

TEXAS TECH UNIVERSITY"



Optimizing Inferences using Principled Missing Data Treatments

Todd D. Little

Director, Institute for Measurement, Methodology, Analysis and Policy Director & Founder, Stats Camp (Stats Camp.org)

NIH Webinar June 29th, 2016

Planning for Missing Data

- STATS CAMP
- Learn about the different types of missing data
- Learn about ways in which the missing data process can be recovered
- Understanding why imputing missing data is not cheating
 - Learn why NOT imputing missing data is more likely to lead to errors in generalization!
- Learn about intentionally missing designs

Key Considerations



• **Recoverability**

- Is it possible to recover what the *sufficient statistics* would have been if there was no missing data?

• (sufficient statistics = means, variances, and covariance)

- Is it possible to recover what the *parameter estimates* of a model would have been if there was no missing data.

Bias

- Are the sufficient statistics/parameter estimates systematically different than what they would have been had there not been any missing data?

• Power

- Do we have the same or similar rates of power (1-Type II error rate) as we would without missing data?

Types of Missing Data



• Missing Completely at Random (MCAR)

- No association with unobserved variables (selective process) and no association with observed variables

- Missing at Random (MAR)
 - No association with unobserved variables, but maybe related to observed variables
 - Random in the statistical sense of predictable
- Non-Random (Selective) Missing (MNAR)

- Some association with unobserved variables and maybe with observed variables

Effects of imputing missing data



immap.educ.ttu.edu

STATS CAMP

statscamp.ord

Missing Data Mechanisms

MCAR: $P(R | Y_{mis}, Y_{obs}, \theta) = P(R | \theta)$

MAR: $P(R | Y_{mis}, Y_{obs}, \theta) = P(R | Y_{obs}, \theta)$

MNAR: $P(R|Y_{mis}, Y_{obs}, \theta) \neq P(R|Y_{obs}, \theta)$ R = Missingness pattern matrix

 Y_{mis} = Missing component of the incomplete data

 Y_{obs} = Observed component of the incomplete data

 θ = Vector of parameters governing the relation between *R* and the incomplete data



Effects of imputing missing data

No Association with Unobserved /Unmeasured Variable(s)

An Association with Unobserved /Unmeasured Variable(s)

No Association with ANY Observed Variable	An Association with <u>Analyzed</u> Variables	An Association with <u>Unanalyzed</u> Variables			
MCAR •Fully recoverable •Fully unbiased	MAR • Partly to fully recoverable • Less biased to unbiased	MAR Partly to fully recoverable Less biased to unbiased 			
NMAR • Unrecoverable • Biased (same bias as not estimating)	 MAR/NMAR Partly to fully recoverable Same to unbiased 	MAR/NMAR Partly to fully recoverable Same to unbiased 			

Modern Missing Data Analysis



MI or FIML

In 1978, Rubin proposed Multiple Imputation (MI)

- An approach especially well suited for use with large public-use databases.

- First suggested in 1978 and developed more fully in 1987.
- MI primarily uses the Expectation Maximization (EM) algorithm and/or the Markov Chain Monte Carlo (MCMC) algorithm.

Beginning in the 1980's, likelihood approaches developed.

- Multiple group SEM
- Full Information Maximum Likelihood (FIML).
 - An approach well suited to more circumscribed models

Multiple Imputation Generalizability



Impute if the population to which you want to generalize is the original sample

Aging studies want to generalize to successful agers.

• Drop outs who've died weren't successful....

Skip and Fill patterns that can't happen.

• Males who are pregnant....

Need for Auxiliary Variables



60% MAR correlation estimates with no auxiliary variables



Simulation results showing XY correlation estimates (with 95 and 99% confidence intervals) associated with a 60% MAR Situation.

Need for Auxiliary Variables







Simulation results showing XY correlation estimates (with 95 and 99% confidence intervals) associated with a 60% MAR Situation and 8 auxiliary variables.

Need for Auxiliary Variables



Simulation results showing XY correlation estimates (with 95 and 99% confidence intervals) associated with a 60% MAR Situation and 1 PCA auxiliary variable.

