137	-.4998451	10.38471	-31.14171	1352.981	30.58451	3.608253
lenroll	-.4998451	.9936801	-.0188409	1.382231	23.23804	.4695323	-.2964039
rmbrd	10.38471	-.0188409	1.32903	-1.685404	67.11875	1.514341	.1885311
stufac	-31.14171	1.382231	-1.685404	26.88555	-198.4039	-4.121786	-.9156043
csat	1352.981	23.23804	67.11875	-198.4039	14745.07	298.9068	9.381542
act	30.58451	.4695323	1.514341	-4.121786	298.9068	7.353064	.29118
private	3.608253	-.2964039	.1885311	-.9156043	9.381542	.29118	.2306743

These are the maximum likelihood estimates of the means and the covariance matrix.

28

Convert Covariances to Correlations

ML covariance matrix → ML correlation matrix

```
matrix Sigma=r(Sigma_em)
matrix M=r(Beta_em)
*we'll need these means later
_getcovcorr Sigma, corr
matrix C = r(C)
matlist C
```

	gradrat	lenroll	rmbrd	stufac	csat	act	private
gradrat	1						
lenroll	-.0265865	1					
rmbrd	.4776137	-.016395	1				
stufac	-.3184437	.2674224	-.2819532	1			
csat	.5907693	.1919786	.4794608	-.3151137	1		
act	.598022	.1737033	.4844202	-.2931513	.907775	1	
private	.3983337	-.6191004	.3404992	-.367662	.1608612	.2235773	1

29

EM As Input to regress

```
corr2data gradrat lenroll rmbrd stufac csat act
private, cov(Sigma) mean(M) clear
regress gradrat lenroll rmbrd stufac csat private
```

This produces ML estimates of the regression coefficients. But standard errors and associated statistics are incorrect because the sample size is taken to be 1302.

gradrat	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
lenroll	2.083176	.5393847	3.86	0.000	1.025013 3.141339
rmbrd	2.403941	.4000983	6.01	0.000	1.61903 3.188852
stufac	-.1813901	.0841226	-2.16	0.031	-.3464216 -.0163587
csat	.066875	.0039007	17.14	0.000	.0592227 .0745273
private	12.91442	1.146564	11.26	0.000	10.66509 15.16374
_cons	-32.39475	4.354628	-7.44	0.000	-40.93764 -23.85186

These are ML estimates

These are biased estimates

30

Direct ML

Also known as “raw ML” or “full information ML” (FIML)

Directly maximize the likelihood for the specified model

Several structural equation modeling (SEM) packages can do this for a large class of linear models.

- Amos
(www-03.ibm.com/software/products/en/spss-amos)
- Mplus (www.statmodel.com)
- LISREL (www.ssicentral.com/lisrel)
- OpenMX (R package) (openmx.psyc.virginia.edu)
- EQS (www.mvsoft.com)
- PROC CALIS (support.sas.com)
- Stata **sem** (www.stata.com)
- lavaan (R package) (lavaan.ugent.be)

31

Direct ML

With no missing data, the multivariate normal likelihood is

$$L(\theta) = \prod_i f(\mathbf{y}_i | \boldsymbol{\mu}(\theta), \boldsymbol{\Sigma}(\theta))$$

where

$$f(\mathbf{y}) = \frac{\exp[-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1}(\mathbf{y} - \boldsymbol{\mu})]}{(2\pi)^{k/2} |\boldsymbol{\Sigma}|^{1/2}}$$

32