

# Missing Data & How to Deal: An overview of missing data

Melissa Humphries  
Population Research Center

# Goals

---

- ▶ Discuss ways to evaluate and understand missing data
- ▶ Discuss common missing data methods
- ▶ Know the advantages and disadvantages of common methods
- ▶ Review useful commands in Stata for missing data



# General Steps for Analysis with Missing Data

---

- ▶ 1. Identify patterns/reasons for missing and recode correctly
- ▶ 2. Understand distribution of missing data
- ▶ 3. Decide on best method of analysis



# Step One: Understand your data

---

- ▶ **Attrition due to social/natural processes**
  - ▶ Example: School graduation, dropout, death
- ▶ **Skip pattern in survey**
  - ▶ Example: Certain questions only asked to respondents who indicate they are married
- ▶ **Intentional missing as part of data collection process**
- ▶ **Random data collection issues**
- ▶ **Respondent refusal/Non-response**



# Find information from survey (codebook, questionnaire)

- ▶ Identify skip patterns and/or sampling strategy from documentation

86A. Have you ever worked for pay, not counting work around the house?



(CIRCLE ONE)

No ..... 1 (SKIP TO QUESTION 93 ON PAGE 59)

Yes, and I am currently employed ..... 2 (SKIP TO QUESTION 87)

Yes, but I am not currently employed ..... 3 (GO TO QUESTION 86B)

86B. When did you last work for pay, not counting work around the house?



(WRITE IN)

|\_|\_|  
Month

19 |\_|\_|  
Year

```
. tab F2S86BYR
```

YEAR, LAST TIME R WORKED	Freq.	Percent	Cum.
81	1	0.10	0.10
87	1	0.10	0.19
88	4	0.39	0.58
89	7	0.68	1.25
90	35	3.38	4.63
91	234	22.59	27.22
92	81	7.82	35.04
97: REFUSED	5	0.48	35.52
98: MISSING	19	1.83	37.36
99: LEGITIMATE SKIP/NOT IN WAVE	649	62.64	100.00
Total	1,036	100.00	

# Recode for analysis: mvdecode command

```
. tab worked,m
```

HAS R EVER WORKED FOR PAY OUTSIDE HOME	Freq.	Percent	Cum.
1	160	13.41	13.41
2	488	40.91	54.32
3	383	32.10	86.42
7	1	0.08	86.50
8	8	0.67	87.18
.	153	12.82	100.00
Total	1,193	100.00	

```
. mvdecode worked, mv(7 8)  
worked: 9 missing values generated
```

```
. tab worked,m
```

HAS R EVER WORKED FOR PAY OUTSIDE HOME	Freq.	Percent	Cum.
1	160	13.41	13.41
2	488	40.91	54.32
3	383	32.10	86.42
.	162	13.58	100.00
Total	1,193	100.00	

## Recode for analysis: mvdecode command

Note: Stata reads missing (.) as a value greater than any number.

```
. tab worked,m
```

HAS R EVER WORKED FOR PAY OUTSIDE HOME	Freq.	Percent	Cum.
1	160	13.41	13.41
2	488	40.91	54.32
3	383	32.10	86.42
7	1	0.08	86.50
8	8	0.67	87.18
.	153	12.82	100.00
Total	1,193	100.00	

```
. mvdecode worked, mv(7 8)  
worked: 9 missing values generated
```

```
. tab worked,m
```

HAS R EVER WORKED FOR PAY OUTSIDE HOME	Freq.	Percent	Cum.
1	160	13.41	13.41
2	488	40.91	54.32
3	383	32.10	86.42
.	162	13.58	100.00
Total	1,193	100.00	

```
. tab worked if worked>2,m
```

HAS R EVER WORKED FOR PAY OUTSIDE HOME	Freq.	Percent	Cum.
3	383	70.28	70.28
.	162	29.72	100.00
Total	545	100.00	

# Analyze missing data patterns: misstable command

```
. misstable sum
```

Obs<.

Variable	Obs=.	Obs>.	Obs<.	Unique values	Min	Max
SCH_ID	91		1,102	>500	1249	91991
BYS75	91		1,102	6	0	8
BYS81B	91		1,102	8	1	98
BY2XRIRR	91		1,102	>500	10.91	99.99
BY2XMIRR	91		1,102	>500	16.5	99.99
F1S74	92		1,101	4	1	8
F1S80AA	92		1,101	6	0	8
F1S82	92		1,101	7	1	98

```
. misstable pattern, freq
```

Missing-value patterns

(1 means complete)

Frequency	Pattern							
	1	2	3	4	5	6	7	8
1,028	1	1	1	1	1	1	1	1
74	1	1	1	1	1	0	0	0
73	0	0	0	0	0	1	1	1
18	0	0	0	0	0	0	0	0
1,193								

Variables are (1) BY2XMIRR (2) BY2XRIRR (3) BYS75 (4) BYS81B (5) SCH\_ID (6) F1S74 (7) F1S80AA (8) F1S82

## Step Two: Missing data Mechanism (or probability distribution of missingness)

---

- ▶ Consider the probability of *missingness*
  - ▶ Are certain groups more likely to have missing values?
    - ▶ Example: Respondents in service occupations less likely to report income
  - ▶ Are certain responses more likely to be missing?
    - ▶ Example: Respondents with high income less likely to report income
- ▶ Certain analysis methods assume a certain probability distribution



# Missing Data Mechanisms

---

- ▶ **Missing Completely at Random (MCAR)**
  - ▶ Missing value ( $y$ ) neither depends on  $x$  nor  $y$ 
    - ▶ Example: some survey questions asked of a simple random sample of original sample
- ▶ **Missing at Random (MAR)**
  - ▶ Missing value ( $y$ ) depends on  $x$ , but not  $y$ 
    - ▶ Example: Respondents in service occupations less likely to report income
- ▶ **Missing not at Random (NMAR)**
  - ▶ The probability of a missing value depends on the variable that is missing
    - ▶ Example: Respondents with high income less likely to report income



# Exploring missing data mechanisms

---

- ▶ Can't be 100% sure about probability of missing (since we don't actually know the missing values)
- ▶ Could test for MCAR (t-tests)—but not totally accurate
- ▶ Many missing data methods assume MCAR or MAR but our data often are MNAR
  - ▶ Some methods specifically for MNAR
    - ▶ Selection model (Heckman)
    - ▶ Pattern mixture models