Auxiliary Variable Power Comparison



ATS CAMP

Faster and more reliable convergence





Macros to create PCA Auxiliaries



2) Capture non-linear information

a) Step 1:

Use the SAS MACRO below to obtain squares and interaction terms for any set of variables. First, run the following code to load the SAS MACRO. (*Do not edit %MACRO code below*).

```
/*Do not edit this syntax*/
```

b) Step 2:

Next, enter your <u>data set name</u> and your <u>list of variables</u>. This step calls the SAS MACRO and uses it with your data.

your analysis and auxiliary variable names, (3) set quadr option: 1 = all possible quadratic terms and 2-way interactions OR 2 = all possible cubic terms and 3-way interactions */ RUN;

c) Note:

Again, check your data to ensure that it is correct (see illustration below). Consider the variables INT_a1_a1 and INT_a1_a2 in this example. INT_a1_a1 is the a1 variables squared (i.e., $2^2 = 4$) and INT_a1_a2 is the interaction term for a1 and a2 (i.e., 2 * 3 = 6). Also, note that missing values are generated when at least one term is missing (see INT_a1_a3 below).

Nonlinear Terms Data (data=myinterdata)

	aux1	aux2	aux3	INT_a1_a1	INT_a1_a2	INT_a1_a3	INT_a1_a4	INT_a2_a2
1	4	5	4	4	6		6	9
2	3	3	2	16	4	12	4	1
3	4	3	4	25	15	5	20	9
4		4	4	25	15	5	20	9
5	5	3	3	25	10		5	4
6	5	2	4	25	5	5	10	1

R packages simsem, semTools and the routine Quark.

simsem: http://simsem.org/

semTools: https://github.com/simsem/semTools/wiki

Planned missing data designs

STATS CAMP

In planned missing data designs, participants are *randomly assigned* to conditions in which they do not respond to all items, all measures, and/or all measurement occasions

Why would you want to do this?

- 1. Long assessments can reduce data quality
- 2. Repeated assessments can induce practice effects
- 3. Collecting data can be time- and cost-intensive
- 4. Less taxing assessments may reduce unplanned missingness

Planned Missing Design Examples



• Multi-form designs, specifically the 3-form Planned Missing Design

- 2-method Planned Missing Design
 - Guaranteed to blow your mind



Form	Common Set X	Variable Set A	Variable Set B	Variable Set C
1	¹ ⁄4 of items	¹ ⁄4 of items	¹ ⁄4 of items	missing
2	¹ ⁄4 of items	¹ ⁄4 of items	missing	¹ / ₄ of items
3	¹ ⁄4 of items	missing	¹ ⁄4 of items	¹ / ₄ of items
Form	Common Set X	Variable Set A	Variable Set B	Variable Set C
Form 1	Common Set X A few items	Variable Set A 1/3 of items	Variable Set B 1/3 of items	Variable Set C missing
Form 1 2	Common Set X A few items A few items	Variable Set A 1/3 of items 1/3 of items	Variable Set B 1/3 of items missing	Variable Set C missing 1/3 of items
Form	Common Set X	Variable Set A	Variable Set B	Variable Set C



How it works

- ALL participants receive Measure 2 (the cheap one)
- A subset of participants *also* receive Measure 1 (the gold standard)
- Using both measures (on a subset of participants) enables us to estimate and remove the bias from the inexpensive measure (for all participants) using a latent variable model



21 questions made up of 7 3-question subtests

Subtest	Item	Subtest	Item
Demographics	How old are you? Are you male or female? What is your occupation?	Extraversion	I start conversations. I am the life of the party. I am comfortable around people.
Musical Taste	What is your favorite genre of music?Do you like to listen to music while you work?Do you prefer music played	Neuroticism	I get stressed out easily. I get irritated easily. I have frequent mood swings.
	loud or softly?	Conscientiousness	I am always prepared.
Openness	I have a rich vocabulary. I have excellent ideas.		I like order. I pay attention to details.
	I have a vivid imagination.	Agreeableness	I am interested in people. I have a soft heart. I take time out for others.

