*Missing Data Workshop*
Joint Doctoral Program in Clinical Psyc

# Missing Data Workshop

Jonathan Lee Helm
Friday May 17th, 2019

@jonathanleehelm 1

---

*Grand Overview*

- Introduction to missing data
- Review of sampling
- Patterns, causes, and mechanisms of missing data
- The problem with missing data

@jonathanleehelm 2

---

*Grand Overview*

- Introduction to missing data
- Review of sampling
- Patterns, causes, and mechanisms of missing data
- The problem with missing data

@jonathanleehelm 3

---

*Missing Data Workshop*
Joint Doctoral Program in Clinical Psyc

# Introduction to Missing Data

@jonathanleehelm 4

## Missing Data
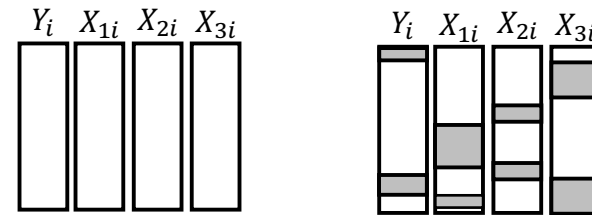
- Missing data occur when we do not have one or more observations for a given variable in our data

## Missing Data

- Missing data occur when we do not have one or more observations for a given variable in our data

$$Y_i \quad X_{1i} \quad X_{2i} \quad X_{3i} \qquad Y_i \quad X_{1i} \quad X_{2i} \quad X_{3i}$$

## Missing Data

- Missing data are ubiquitous in psychological science

- Can anyone think of a real world example that has complete data in psychological science?

## Missing Data

- Missing data pose potential problems

## Missing Data

- Missing data pose potential problems

- Potential for biased estimates
- Potentially increase standard errors
  - Smaller sample size
  - Higher Type 2 error rate (i.e., harder to detect sig. effects)

## Missing Data

- Missing data pose potential problems

- Potential for biased estimates
- Potentially increase standard errors
  - Smaller sample size
  - Higher Type 2 error rate (i.e., harder to detect sig. effects)

- These potential problems can be mitigated by different analyses
  - Multiple imputation

## Missing Data

- The reasons that missing data may cause problems are closely linked to the importance of random sampling

- Let's review random sampling as a segue into missing data

## Grand Overview

- Introduction to missing data
- Review of sampling
- Patterns, causes, and mechanisms of missing data
- The problem with missing data

5/16/19

---

*Missing Data Workshop*
Joint Doctoral Program in Clinical Psyc

# Review of Sampling

@jonathanleehelm
13

---

*Sampling*

- The problems associated with missing data can be conceptualized in terms of sampling

@jonathanleehelm
14

---

*Sampling*

- Scientific process:
  1. Random sample
  2. Measure
  3. Analyze
  4. Draw conclusions

@jonathanleehelm
15

---

*Sampling*

- Scientific process:
  1. Random sample
  2. Measure
  3. Analyze
  4. Draw conclusions

This is an important feature!

Random sampling implies generalizability

@jonathanleehelm
16

---

4

## Sampling



Random Sample
- All members of pop. have equal chance of being selected
- The results from the analysis should generalize to the population

## Sampling

- Example of a random sample
- And a random sample of a random sample

### Population Data

| i | JS | IQ |
|---|-----|-----|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| N-1 | 16 | 118 |
| N | 12 | 134 |

### Population Data

| i | JS | IQ | |
|---|-----|-----|---|
| 1 | 9 | 78 | Pop. mean for JS = 10 |
| 2 | 13 | 84 | Pop. SD for JS = 3 |
| 3 | 10 | 84 | |
| 4 | 8 | 85 | |
| 5 | 7 | 87 | |
| 6 | 7 | 91 | |
| 7 | 9 | 92 | |
| 8 | 9 | 94 | |
| 9 | 11 | 94 | |
| 10 | 7 | 96 | |
| 11 | 7 | 99 | |
| ⋮ | ⋮ | ⋮ | |
| N-1 | 16 | 118 | |
| N | 12 | 134 | |

**Slide 21:**

**Population Data**

| i | JS | IQ |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| N-1 | 16 | 118 |
| N | 12 | 134 |

Pop. mean for JS = 10
Pop. SD for JS = 3

**Random Sample**

| i | JS | IQ |
|---|---|---|
| 1 | 10 | 72 |
| 2 | 12 | 74 |
| 3 | 7 | 90 |
| ⋮ | ⋮ | ⋮ |
| 99 | 9 | 118 |
| 100 | 14 | 134 |

@jonathanleehelm
21

**Slide 22:**

**Population Data**

| i | JS | IQ |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| N-1 | 16 | 118 |
| N | 12 | 134 |

Pop. mean for JS = 10
Pop. SD for JS = 3

**Random Sample**

| i | JS | IQ |
|---|---|---|
| 1 | 10 | 72 |
| 2 | 12 | 74 |
| 3 | 7 | 90 |
| ⋮ | ⋮ | ⋮ |
| 99 | 9 | 118 |
| 100 | 14 | 134 |

The sample mean for JS = 10.3
The sample SD for JS = 2.8

@jonathanleehelm
22

**Slide 23:**

**Population Data**

| i | JS | IQ |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| N-1 | 16 | 118 |
| N | 12 | 134 |