# Good News!!

---

- ▶ Some MAR analysis methods using MNAR data are still pretty good.
  - ▶ May be another measured variable that indirectly can predict the probability of missingness
    - ▶ Example: those with higher incomes are less likely to report income BUT we have a variable for years of education and/or number of investments
  - ▶ ML and MI are often unbiased with NMAR data even though assume data is MAR
    - ▶ See Schafer & Graham 2002



# Step 3: Deal with missing data

---

- ▶ Use what you know about
  - ▶ Why data is missing
  - ▶ Distribution of missing data
- ▶ Decide on the best analysis strategy to yield the least biased estimates
  - ▶ Deletion Methods
    - ▶ Listwise deletion, pairwise deletion
  - ▶ Single Imputation Methods
    - ▶ Mean/mode substitution, dummy variable method, single regression
  - ▶ Model-Based Methods
    - ▶ Maximum Likelihood, Multiple imputation



# Deletion Methods

---

- ▶ Listwise deletion
  - ▶ AKA complete case analysis
- ▶ Pairwise deletion



# Listwise Deletion (Complete Case Analysis)

- ▶ Only analyze cases with available data on each variable
  - ▶ Advantages:
    - ▶ Simplicity
    - ▶ Comparability across analyses
  - ▶ Disadvantages:
    - ▶ Reduces statistical power (because lowers n)
    - ▶ Doesn't use all information
    - ▶ Estimates may be biased if data not MCAR\*

Gender	8 <sup>th</sup> grade math test score	12 <sup>th</sup> grade math score
F	45	.
M	.	99
F	55	86
F	85	88
F	80	75
.	81	82
F	75	80
M	95	.
M	86	90
F	70	75
F	85	.

\*NOTE: List-wise deletion often produces *unbiased regression slope estimates* as long as missingness is not a function of outcome variable.



# Pairwise deletion (Available Case Analysis)

- ▶ Analysis with all cases in which the variables of interest are present.
  - ▶ Advantage:
    - ▶ Keeps as many cases as possible for each analysis
    - ▶ Uses all information possible with each analysis
  - ▶ Disadvantage:
    - ▶ Can't compare analyses because sample different each time

Gender	8 <sup>th</sup> grade math test score	12 <sup>th</sup> grade math score
F	45	.
M	.	99
F	55	86
F	85	88
F	80	75
.	81	82
F	75	80
M	95	.
M	86	90
F	70	75
F	85	.

# Single imputation methods

---

- ▶ Mean/Mode substitution
- ▶ Dummy variable control
- ▶ Conditional mean substitution

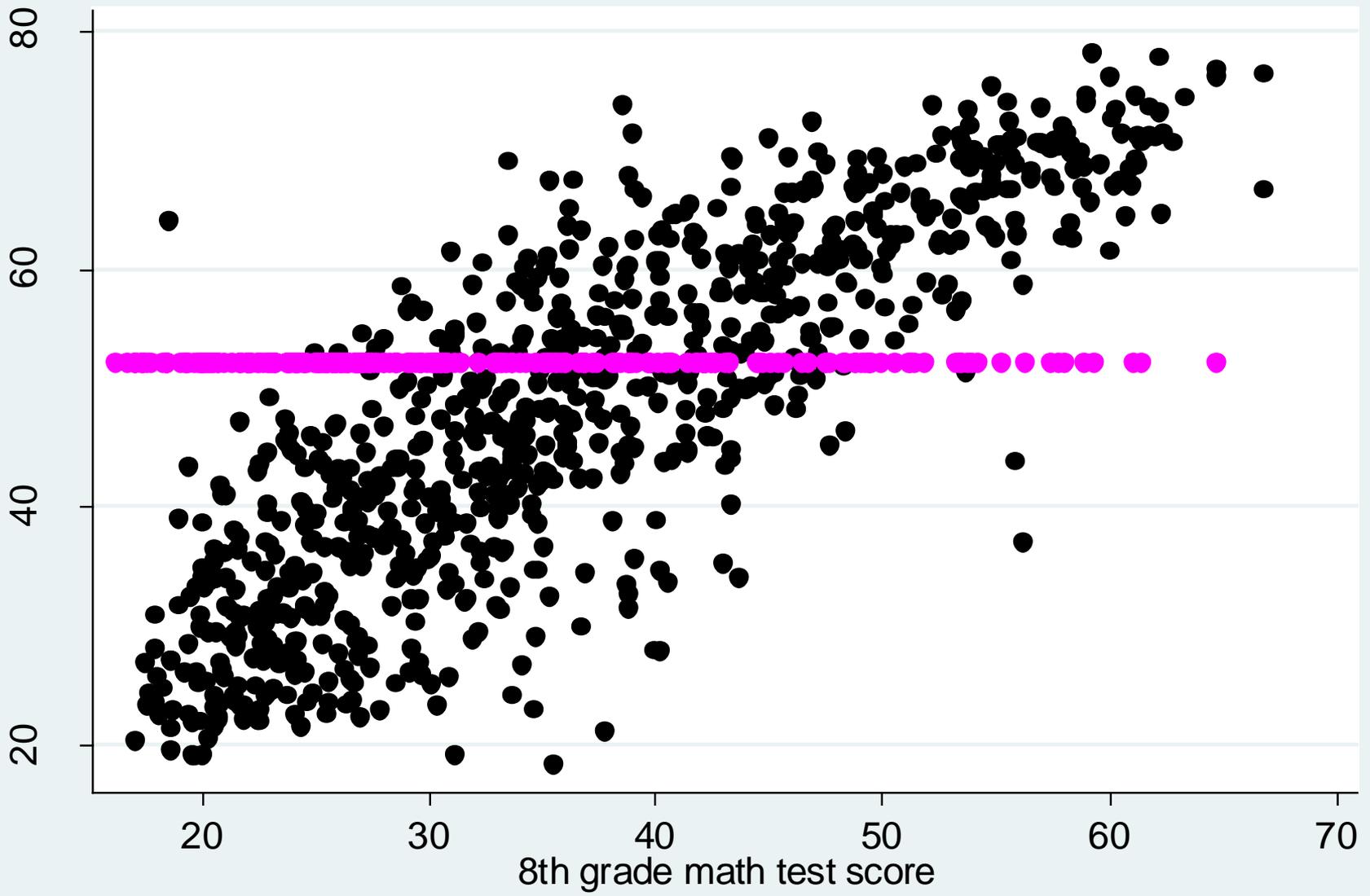


# Mean/Mode Substitution

---

- ▶ Replace missing value with sample mean or mode
- ▶ Run analyses as if all complete cases
- ▶ Advantages:
  - ▶ Can use complete case analysis methods
- ▶ Disadvantages:
  - ▶ Reduces variability
  - ▶ Weakens covariance and correlation estimates in the data (because ignores relationship between variables)