STATS CAMP

• Common Set (X)

Subtest	Item	Subtest	Item		
Demographics	Demographics How old are you? Are you male or female? What is your occupation?		I start conversations. I am the life of the party. I am comfortable around		
Musical Taste	What is your favorite genre of music?Do you like to listen to music while you work?Do you prefer music played	Neuroticism	I get stressed out easily. I get irritated easily. I have frequent mood swings.		
	loud or softly?	Conscientiousness	I am always prepared.		
Openness	I have a rich vocabulary. I have excellent ideas.		I like order. I pay attention to details.		
	I have a vivid imagination.	Agreeableness	I am interested in people. I have a soft heart. I take time out for others.		

STATS CAMP

• Common Set (X)

Subtest	Item	Subtest	Item		
Demographics	How old are you? Are you male or female? What is your occupation?	Extraversion	I start conversations. I am the life of the party. I am comfortable around people.		
Musical Taste	What is your favorite genre of music?Do you like to listen to music while you work?Do you prefer music played	Neuroticism	I get stressed out easily. I get irritated easily. I have frequent mood swings.		
	loud or softly?	Conscientiousness	I am always prepared.		
Openness	I have a rich vocabulary. I have excellent ideas.		I like order. I pay attention to details.		
	I have a vivid imagination.	Agreeableness	I am interested in people. I have a soft heart. I take time out for others.		



• Set A

Subtest	Item	Subtest	Item		
Demographics	How old are you? Are you male or female? What is your occupation?	Extraversion	I start conversations. I am the life of the party. I am comfortable around people		
Musical Taste	What is your favorite genre of music?Do you like to listen to music while you work?Do you prefer music played	Neuroticism	I get stressed out easily. I get irritated easily. I have frequent mood swings.		
Openness	loud or softly? I have a rich vocabulary.	Conscientiousness	I am always prepared. I like order. I pay attention to details.		
	I have a vivid imagination.	Agreeableness	I am interested in people. I have a soft heart. I take time out for others.		



• Set B

Subtest	Item	Subtest	Item		
Demographics	How old are you? Are you male or female? What is your occupation?	Extraversion	I start conversations. I am the life of the party. I am comfortable around people		
Musical Taste	What is your favorite genre of music?Do you like to listen to music while you work?Do you prefer music played	Neuroticism	I get stressed out easily. I get irritated easily. I have frequent mood swings.		
	loud or softly?	Conscientiousness	I am always prepared.		
Openness	I have a rich vocabulary. I have excellent ideas.		I like order. I pay attention to details.		
	I have a vivid imagination.	Agreeableness	I am interested in people.		
			I have a soft heart. I take time out for others		



• Set C

Subtest	Item	Subtest	Item
Demographics	How old are you? Are you male or female? What is your occupation?	Extraversion	I start conversations. I am the life of the party. I am comfortable around people
Musical Taste	What is your favorite genre of music?Do you like to listen to music while you work?Do you prefer music played	Neuroticism	I get stressed out easily. I get irritated easily. I have frequent mood swings.
Openness	I have a rich vocabulary. I have excellent ideas.	Conscientiousness	I am always prepared. I like order. I pay attention to details.
	I have a vivid imagination.	Agreeableness	I am interested in people. I have a soft heart.

I take time out for others.

Form 1 (XAB)	Form 2 (XAC)	Form 3 (XBC)
How old are you?	How old are you?	How old are you?
Are you male or female?	Are you male or female?	Are you male or female?
What is your occupation?	What is your occupation?	What is your occupation?
What is your favorite genre of music?Do you like to listen to music while you work?Do you prefer music played loud or softly?	What is your favorite genre of music?Do you like to listen to music while you work?Do you prefer music played loud or softly?	What is your favorite genre of music?Do you like to listen to music while you work?Do you prefer music played loud or softly?
I have a rich vocabulary.	I have a rich vocabulary.	I have excellent ideas.
I have excellent ideas.	I have a vivid imagination.	I have a vivid imagination.
I start conversations.	I start conversations.	I am the life of the party.
I am the life of the party.	I am comfortable around people.	I am comfortable around people.
I get stressed out easily.	I get stressed out easily.	I get irritated easily.
I get irritated easily.	I have frequent mood swings.	I have frequent mood swings.
I am always prepared.	I am always prepared.	I like order.
I like order.	I pay attention to details.	I pay attention to details.
I am interested in people.	I am interested in people.	I have a soft heart.
I have a soft heart.	I take time out for others.	I take time out for others.