Pop. mean for JS = 10
Pop. SD for JS = 3

Generalizability is a function of sampling

**Random Sample**

| i | JS | IQ |
|---|---|---|
| 1 | 10 | 72 |
| 2 | 12 | 74 |
| 3 | 7 | 90 |
| ⋮ | ⋮ | ⋮ |
| 99 | 9 | 118 |
| 100 | 14 | 134 |

The sample mean for JS = 10.3
The sample SD for JS = 2.8

@jonathanleehelm
23

**Slide 24:**

**Population Data**

| i | JS | IQ |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| N-1 | 16 | 118 |
| N | 12 | 134 |

Pop. mean for JS = 10
Pop. SD for JS = 3

Generalizability is a function of sampling

Sample mean for JS = 10.3
Sample SD for JS = 2.8

**Random Sample**

| i | JS | IQ |
|---|---|---|
| 1 | 10 | 72 |
| 2 | 12 | 74 |
| 3 | 7 | 90 |
| ⋮ | ⋮ | ⋮ |
| 99 | 9 | 118 |
| 100 | 14 | 134 |

Sample mean for JS = 10.2
Sample SD for JS = 3.1

**Random Sample of Random Sample**

| i | JS | IQ |
|---|---|---|
| 1 | 10 | 72 |
| ⋮ | ⋮ | ⋮ |
| 49 | 9 | 118 |
| 50 | 14 | 134 |

@jonathanleehelm
24

## Slide 25

### *Sampling*

- Example of a non-random sample
- And a non-random sample of a random sample

@jonathanleehelm                                          25

## Slide 26

**Population Data**

| i | JS | IQ |
|---|----|----|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| N-1 | 16 | 118 |
| N | 12 | 134 |

Pop. mean for JS = 10
Pop. SD for JS = 3

@jonathanleehelm                                          26

## Slide 27

**Population Data**

| i | JS | IQ |
|---|----|----|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| N-1 | 16 | 118 |
| N | 12 | 134 |

Pop. mean for JS = 10
Pop. SD for JS = 3

**Non Random Sample**

| i | JS | IQ |
|---|----|----|
| 1 | 10 | 72 |
| 2 | 12 | 74 |
| 3 | 7 | 90 |
| ⋮ | ⋮ | ⋮ |
| 99 | 4 | 80 |
| 100 | 6 | 73 |

@jonathanleehelm                                          27

## Slide 28

**Population Data**

| i | JS | IQ |
|---|----|----|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| N-1 | 16 | 118 |
| N | 12 | 134 |

Pop. mean for JS = 10
Pop. SD for JS = 3

**Non Random Sample**

| i | JS | IQ |
|---|----|----|
| 1 | 10 | 72 |
| 2 | 12 | 74 |
| 3 | 7 | 90 |
| ⋮ | ⋮ | ⋮ |
| 99 | 4 | 80 |
| 100 | 6 | 73 |

Sample mean for JS = 7.1
Sample SD for JS = 1.1

@jonathanleehelm                                          28

7

## Slide 29

**Population Data**

| i | JS | IQ |
|---|----|----|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| N-1 | 16 | 118 |
| N | 12 | 134 |

Pop. mean for JS = 10
Pop. SD for JS = 3

Generalizability does not hold!

**Non Random Sample**

| i | JS | IQ |
|---|----|----|
| 1 | 10 | 72 |
| 2 | 12 | 74 |
| 3 | 7 | 90 |
| ⋮ | ⋮ | ⋮ |
| 99 | 4 | 80 |
| 100 | 6 | 73 |

Sample mean for JS = 7.1
Sample SD for JS = 1.1

@jonathanleehelm    29

## Slide 30

**Population Data**

| i | JS | IQ |
|---|----|----|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| N-1 | 16 | 118 |
| N | 12 | 134 |

Pop. mean for JS = 10
Pop. SD for JS = 3

Generalizability is a function of sampling

Sample mean for JS = 10.3
Sample SD for JS = 2.8

**Random Sample**

| i | JS | IQ |
|---|----|----|
| 1 | 10 | 72 |
| 2 | 12 | 74 |
| 3 | 7 | 90 |
| ⋮ | ⋮ | ⋮ |
| 99 | 9 | 118 |
| 100 | 14 | 134 |

Sample mean for JS = 7.1
Sample SD for JS = 1.1

**Non-Random Sample of Random Sample**

| i | JS | IQ |
|---|----|----|
| 1 | 10 | 72 |
| ⋮ | ⋮ | ⋮ |
| 49 | 9 | 80 |
| 50 | 14 | 73 |

@jonathanleehelm    30

## Slide 31

### *Sampling*

- Missing data, in some cases can be conceptualized as a sample of a sample
- This occurs when we use listwise deletion
  - *Delete any row that has a missing value*

@jonathanleehelm    31

## Slide 32

**Population Data**

| i | JS | IQ |
|---|----|----|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| N-1 | 16 | 118 |
| N | 12 | 134 |

**Random Sample**

| i | JS | IQ |
|---|----|----|
| 1 | 10 | 72 |
| 2 | 12 | 74 |
| 3 | 7 | 90 |
| ⋮ | ⋮ | ⋮ |
| 99 | 9 | 118 |
| 100 | 14 | 134 |

**Sample Remaining after Listwise Deletion**

| i | JS | IQ |
|---|----|----|
| 1 | 10 | 72 |
| ⋮ | ⋮ | ⋮ |
| 49 | 9 | 80 |
| 50 | 14 | 73 |