● imputed 12th grade math test score (mean sub)

# Dummy variable adjustment

---

- ▶ Create an indicator for missing value (1=value is missing for observation; 0=value is observed for observation)
- ▶ Impute missing values to a constant (such as the mean)
- ▶ Include missing indicator in regression
- ▶ Advantage:
  - ▶ Uses all available information about missing observation
- ▶ Disadvantage:
  - ▶ Results in biased estimates
  - ▶ Not theoretically driven
- ▶ **NOTE:** Results not biased if value is missing because of a legitimate skip

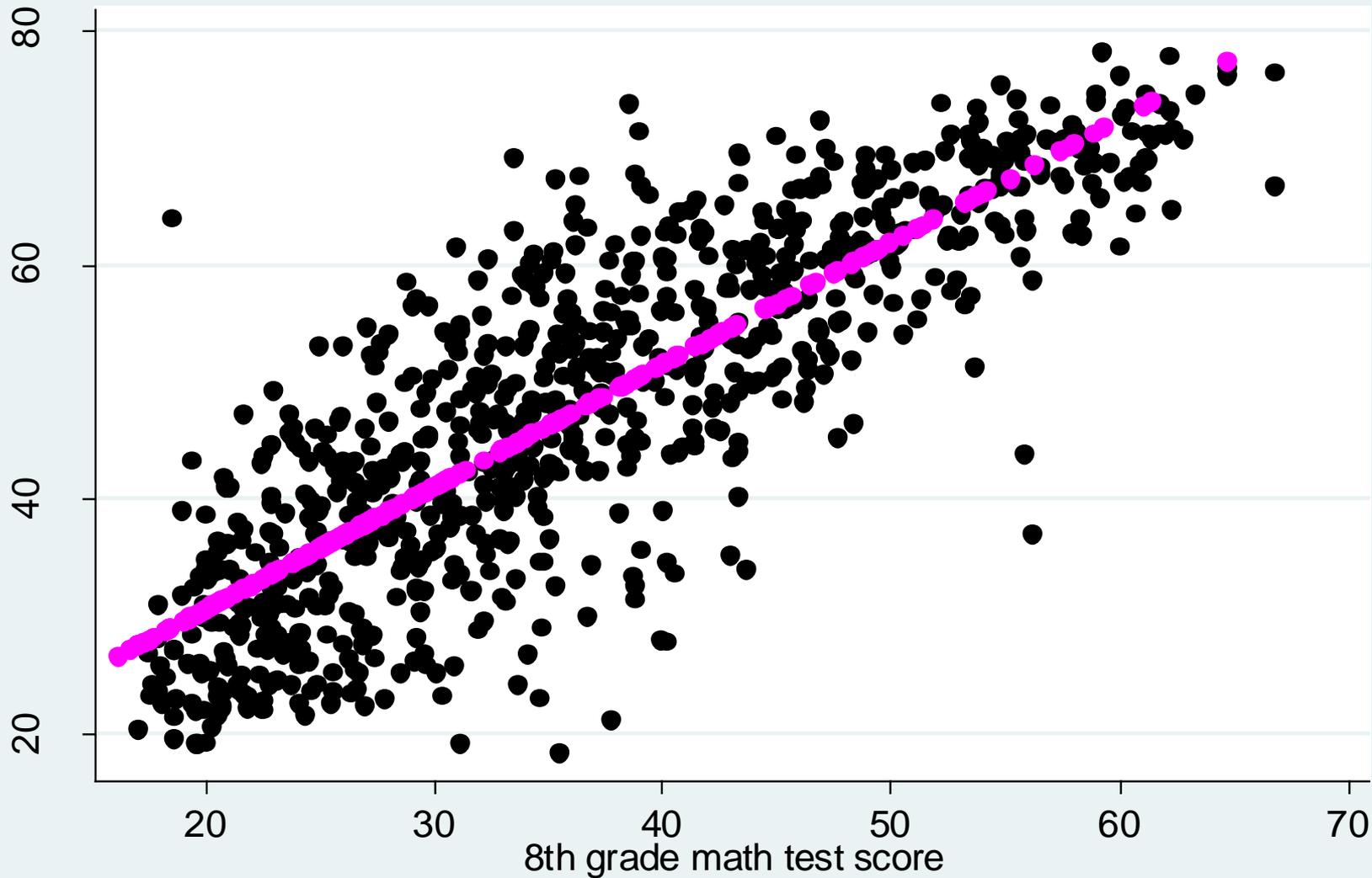


# Regression Imputation

---

- ▶ Replaces missing values with predicted score from a regression equation.
  - ▶ Advantage:
    - ▶ Uses information from observed data
  - ▶ Disadvantages:
    - ▶ Overestimates model fit and correlation estimates
    - ▶ Weakens variance





# Model-based methods

---

- ▶ Maximum Likelihood
- ▶ Multiple imputation



# Model-based Methods: Maximum Likelihood Estimation

---

- ▶ Identifies the set of parameter values that produces the highest log-likelihood.
  - ▶ ML estimate: value that is most likely to have resulted in the observed data
- ▶ Conceptually, process the same with or without missing data
  - ▶ Advantages:
    - ▶ Uses full information (both complete cases and incomplete cases) to calculate log likelihood
    - ▶ Unbiased parameter estimates with MCAR/MAR data
  - ▶ Disadvantages
    - ▶ SEs biased downward—can be adjusted by using observed information matrix