Participant	Form	Age	Sex	Occupation	Genre	Work Music	Volume	Open1	Open2	Open3	Extra1	Extra2	Extra3	Neuro1	Neuro2	Neuro3	Consc1	Consc2	Consc3	Agree1	Agree2	Agree3
1	1	47	F	professor	Classical	Ν	loud	4	4		1	5		1	2		4	2		3	2	
2	1	42	F	musician	Funk	Ν	soft	1	3		2	2		5	3		4	1		2	1	
3	1	27	Μ	student	Jazz	Ν	soft	2	4		5	5		2	4		5	1		4	2	
4	1	29	Μ	server	Metal	Ν	soft	1	3		5	2		2	1		1	1		4	2	
5	1	27	Μ	chef	Rock	Ν	soft	1	4		5	1		2	2		5	3		2	2	
6	2	21	F	painter	Pop	Y	loud	4		4	2		1	1		5	1		5	5		3
7	2	39	F	librarian	Alt	Ν	loud	1		4	4		3	4		3	4		2	4		3
8	2	22	F	server	Ska	Ν	soft	4		2	3		3	3		3	1		2	5		5
9	2	38	Μ	doctor	Punk	Ν	loud	1		3	2		2	2		4	4		1	3		2
10	2	29	F	statistician	Pop	Ν	loud	4		5	3		4	5		4	3		2	3		1
11	3	28	F	chef	Rock	Y	loud		3	3		5	5		5	4		3	3		2	5
12	3	25	Μ	nurse	Rock	N	soft		4	5		2	2		2	5		4	5		3	5
13	3	29	Μ	lawyer	Jazz	Y	soft		3	4		3	2		4	5		4	5		1	2
14	3	38	F	accountant	Metal	Ν	soft		3	1		1	2		3	3		4	4		5	4
15	3	21	F	secretary	Alt	Ν	loud		4	4		1	2		1	1		5	3		4	5

Expansions of 3-form Design



Table 3

Ten-Form, Six-Set Variation of the Split Questionnaire Survey Design, With X Set

Form 1	Item set												
	Х	А	В	С	D	E							
	1	1	1	0	0	0							
2	1	1	0	1	0	0							
3	1	1	0	0	1	0							
4	1	1	0	0	0	1							
5	1	0	1	1	0	0							
6	1	0	1	0	1	0							
7	1	0	1	0	0	1							
8	1	0	0	1	1	0							
9	1	0	0	1	0	1							
10	1	0	0	0	1	1							

Note. 1 = questions asked; 0 = questions not asked.



Table 3. A two-method planned missing design

Prop. Cheap Items Expensive Items

Administered 2/3

Planned Missing

Administered 1/3

Administered

Note. The proportions of subjects (prop.) who receive variables across the different methods can vary as needed.

2-Method Measurement



Expensive Measure 1

- Gold standard highly valid (unbiased) measure of the construct under investigation
- Problem: Measure 1 is time-consuming and/or costly to collect, so it is not feasible to collect from a large sample
 Inexpensive Measure 2
 - Practical inexpensive and/or quick to collect on a large sample
 - Problem: Measure 2 is systematically biased so not ideal

2-Method Measurement

e.g., measuring stress

- Expensive Measure 1 = collect spit samples, measure cortisol
- Inexpensive Measure 2 = survey querying stressful thoughts
- e.g., measuring intelligence
 - Expensive Measure 1 = WAIS IQ scale
 - Inexpensive Measure 2 = multiple choice IQ test
- e.g., measuring smoking
 - Expensive Measure 1 = carbon monoxide measure
 - Inexpensive Measure 2 = self-report
- e.g., measuring student attention
 - Expensive Measure 1 = Classroom observations
 - Inexpensive Measure 2 = Teacher report



2 Method Planned Missing Data



A) Smoking Cessation as an example

B) Stress as an example