@jonathanleehelm    32

8

**Slide 33**

**Population Data**

| i | JS | IQ |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| N-1 | 16 | 118 |
| N | 12 | 134 |

Listwise deletion will produce unbiased results if the remaining sample is still a random sample

**Random Sample**

| i | JS | IQ |
|---|---|---|
| 1 | 10 | 72 |
| 2 | 12 | 74 |
| 3 | 7 | 90 |
| ⋮ | ⋮ | ⋮ |
| 99 | 9 | 118 |
| 100 | 14 | 134 |

**Sample Remaining after Listwise Deletion**

| i | JS | IQ |
|---|---|---|
| 1 | 10 | 72 |
| ⋮ | ⋮ | ⋮ |
| 49 | 9 | 80 |
| 50 | 14 | 73 |

@jonathanleehelm    33

---

**Slide 34**

## *Sampling*

- Take aways:
  1. We have to perform sampling
  2. Our results will generalize if we have a random sample
  3. Missing data with listwise deletion can be conceptualized as a sample of a sample

@jonathanleehelm    34

---

**Slide 35**

## *Grand Overview*

- Introduction to missing data
- Review of sampling
- Patterns, causes, and mechanisms of missing data
- The problem with missing data

@jonathanleehelm    35

---

**Slide 36**

## *Missing Data Workshop*
## Joint Doctoral Program in Clinical Psyc

# Patterns, Causes, and Mechanisms of Missing Data

@jonathanleehelm    36

## Patterns vs Causes vs Mechanisms

- Patterns of missingness
- Causes of missingness
- Missing data mechanisms
  - *Missing data assumptions*

@jonathanleehelm 37

## Patterns vs Causes vs Mechanisms

- Patterns of missingness
- Causes of missingness
- Missing data mechanisms
  - *Missing data assumptions*

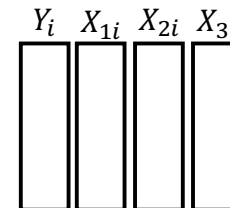@jonathanleehelm 38

## Patterns of Missingness

- A pattern of missingness is a pattern that occurs in a specific data set

@jonathanleehelm 39

## Patterns of Missingness

A pattern of missingness is a pattern that occurs in a specific data set

$$Y_i \quad X_{1i} \quad X_{2i} \quad X_{3i}$$

Complete data
No missingness
No missingness pattern

@jonathanleehelm 40

## Patterns of Missingness

A pattern of missingness is a pattern that occurs in a specific data set

$Y_i$ $X_{1i}$ $X_{2i}$ $X_{3i}$

Missing on $X_{3i}$
Univariate pattern

@jonathanleehelm                                                           41

## Patterns of Missingness

A pattern of missingness is a pattern that occurs in a specific data set

$Y_i$ $X_{1i}$ $X_{2i}$ $X_{3i}$

Missing on more than one variable
General pattern

@jonathanleehelm                                                           42

## Patterns of Missingness

In psychology, we usually have general patterns

We would like methods that can account for general patterns

$Y_i$ $X_{1i}$ $X_{2i}$ $X_{3i}$

@jonathanleehelm                                                           43

## Patterns vs Causes vs Mechanisms

- Patterns of missingness
- Causes of missingness
- Missing data mechanisms
  - *Missing data assumptions*

@jonathanleehelm                                                           44

## Cause of Missingness

- The **_true_** (not assumed) reason why the data are missing

## Cause of Missingness

- Example 1
- A researcher asks individuals to report their income
- Individuals with lower income tend to not report their income

- Cause of missingness: Those with lower income do not report their income

## Cause of Missingness

- Example 2
- A researcher performs an intervention on smoking
- Individuals that find the intervention challenging drop out
- There is missing data on smoking behavior at follow up

- Cause of missingness: Those that find the intervention to be challenging drop out

## Cause of Missingness

- Example 3
- A researcher performs an intervention on smoking
- Individuals that find the intervention challenging drop out
- There is missing data on smoking behavior at follow up

- Cause of missingness: Those that find the intervention to be challenging drop out

## Cause of Missingness

- Example 4
- A researcher performs measures intelligence longitudinally
- The researcher randomly assigns half of the sample to be measured at ages 5, and 7; and half at 6 and 8

- Cause of missingness: Planned missingness: The researcher creates the missingness by design

@jonathanleehelm                                                                 49

## Cause of Missingness

- The **_true_** (not assumed) reason why the data are missing

- Typically not known to the researcher
  - Counter: Planned missingness

@jonathanleehelm                                                                 50

## Patterns vs Causes vs Mechanisms

- Patterns of missingness
- Causes of missingness
- Missing data mechanisms
  - *Missing data assumptions*

@jonathanleehelm                                                                 51

## Missing Data Mechanisms

- A.K.A:
  1. Categories of missingness
  2. Types of missingness

- These are **_assumptions_** regarding the missing data
- Not known causes of the missingness

@jonathanleehelm                                                                 52

## Missing Data Mechanisms

- Three Missing Data Mechanisms

1. Missing Completely at Random (MCAR)
2. Missing at Random (MAR)
3. Missing not at Random (MNAR)

@jonathanleehelm
53

## Missing Data Mechanisms

- Three Missing Data Mechanisms

1. Missing Completely at Random (MCAR)
2. Missing at Random (MAR)
3. Missing not at Random (MNAR)

@jonathanleehelm
54

## Missing Data Mechanisms

- A given analysis (e.g., *t*-test) *__inherently assumes__* a missing data mechanism
- If the assumption is incorrect, the parameter estimates may be biased