# Multiple Imputation

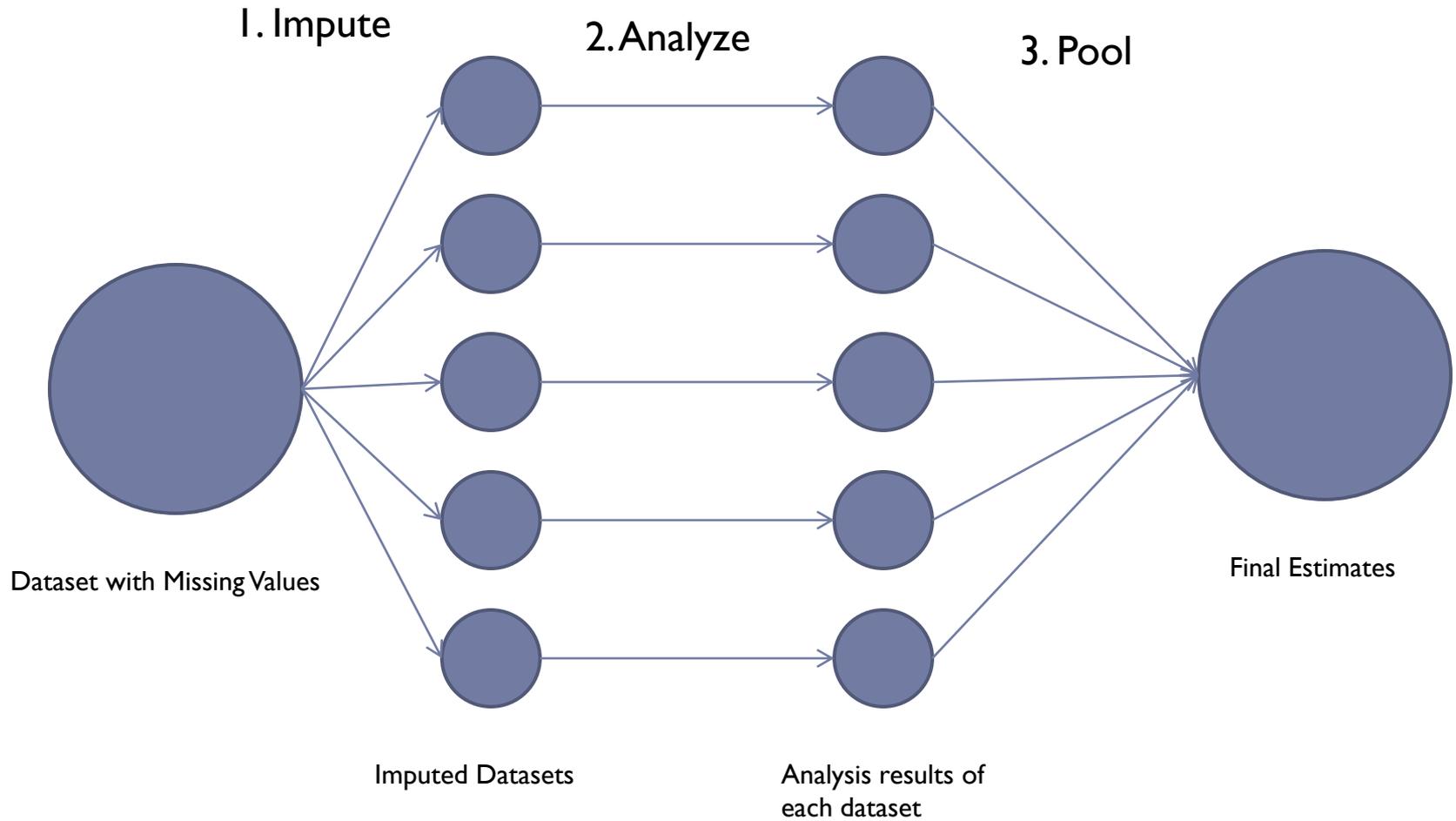
---

- ▶ 1. Impute: Data is 'filled in' with imputed values using specified regression model
  - ▶ This step is repeated  $m$  times, resulting in a separate dataset each time.
- ▶ 2. Analyze: Analyses performed within each dataset
- ▶ 3. Pool: Results pooled into one estimate
  - ▶ Advantages:
    - ▶ Variability more accurate with multiple imputations for each missing value
      - Considers variability due to sampling AND variability due to imputation
  - ▶ Disadvantages:
    - ▶ Cumbersome coding
    - ▶ Room for error when specifying models



# Multiple Imputation Process

---



# Multiple Imputation: Stata & SAS

---

- ▶ **SAS:**

- ▶ Proc mi

- ▶ **Stata:**

- ▶ ice (imputation using chained equations) & mim (analysis with multiply imputed dataset)
  - ▶ mi commands
    - ▶ mi set
    - ▶ mi register
    - ▶ mi impute
    - ▶ mi estimate
  - ▶ **NOTE:** the ice command is the only chained equation method until Stata 12. Chained equations can be used as an option of mi impute since Stata 12.