@jonathanleehelm
55

## Missing Data Mechanisms: MCAR

- Missingness on a variable $Y_i$ is MCAR if the *probability* of missingness is *unrelated* to
  1. The values $Y_i$ (including the missing values!)
  2. Or to any other variable in the analysis*

\* You can have variables in your data set that are not in the analysis

@jonathanleehelm
56

## Missing Data Mechanisms: MCAR

- Missingness on a variable $Y_i$ is MCAR if the *probability* of missingness is *unrelated* to
  1. The values $Y_i$ (including the missing values!)
  2. Or to any other variable in the analysis*

- ***The observed are a random sub-sample of the complete sample***

  * You can have variables in your data set that are not in the analysis

@jonathanleehelm 57

### Complete Data

| i | JS$^{com}$ | IQ$^{com}$ |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | 12 | 134 |

@jonathanleehelm 58

### Complete Data

| i | JS$^{com}$ | IQ$^{com}$ |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | 12 | 134 |

### Observed Data

| i | JS$^{obs}$ | IQ$^{obs}$ |
|---|---|---|
| 1 | -- | 78 |
| 2 | 13 | 84 |
| 3 | -- | 84 |
| 4 | 8 | 85 |
| 5 | | |
| | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | -- | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | -- | 134 |

@jonathanleehelm 59

### Complete Data

| i | JS$^{com}$ | IQ$^{com}$ |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | 12 | 134 |

### Observed Data

| i | JS$^{obs}$ | IQ$^{obs}$ | JS$^{ind}$ |
|---|---|---|---|
| 1 | -- | 78 | 1 |
| 2 | 13 | 84 | 0 |
| 3 | -- | 84 | 1 |
| 4 | 8 | 85 | 0 |
| 5 | | | 0 |
| | 7 | 87 | 0 |
| 6 | 7 | 91 | 0 |
| 7 | 9 | 92 | 0 |
| 8 | 9 | 94 | 0 |
| 9 | 11 | 94 | 0 |
| 10 | -- | 96 | 1 |
| 11 | 7 | 99 | 0 |
| ⋮ | ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 | 0 |
| 20 | -- | 134 | 1 |

@jonathanleehelm 60

## Slide 61

| Complete Data | | | Observed Data | | | | MCAR: |
|---|---|---|---|---|---|---|---|
| $i$ | JS$^{com}$ | IQ$^{com}$ | $i$ | JS$^{obs}$ | IQ$^{obs}$ | JS$^{ind}$ | $Y^{com}$ is not related to $Y^{ind}$ |
| 1 | 9 | 78 | 1 | -- | 78 | 1 | or any other variable in |
| 2 | 13 | 84 | 2 | 13 | 84 | 0 | the analysis |
| 3 | 10 | 84 | 3 | -- | 84 | 1 | |
| 4 | 8 | 85 | 4 | 8 | 85 | 0 | |
| 5 | 7 | 87 | 5 | 7 | 87 | 0 0 | |
| 6 | 7 | 91 | 6 | 7 | 91 | 0 | |
| 7 | 9 | 92 | 7 | 9 | 92 | 0 | |
| 8 | 9 | 94 | 8 | 9 | 94 | 0 | |
| 9 | 11 | 94 | 9 | 11 | 94 | 0 | |
| 10 | 7 | 96 | 10 | -- | 96 | 1 | |
| 11 | 7 | 99 | 11 | 7 | 99 | 0 | |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | |
| 19 | 16 | 118 | 19 | 16 | 118 | 0 | |
| 20 | 12 | 134 | 20 | -- | 134 | 1 | |

61

## Slide 62

| Complete Data | | | Observed Data | | | | MCAR: |
|---|---|---|---|---|---|---|---|
| $i$ | JS$^{com}$ | IQ$^{com}$ | $i$ | JS$^{obs}$ | IQ$^{obs}$ | JS$^{ind}$ | $Y^{com}$ is not related to $Y^{ind}$ |
| 1 | 9 | 78 | 1 | -- | 78 | 1 | |
| 2 | 13 | 84 | 2 | 13 | 84 | 0 | ***t-test:*** |
| 3 | 10 | 84 | 3 | -- | 84 | 1 | Mean JS$^{com}$ for JS$^{ind}$ = 0 |
| 4 | 8 | 85 | 4 | 8 | 85 | 0 | 10.6 |
| 5 | 7 | 87 | 5 | 7 | 87 | 0 0 | |
| 6 | 7 | 91 | 6 | 7 | 91 | 0 | Mean JS$^{com}$ for JS$^{ind}$ = 1 |
| 7 | 9 | 92 | 7 | 9 | 92 | 0 | 9.6 |
| 8 | 9 | 94 | 8 | 9 | 94 | 0 | |
| 9 | 11 | 94 | 9 | 11 | 94 | 0 | |
| 10 | 7 | 96 | 10 | -- | 96 | 1 | $p$-value for difference |
| 11 | 7 | 99 | 11 | 7 | 99 | 0 | .3857 |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | |
| 19 | 16 | 118 | 19 | 16 | 118 | 0 | |
| 20 | 12 | 134 | 20 | -- | 134 | 1 | |