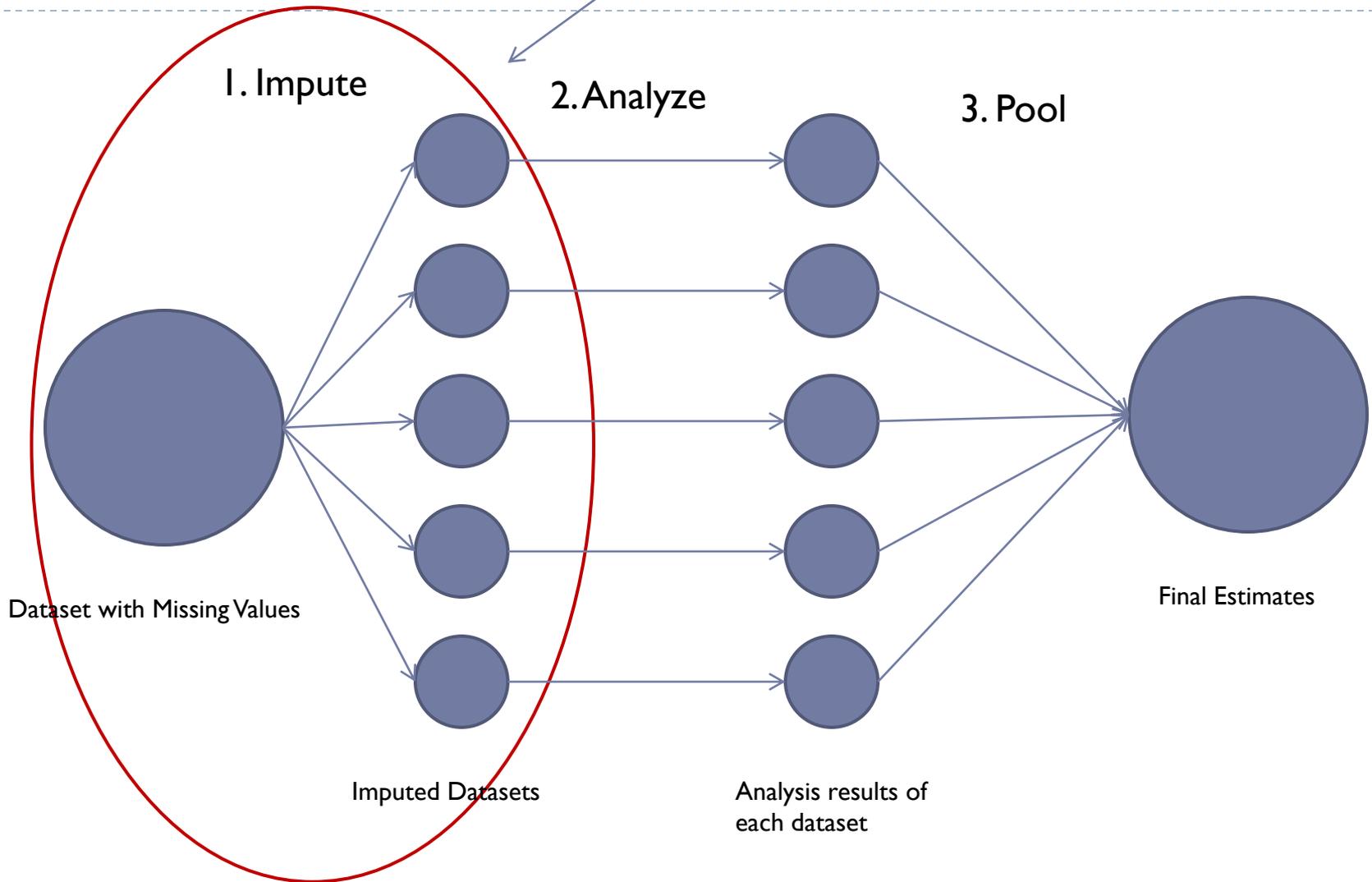
# ice & mim

---

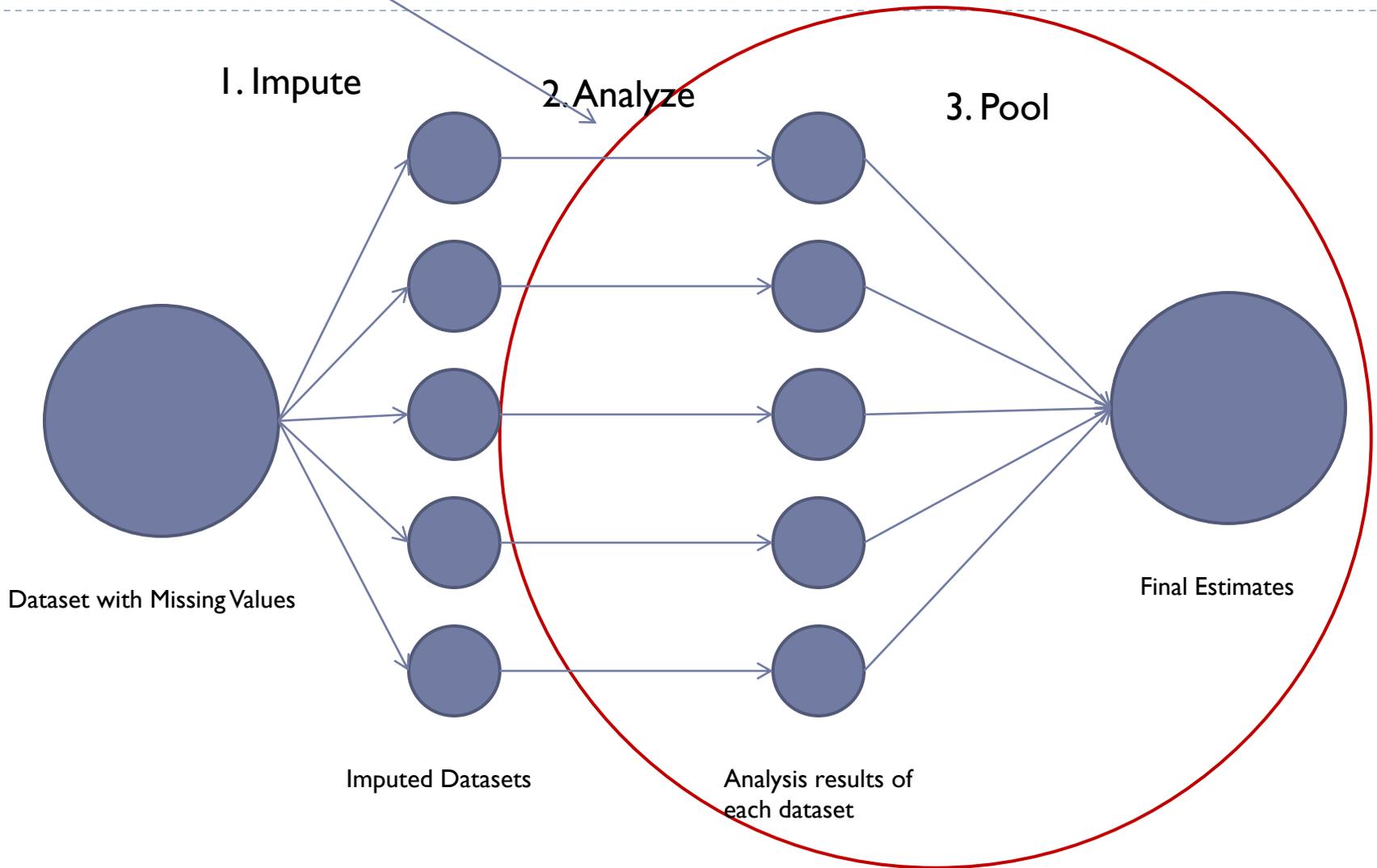
- ▶ **ice: Imputation using chained equations**
  - ▶ Series of equations predicting one variable at a time
  - ▶ Creates as many datasets as desired
- ▶ **mim: prefix used before analysis that performs analyses across datasets and pools estimates**



# ice command



# mim command



```
ice female lm latino black asian other F1PARED AGE1 intact bymirt ESL2 ALG2OH acgpa ac_engall
hardwtr Lksch MAE10 RAE10 hilep midw south public catholic colltype aceng_ESL Lksch_ESL, ///
saving(imputed2) m(5) cmd (Lksch:ologit)
```

Variable	Command	Prediction equation
female		[No missing data in estimation sample]
lm		[No missing data in estimation sample]
latino		[No missing data in estimation sample]
black		[No missing data in estimation sample]
ALG2OH	logit	female lm latino black asian other F1PARED AGE1 intact bymirt ESL2 acgpa ac_engall hardwtr Lksch MAE10 RAE10 hilep midw south public catholic colltype aceng_ESL Lksch_ESL
acgpa	regress	female lm latino black asian other F1PARED AGE1 intact bymirt ESL2 ALG2OH ac_engall hardwtr Lksch MAE10 RAE10 hilep midw south public catholic colltype aceng_ESL Lksch_ESL
ac_engall	regress	female lm latino black asian other F1PARED AGE1 intact bymirt ESL2 ALG2OH acgpa hardwtr Lksch MAE10 RAE10 hilep midw south public catholic colltype Lksch_ESL
hardwtr	logit	female lm latino black asian other F1PARED AGE1 intact bymirt ESL2 ALG2OH acgpa ac_engall Lksch MAE10 RAE10 hilep midw south public catholic colltype aceng_ESL Lksch_ESL
Lksch	ologit	female lm latino black asian other F1PARED AGE1 intact bymirt ESL2 ALG2OH acgpa ac_engall hardwtr MAE10 RAE10 hilep midw south public catholic colltype aceng_ESL
MAE10	regress	female lm latino black asian other F1PARED AGE1 intact bymirt ESL2 ALG2OH acgpa ac_engall hardwtr Lksch RAE10 hilep midw south public catholic colltype aceng_ESL Lksch_ESL
RAE10	regress	female lm latino black asian other F1PARED AGE1 intact bymirt ESL2 ALG2OH acgpa ac_engall hardwtr Lksch MAE10 hilep midw south public catholic colltype aceng_ESL Lksch_ESL
hilep	logit	female lm latino black asian other F1PARED AGE1 intact bymirt ESL2 ALG2OH acgpa ac_engall hardwtr Lksch MAE10 RAE10 midw south public catholic colltype aceng_ESL Lksch_ESL

```
Imputing .....1.....2.....3.....4.....5
file imputed2.dta saved
```



# mi commands

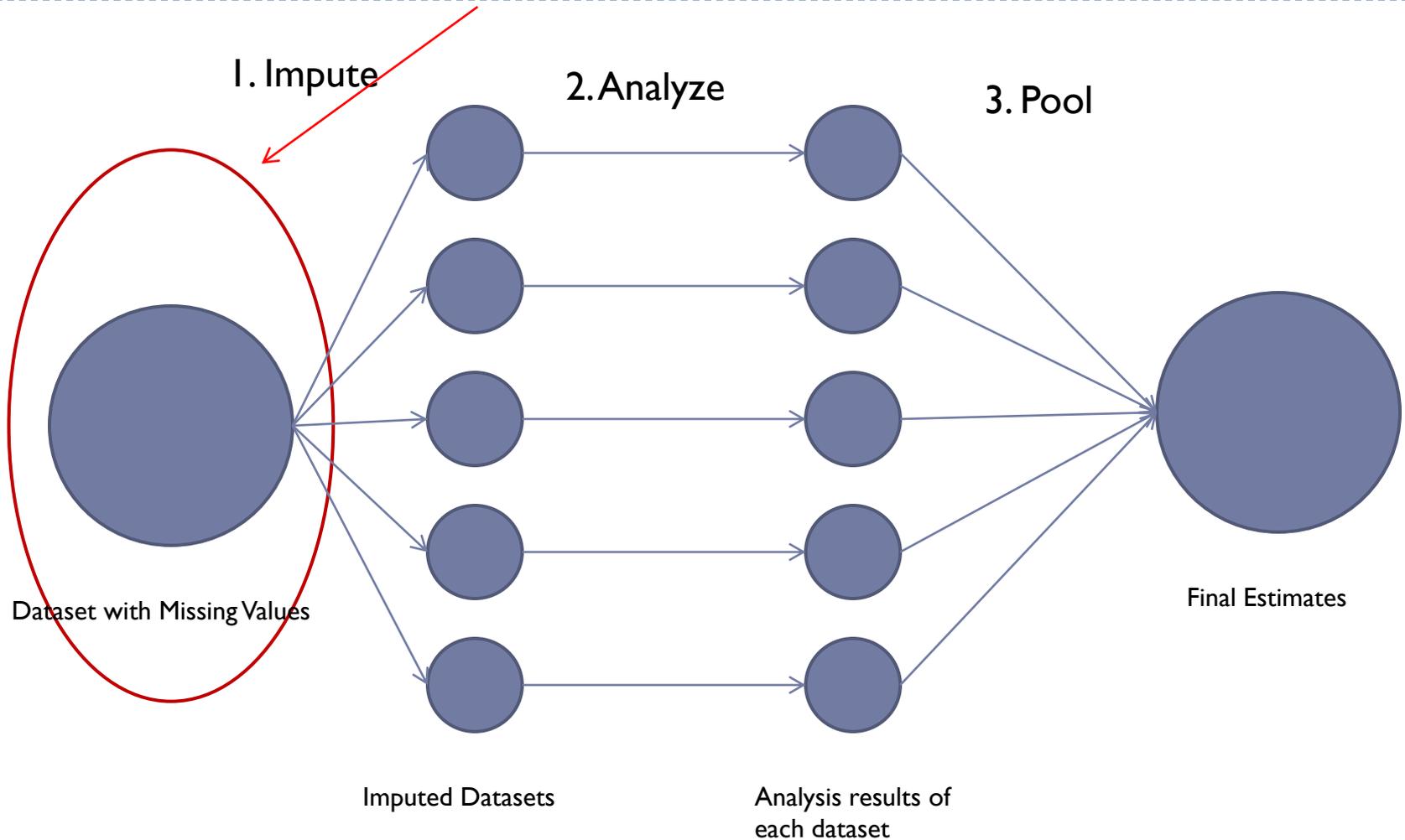
---

- ▶ Included in Stata 11
- ▶ Includes univariate multiple imputation (impute only one variable)
- ▶ Multivariate imputation probably more useful for our data
- ▶ Specific order:
  - ▶ mi set
  - ▶ mi register
  - ▶ mi impute
  - ▶ mi estimate