62

## Slide 63

| Complete Data | | | Observed Data | | | | MCAR: |
|---|---|---|---|---|---|---|---|
| $i$ | JS$^{com}$ | IQ$^{com}$ | $i$ | JS$^{obs}$ | IQ$^{obs}$ | JS$^{ind}$ | $Y^{com}$ is not related to $Y^{ind}$ |
| 1 | 9 | 78 | 1 | -- | 78 | 1 | Or any other variable |
| 2 | 13 | 84 | 2 | 13 | 84 | 0 | ***t-test:*** |
| 3 | 10 | 84 | 3 | -- | 84 | 1 | Mean IQ$^{com}$ for JS$^{ind}$ = 0 |
| 4 | 8 | 85 | 4 | 8 | 85 | 0 | 99.73 |
| 5 | 7 | 87 | 5 | 7 | 87 | 0 0 | |
| 6 | 7 | 91 | 6 | 7 | 91 | 0 | |
| 7 | 9 | 92 | 7 | 9 | 92 | 0 | Mean IQ$^{com}$ for JS$^{ind}$ = 1 |
| 8 | 9 | 94 | 8 | 9 | 94 | 0 | 100.8 |
| 9 | 11 | 94 | 9 | 11 | 94 | 0 | |
| 10 | 7 | 96 | 10 | -- | 96 | 1 | |
| 11 | 7 | 99 | 11 | 7 | 99 | 0 | $p$-value for difference |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | .9235 |
| 19 | 16 | 118 | 19 | 16 | 118 | 0 | |
| 20 | 12 | 134 | 20 | -- | 134 | 1 | |

63

## Slide 64

### *Missing Data Mechanisms: MCAR*

- MCAR: missingness is not related to any variable in the data set
  - The observed are a random sample of your sample

64

## Missing Data Mechanisms: MCAR

- MCAR: missingness is not related to any variable in the data set
  - The observed are a random sample of your sample

- Can we ever know that data are MCAR in practice?

65

## Missing Data Mechanisms: MCAR

- MCAR: missingness is not related to any variable in the data set
  - The observed are a random sample of your sample

- Can we ever know that data are MCAR in practice?
  - *No, we would need the complete data*

66

## Missing Data Mechanisms

- Three Missing Data Mechanisms

1. Missing Completely at Random (MCAR)
2. Missing at Random (MAR)
3. Missing not at Random (MNAR)

67

## Missing Data Mechanisms: MAR

- Missingness on a variable $Y_i$ is ***Missing at Random*** if the probability of missingness is unrelated to $Y_i$ ***after controlling for other variables*** in the analysis

68

17

## Missing Data Mechanisms: MAR

- Missingness on a variable $Y_i$ is **_Missing at Random_** if the probability of missingness is unrelated to $Y_i$ **_after controlling for other variables_** in the analysis

- Missingness on $Y_i$ is related to another variable in the analysis
- You have other measures of the other variable

@jonathanleehelm 69

**Complete Data**

| $i$ | JS$^{com}$ | IQ$^{com}$ |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | 12 | 134 |

@jonathanleehelm 70

**Complete Data**

| $i$ | JS$^{com}$ | IQ$^{com}$ |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | 12 | 134 |

**Observed Data**

| $i$ | JS$^{obs}$ | IQ$^{obs}$ |
|---|---|---|
| 1 | -- | 78 |
| 2 | -- | 84 |
| 3 | -- | 84 |
| 4 | -- | 85 |
| 5 | -- | 87 |
| 6 | -- | 91 |
| 7 | -- | 92 |
| 8 | -- | 94 |
| 9 | -- | 94 |
| 10 | -- | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | 12 | 134 |

@jonathanleehelm 71

**Complete Data**

| $i$ | JS$^{com}$ | IQ$^{com}$ |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | 12 | 134 |

**Observed Data**

| $i$ | JS$^{obs}$ | IQ$^{obs}$ | JS$^{ind}$ |
|---|---|---|---|
| 1 | -- | 78 | 1 |
| 2 | -- | 84 | 1 |
| 3 | -- | 84 | 1 |
| 4 | -- | 85 | 1 |
| 5 | -- | 87 | 1 1 |
| 6 | -- | 91 | 1 |
| 7 | -- | 92 | 1 |
| 8 | -- | 94 | 1 |
| 9 | -- | 94 | 1 |
| 10 | -- | 96 | 1 |
| 11 | 7 | 99 | 0 |
| ⋮ | ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 | 0 |
| 20 | 12 | 134 | 0 |

@jonathanleehelm 72

18

## Slide 73

**Complete Data**

| i | JS^com | IQ^com |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | 12 | 134 |

**Observed Data**

| i | JS^obs | IQ^obs | JS^ind |
|---|---|---|---|
| 1 | -- | 78 | 1 |
| 2 | -- | 84 | 1 |
| 3 | -- | 84 | 1 |
| 4 | -- | 85 | 1 |
| 5 | -- | 87 | 1, 1 |
| 6 | -- | 91 | 1 |
| 7 | -- | 92 | 1 |
| 8 | -- | 94 | 1 |
| 9 | -- | 94 | 1 |
| 10 | -- | 96 | 1 |
| 11 | 7 | 99 | 0 |
| ⋮ | ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 | 0 |
| 20 | 12 | 134 | 0 |

**MAR:**
$Y^{com}$ is not related to $Y^{ind}$ after controlling for other variables in the analysis

@jonathanleehelm — 73

## Slide 74

**Complete Data**

| i | JS^com | IQ^com |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | 12 | 134 |