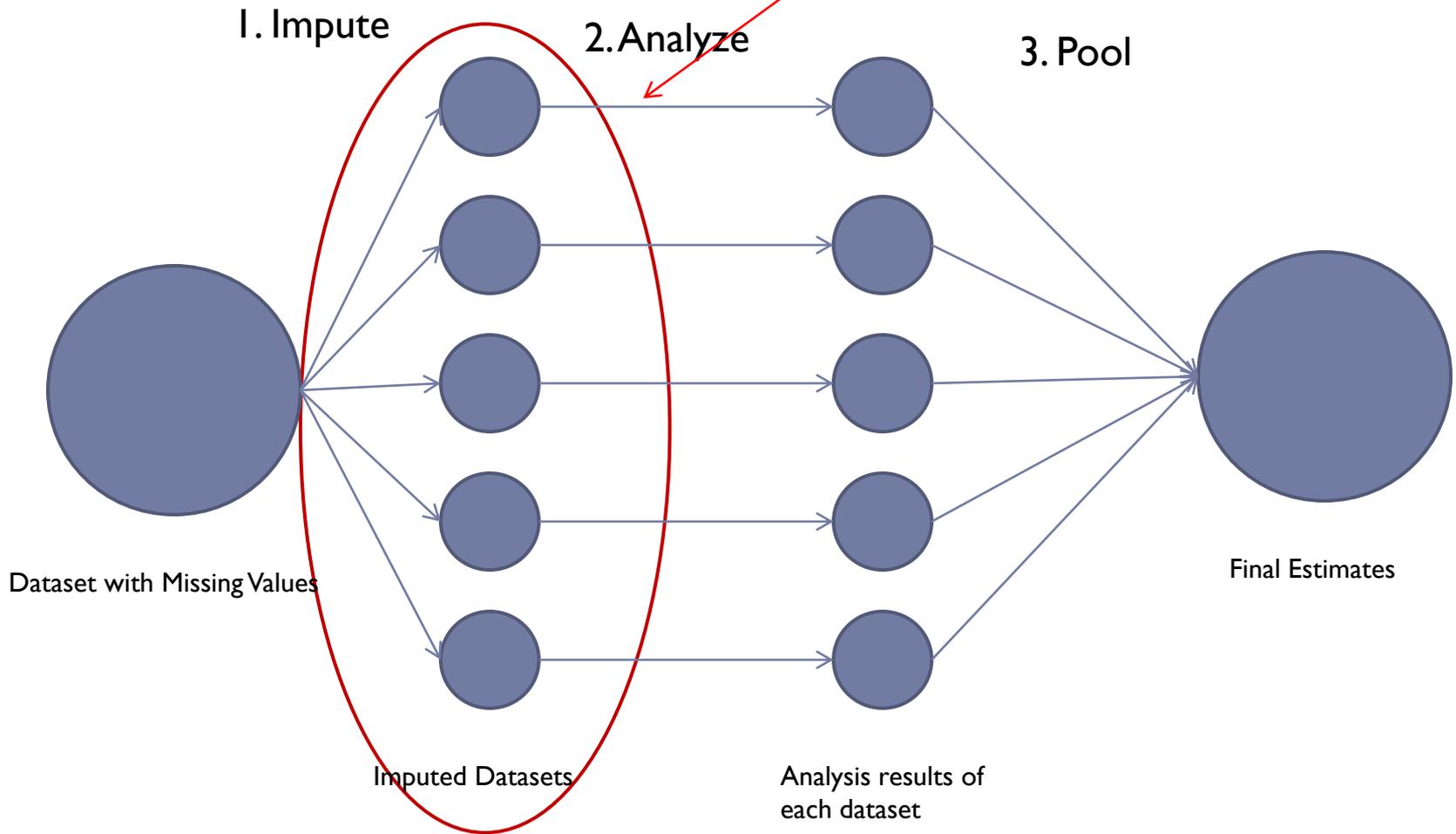


# mi set and mi register commands

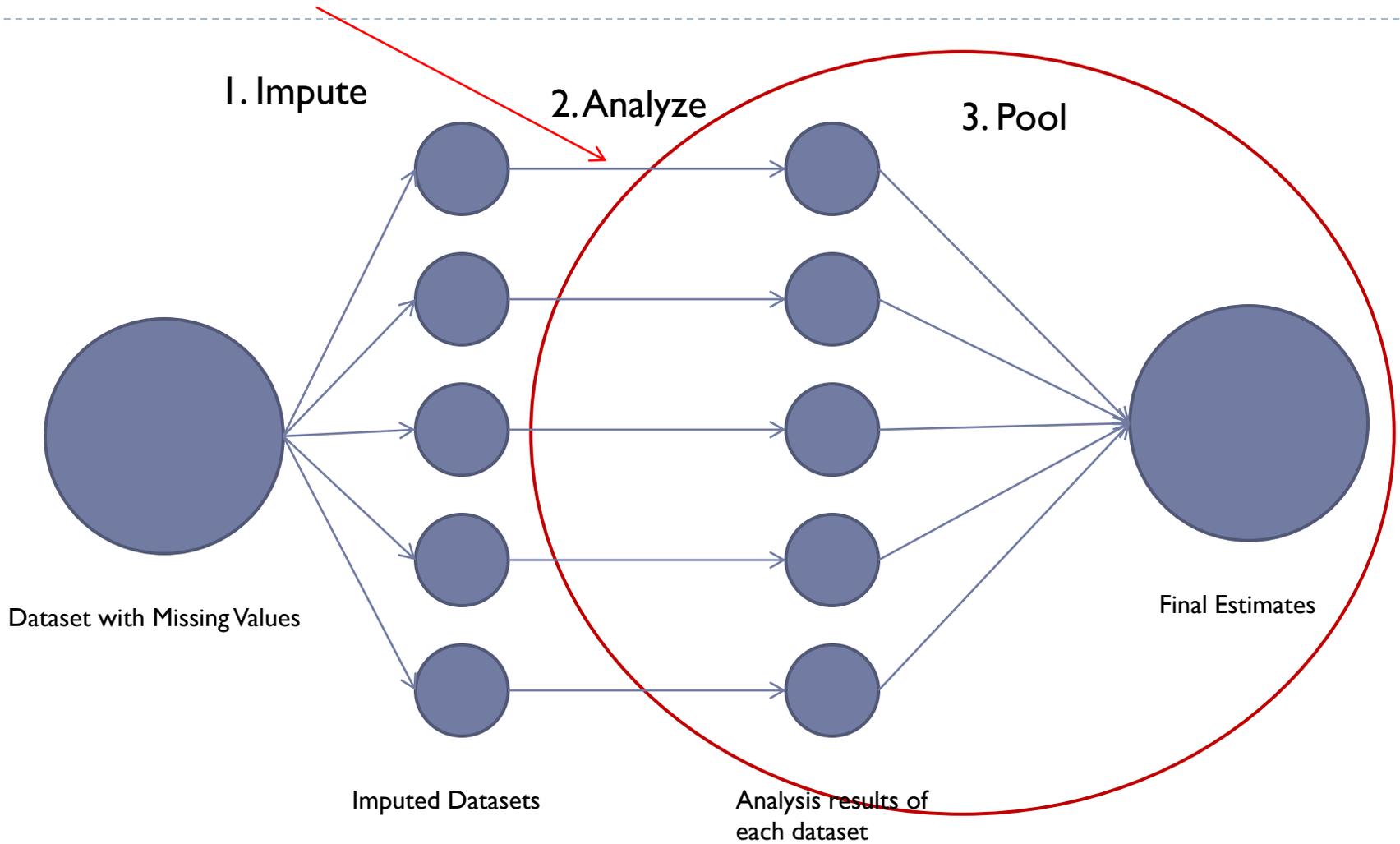
---



# mi impute command



# mi estimate command



```
*****set data to be multiply imputed (can set to 'wide' format also)
mi set flong
```

---

```
*****register variables as "imputed" (variables with missing data that you want imputed)
or "regular"
mi register imputed readtest8 worked mathtest8
mi register regular sex race
```

```
*****describing data
mi describe
```

```
*****setting seed so results are replicable
set seed 8945
```

```
*****imputing using chained equations—using ols regression for predicting read and math
test using mlogit to predict worked
mi impute chained (regress) readtest8 mathtest8 (mlogit) worked=sex i.race, add(10)
```

```
*****check new imputed dataset
mi describe
```

```
*****estimating model using imputed values
mi estimate:regress mathtest12 mathtest8 sex race
```

---



```

. *set data to be multiply imputed
. mi set flong

. mi svyset
no survey characteristics are set

. *register variables as "imputed" (variables with missing data that you want imputed)
. mi register imputed readtest8 worked mathtest8
(165 m=0 obs. now marked as incomplete)

. mi register regular sex race

.
. mi describe

Style:  flong
      last mi update 29mar2012 10:36:55, 0 seconds ago

Obs.:  complete           955
      incomplete         165  (M = 0 imputations)
      -----
      total              1,120

Vars.:  imputed:  3; readtest8(40) worked(132) mathtest8(41)

      passive:  0

      regular:  2; sex race

      system:   3; _mi_m _mi_id _mi_miss

      (there are 20 unregistered variables)

```

```
. *****setting seed so results are replicable
. set seed 8945

. *****imputing using chained equations
. mi impute chained (regress) readtest8 mathtest8 (mlogit) worked=sex i.race, add(10)
```

Conditional models:

```
readtest8: regress readtest8 mathtest8 i.worked sex i.race
mathtest8: regress mathtest8 readtest8 i.worked sex i.race
worked: mlogit worked readtest8 mathtest8 sex i.race
```

Performing chained iterations ...

```
Multivariate imputation          Imputations =          10
Chained equations                added =          10
Imputed: m=1 through m=10       updated =           0
```

```
Initialization: monotone        Iterations =         100
                                burn-in =          10
```

```
readtest8: linear regression
mathtest8: linear regression
worked: multinomial logistic regression