**Observed Data**

| i | JS^obs | IQ^obs | JS^ind |
|---|---|---|---|
| 1 | -- | 78 | 1 |
| 2 | -- | 84 | 1 |
| 3 | -- | 84 | 1 |
| 4 | -- | 85 | 1 |
| 5 | -- | 87 | 1, 1 |
| 6 | -- | 91 | 1 |
| 7 | -- | 92 | 1 |
| 8 | -- | 94 | 1 |
| 9 | -- | 94 | 1 |
| 10 | -- | 96 | 1 |
| 11 | 7 | 99 | 0 |
| ⋮ | ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 | 0 |
| 20 | 12 | 134 | 0 |

**MAR:**
$Y^{com}$ is not related to $Y^{ind}$ after controlling for other variables in the analysis

$Y^{com} = b_0 + b_1 Y^{ind}$

| | Est. | s.e. | p |
|---|---|---|---|
| $b_0$ | 11.7 | .75 | <.01 |
| $b_1$ | -2.7 | 1.05 | .02 |

@jonathanleehelm — 74

## Slide 75

**Complete Data**

| i | JS^com | IQ^com |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | 12 | 134 |

**Observed Data**

| i | JS^obs | IQ^obs | JS^ind |
|---|---|---|---|
| 1 | -- | 78 | 1 |
| 2 | -- | 84 | 1 |
| 3 | -- | 84 | 1 |
| 4 | -- | 85 | 1 |
| 5 | -- | 87 | 1, 1 |
| 6 | -- | 91 | 1 |
| 7 | -- | 92 | 1 |
| 8 | -- | 94 | 1 |
| 9 | -- | 94 | 1 |
| 10 | -- | 96 | 1 |
| 11 | 7 | 99 | 0 |
| ⋮ | ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 | 0 |
| 20 | 12 | 134 | 0 |

**MAR:**
$Y^{com}$ is not related to $Y^{ind}$ after controlling for other variables in the analysis

$Y^{com} = b_0 + b_1 IQ^{obs}$

| | Est. | s.e. | p |
|---|---|---|---|
| $b_0$ | 0.07 | 3.79 | .98 |
| $b_1$ | .10 | .04 | .01 |

@jonathanleehelm — 75

## Slide 76

**Complete Data**

| i | JS^com | IQ^com |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | 12 | 134 |

**Observed Data**

| i | JS^obs | IQ^obs | JS^ind |
|---|---|---|---|
| 1 | -- | 78 | 1 |
| 2 | -- | 84 | 1 |
| 3 | -- | 84 | 1 |
| 4 | -- | 85 | 1 |
| 5 | -- | 87 | 1, 1 |
| 6 | -- | 91 | 1 |
| 7 | -- | 92 | 1 |
| 8 | -- | 94 | 1 |
| 9 | -- | 94 | 1 |
| 10 | -- | 96 | 1 |
| 11 | 7 | 99 | 0 |
| ⋮ | ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 | 0 |
| 20 | 12 | 134 | 0 |

**MAR:**
$Y^{com}$ is not related to $Y^{ind}$ after controlling for other variables in the analysis

$Y^{com} = b_0 + b_1 IQ^{obs} + b_2 Y^{ind}$

| | Est. | s.e. | p |
|---|---|---|---|
| $b_0$ | 3.97 | 7.80 | .62 |
| $b_1$ | .07 | .07 | .33 |
| $b_2$ | -1.11 | 1.92 | .57 |

@jonathanleehelm — 76

## Missing Data Mechanisms: MAR

- Missingness on random indicates that there is some other variable in the analyses that accounting for the missingness

## Missing Data Mechanisms: MAR

- Missingness on random indicates that there is some other variable in the analyses that accounting for the missingness

- *i.e., that variable can be a proxy for the cause of the missingness*
- *Once you account for that variable, you have a random sample again*

## Missing Data Mechanisms: MAR

- Can we ever know for certain that data are MAR?

- *No, we would need the complete data to be certain*
- *Even then, it would still be an assumption*

## Missing Data Mechanisms

- Three Missing Data Mechanisms

1. Missing Completely at Random (MCAR)
2. Missing at Random (MAR)
3. Missing not at Random (MNAR)

## Missing Data Mechanisms: MNAR

- Missingness on a variable $Y_i$ is **_Missing Not at Random_** if the probability of missingness is **_still related_** to $Y_i$ **_after controlling for other variables_** in the analysis

@jonathanleehelm                                                                 81

## Missing Data Mechanisms: MAR

- Missingness on a variable $Y_i$ is **_Missing Not at Random_** if the probability of missingness is **_still related_** to $Y_i$ **_after controlling for other variables_** in the analysis

- Missingness on $Y_i$ is still related to $Y_i$ after controlling for other variables

@jonathanleehelm                                                                 82

**Complete Data**

| $i$ | JS$^{com}$ | IQ$^{com}$ |
|-----|-----|-----|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | 12 | 134 |

@jonathanleehelm                                                                 83

**Complete Data**

| $i$ | JS$^{com}$ | IQ$^{com}$ |
|-----|-----|-----|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | 12 | 134 |

**Observed Data**

| $i$ | JS$^{obs}$ | IQ$^{obs}$ |
|-----|-----|-----|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | -- | 91 |
| 7 | -- | 92 |
| 8 | -- | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | 12 | 134 |

@jonathanleehelm                                                                 84

## Slide 85

**Complete Data**

| i | JS$^{com}$ | IQ$^{com}$ |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | 12 | 134 |

**Observed Data**

| i | JS$^{obs}$ | IQ$^{obs}$ | JS$^{ind}$ |
|---|---|---|---|
| 1 | 9 | 78 | 0 |
| 2 | 13 | 84 | 0 |
| 3 | 10 | 84 | 0 |
| 4 | 8 | 85 | 0 |
| 5 | 7 | 87 | 0 / 0 |
| 6 | -- | 91 | 1 |
| 7 | -- | 92 | 1 |
| 8 | -- | 94 | 1 |
| 9 | 11 | 94 | 0 |
| 10 | 7 | 96 | 0 |
| 11 | 7 | 99 | 0 |
| ⋮ | ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 | 0 |
| 20 | 12 | 134 | 0 |

@jonathanleehelm  85

## Slide 86

**Complete Data**

| i | JS$^{com}$ | IQ$^{com}$ |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | 12 | 134 |

**Observed Data**

| i | JS$^{obs}$ | IQ$^{obs}$ | JS$^{ind}$ |
|---|---|---|---|
| 1 | 9 | 78 | 0 |
| 2 | 13 | 84 | 0 |
| 3 | 10 | 84 | 0 |
| 4 | 8 | 85 | 0 |
| 5 | 7 | 87 | 0 / 0 |
| 6 | -- | 91 | 1 |
| 7 | -- | 92 | 1 |
| 8 | -- | 94 | 1 |
| 9 | 11 | 94 | 0 |
| 10 | 7 | 96 | 0 |
| 11 | 7 | 99 | 0 |
| ⋮ | ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 | 0 |
| 20 | 12 | 134 | 0 |

**MNAR:**
$Y^{com}$ is still related to $Y^{ind}$ after controlling for other variables in the analysis

@jonathanleehelm  86

## Slide 87

**Complete Data**

| i | JS$^{com}$ | IQ$^{com}$ |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | 12 | 134 |

**Observed Data**

| i | JS$^{obs}$ | IQ$^{obs}$ | JS$^{ind}$ |
|---|---|---|---|
| 1 | 9 | 78 | 0 |
| 2 | 13 | 84 | 0 |
| 3 | 10 | 84 | 0 |
| 4 | 8 | 85 | 0 |
| 5 | 7 | 87 | 0 / 0 |
| 6 | -- | 91 | 1 |
| 7 | -- | 92 | 1 |
| 8 | -- | 94 | 1 |
| 9 | 11 | 94 | 0 |
| 10 | 7 | 96 | 0 |
| 11 | 7 | 99 | 0 |
| ⋮ | ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 | 0 |
| 20 | 12 | 134 | 0 |

**MNAR:**
$Y^{com}$ is still related to $Y^{ind}$ after controlling for other variables in the analysis

$Y^{com} = b_0 + b_1 Y^{ind}$

| | Est. | s.e. | p |
|---|---|---|---|
| $b_0$ | 11.4 | .51 | <.01 |
| $b_1$ | -4.2 | 1.02 | <.01 |

@jonathanleehelm  87

## Slide 88

**Complete Data**

| i | JS$^{com}$ | IQ$^{com}$ |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | 12 | 134 |

**Observed Data**

| i | JS$^{obs}$ | IQ$^{obs}$ | JS$^{ind}$ |
|---|---|---|---|
| 1 | 9 | 78 | 0 |
| 2 | 13 | 84 | 0 |
| 3 | 10 | 84 | 0 |
| 4 | 8 | 85 | 0 |
| 5 | 7 | 87 | 0 / 0 |
| 6 | -- | 91 | 1 |
| 7 | -- | 92 | 1 |
| 8 | -- | 94 | 1 |
| 9 | 11 | 94 | 0 |
| 10 | 7 | 96 | 0 |
| 11 | 7 | 99 | 0 |
| ⋮ | ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 | 0 |
| 20 | 12 | 134 | 0 |

**MNAR:**
$Y^{com}$ is still related to $Y^{ind}$ after controlling for other variables in the analysis

$Y^{com} = b_0 + b_1 IQ^{obs}$

| | Est. | s.e. | p |
|---|---|---|---|
| $b_0$ | 0.07 | 3.79 | .98 |
| $b_1$ | .10 | .04 | .01 |

@jonathanleehelm  88

**Complete Data**

| i | JS$^{com}$ | IQ$^{com}$ |
|---|---|---|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 |
| 20 | 12 | 134 |

**Observed Data**

| i | JS$^{obs}$ | IQ$^{obs}$ | JS$^{ind}$ |
|---|---|---|---|
| 1 | 9 | 78 | 0 |
| 2 | 13 | 84 | 0 |
| 3 | 10 | 84 | 0 |
| 4 | 8 | 85 | 0 |
| 5 | 7 | 87 | 0 |
| 6 | -- | 91 | 1 |
| 7 | -- | 92 | 1 |
| 8 | -- | 94 | 1 |
| 9 | 11 | 94 | 0 |
| 10 | 7 | 96 | 0 |
| 11 | 7 | 99 | 0 |
| ⋮ | ⋮ | ⋮ | ⋮ |
| 19 | 16 | 118 | 0 |
| 20 | 12 | 134 | 0 |