```

Variable	Observations per m			Total
	Complete	Incomplete	Imputed	
readtest8	1060	37	37	1097
mathtest8	1060	37	37	1097
worked	958	139	139	1097

(complete + incomplete = total; imputed is the minimum across m of the number of filled-in observations.)

# Dataset after imputation

---

```
. sum readtest8
```

Variable	Obs	Mean	Std. Dev.	Min	Max
readtest8	12030	27.02419	8.862721	-4.311928	52.88457



```
. *****estimating model using imputed values
. mi estimate:regress mathtest12 mathtest8 sex race
```

```
Multiple-imputation estimates      Imputations      =      10
Linear regression                  Number of obs    =      830
                                   Average RVI      =      0.1070
                                   Largest FMI      =      0.2241
                                   Complete DF     =      826
DF adjustment:  Small sample      DF:      min    =      148.93
                                   avg            =      410.45
                                   max            =      704.25
Model F test:      Equal FMI      F(   3,  515.9) =      525.22
Within VCE type:   OLS            Prob > F        =      0.0000
```

mathtest12	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
mathtest8	1.022907	.0276674	36.97	0.000	.968236	1.077579
sex	-2.05258	.5941458	-3.45	0.001	-3.219089	-.8860704
race	-.015127	.3232535	-0.05	0.963	-.6501539	.6198999
_cons	13.78052	1.717096	8.03	0.000	10.39951	17.16154



# Notes and help with mi in stata

---

- ▶ **LOTS of options**
  - ▶ Can specify exactly how you want imputed
  - ▶ Can specify the model appropriately (ex. Using svy command)
  - ▶ `mi impute mvn` (multivariate normal regression) also useful
- ▶ **Help mi is useful**
- ▶ **Also, UCLA has great website about ice and mi**



# General Tips

---

- ▶ Try a few methods: often if result in similar estimates, can put as a footnote to support method
- ▶ Some don't impute dependent variable
  - ▶ But would still use to impute independent variables



# References

---

- ▶ Allison, Paul D. 2001. Missing Data. *Sage University Papers Series on Quantitative Applications in the Social Sciences*. Thousand Oaks: Sage.
- ▶ Enders, Craig. 2010. Applied Missing Data Analysis. Guilford Press: New York.
- ▶ Little, Roderick J., Donald Rubin. 2002. Statistical Analysis with Missing Data. John Wiley & Sons, Inc: Hoboken.
- ▶ Schafer, Joseph L., John W. Graham. 2002. “Missing Data: Our View of the State of the Art.” *Psychological Methods*.