**MNAR:**
Y$^{com}$ is still related to Y$^{ind}$ after controlling for other variables in the analysis

$$Y^{com} = b_0 + b_1 IQ^{obs} + b_2 Y^{ind}$$

| | Est. | s.e. | p |
|---|---|---|---|
| $b_0$ | 4.88 | 3.30 | .16 |
| $b_1$ | .06 | .03 | .06 |
| $b_2$ | -3.48 | 1.01 | <.01 |

@jonathanleehelm

89

---

## Missing Data Mechanisms: MNAR

- Data are missing not at random when the missingness on $Y_i$ is related to the values of $Y_i$, even after controlling for other variables in the analysis

@jonathanleehelm

90

---

## Missing Data Mechanisms: MAR

- Missingness not at random occurs when the missingness on $Y_i$ is related to the values of $Y_i$, even after controlling for other variables in the analysis

- *i.e., the sample is still not a random sample from the population, even after accounting for other variables in the analysis*

@jonathanleehelm

91

---

## Missing Data Mechanisms: MAR

- Can we ever know for certain that data are MNAR?

@jonathanleehelm

92

23

## Missing Data Mechanisms: MAR

- Can we ever know for certain that data are MNAR?

- *No, we would need the complete data to be certain*
- *Even then, it would still be an assumption*

@jonathanleehelm 93

## Missing Data Mechanisms

- Three Missing Data Mechanisms

1. Missing Completely at Random (MCAR)
2. Missing at Random (MAR)
3. Missing not at Random (MNAR)

@jonathanleehelm 94

## Missing Data Mechanisms

1. Missing Completely at Random (MCAR)
2. Missing at Random (MAR)
3. Missing not at Random (MNAR)

- These are **_assumptions_** regarding missingness
- Different statistical analysis will be valid under different assumptions

@jonathanleehelm 95

## Grand Overview

- Introduction to missing data
- Review of sampling
- Patterns, causes, and mechanisms of missing data
- The problem with missing data

@jonathanleehelm 96

*Missing Data Workshop*
Joint Doctoral Program in Clinical Psyc

# The Problem with Missing Data

@jonathanleehelm                                                    97

---

*Problems with Missing Data*

- Most analytic techniques require complete data
  - Descriptive statistics (e.g., means, SDs, correlations)
  - *t*-tests
  - Significance testing of correlation
  - ANOVA
  - Regression

@jonathanleehelm                                                    98

---

*Problems with Missing Data*

- If our data contain missingness, then we typically perform listwise deletion

@jonathanleehelm                                                    99

---

*Problems with Missing Data*

- If our data contain missingness, then we typically perform listwise deletion
- If data are missing completely at random, then we should obtain unbiased estimates

@jonathanleehelm                                                    100

## Slide 101

**Population Data**

| i | JS | IQ |
|---|----|----|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| N-1 | 16 | 118 |
| N | 12 | 134 |

- Missingness is MCAR
- Sub-sample is a random sample
- Parameter estimates are unbiased ☺
- Standard errors are still larger ☹

**Random Sample**

| i | JS | IQ |
|---|----|----|
| 1 | 10 | 72 |
| 2 | 12 | 74 |
| 3 | 7 | 90 |
| ⋮ | ⋮ | ⋮ |
| 99 | 9 | 118 |
| 100 | 14 | 134 |

**Sample Remaining after Listwise Deletion**

| i | JS | IQ |
|---|----|----|
| 1 | 10 | 72 |
| ⋮ | ⋮ | ⋮ |
| 49 | 9 | 80 |
| 50 | 14 | 73 |

@jonathanleehelm

101

## Slide 102

### *Problems with Missing Data*

- If our data contain missingness, then we typically perform listwise deletion
- If data are missing completely at random, then we should obtain unbiased estimated
- If data are missing at random or missing not at random then we will obtain biased estimates

@jonathanleehelm

102

## Slide 103

**Population Data**

| i | JS | IQ |
|---|----|----|
| 1 | 9 | 78 |
| 2 | 13 | 84 |
| 3 | 10 | 84 |
| 4 | 8 | 85 |
| 5 | 7 | 87 |
| 6 | 7 | 91 |
| 7 | 9 | 92 |
| 8 | 9 | 94 |
| 9 | 11 | 94 |
| 10 | 7 | 96 |
| 11 | 7 | 99 |
| ⋮ | ⋮ | ⋮ |
| N-1 | 16 | 118 |
| N | 12 | 134 |

- Missingness is MAR or MNAR
- Sub-sample is not random a sample
- Parameter estimates can be biased ☹
- Standard errors are larger ☹

**Random Sample**

| i | JS | IQ |
|---|----|----|
| 1 | 10 | 72 |
| 2 | 12 | 74 |
| 3 | 7 | 90 |
| ⋮ | ⋮ | ⋮ |
| 99 | 9 | 118 |
| 100 | 14 | 134 |

**Sample Remaining after Listwise Deletion**

| i | JS | IQ |
|---|----|----|
| 1 | 10 | 72 |
| ⋮ | ⋮ | ⋮ |
| 49 | 9 | 80 |
| 50 | 14 | 73 |

@jonathanleehelm

103

## Slide 104

### *Problems with Missing Data*

- Therefore, many analytic techniques assume data are missing completely at random
  - *And they didn't even tell you!*

@jonathanleehelm

104

## Problems with Missing Data

- So the major question is, is there anything we can do to at least assume MAR instead of the more restrictive MCAR?
  - *Multiple Imputation*

## Grand Overview

- Introduction to missing data
- Review of sampling
- Patterns, causes, and mechanisms of missing data
- The problem with missing data