# PCORI Methodology Standards

Standards for Preventing and Handling of Missing Data

**Presented By**

Tianjin Li, MD, PhD
Johns Hopkins University

Daniel Scharfstein, ScD
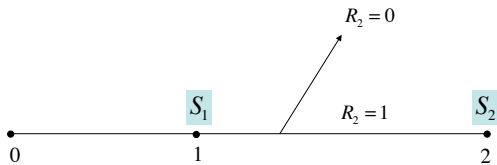Johns Hopkins University

# Module listing

- Module 1: Introduction

- Module 2: What are missing data?

- Module 3: Methods to prevent and monitor missing data

- Module 4: Record and report missing data

- Module 5: Describe statistical methods to handle missing data

- **Module 6: Statistical methods to deal with missing data**

- **Module 7: Examine sensitivity of inferences to missing data methods and assumptions**

## Hypothetical Study - Two Time Points

- Imagine a study in which eligible individuals are to receive a new drug.
- Individuals are expected to return for two post-enrollment visits (V1-V2) at which the presence (1) or absence (0) of symptoms is recorded.
- The goal is to learn about the probability of having symptoms at V2.
- Assume all individuals show up at V1 and some individuals drop out of the study before V2.
- To start, imagine that we conduct this study in "infinite" population so that there is no sampling variability.

# Observed and Unobserved Data



$s_1$

$1$

$p_1$

$0$

$1 - p_1$

# Observed and Unobserved Data

| $S_1$ | $R_2$ |
|-------|-------|
| 1 $p_1$ | 1 <br> ... |
| | 0 |
| 0 $1 - p_1$ | 1 |
| | 0 |

# Observed and Unobserved Data

| $S_1$ | $R_2$ |
|---|---|
| 1<br><br>$p_1$ | 1<br>$q_2(1)$ |
| | 0<br>$1 - q_2(1)$ |
| 0<br><br>$1 - p_1$ | 1<br>$q_2(0)$ |
| | 0<br>$1 - q_2(0)$ |

# Observed and Unobserved Data

| $S_1$ | $R_2$ | $S_2$ |
|---|---|---|
| 1 $p_1$ | 1 $q_2(1)$ | 1 |
| | | 0 |
| | 0 $1 - q_2(1)$ | **1** |
| | | **0** |
| 0 $1 - p_1$ | 1 $q_2(0)$ | 1 |
| | | 0 |
| | 0 $1 - q_2(0)$ | **1** |
| | | **0** |

# Observed and Unobserved Data

| $S_1$ | $R_2$ | $S_2$ |
|---|---|---|
| $1$ $p_1$ | $1$ $q_2(1)$ | $1$ $p_2(1,1)$ |
| | | $0$ $1 - p_2(1,1)$ |
| | $0$ $1 - q_2(1)$ | $\mathbf{1}$ $\boldsymbol{p_2(1,0)}$ |
| | | $\mathbf{0}$ $\mathbf{1 - p_2(1,0)}$ |
| $0$ $1 - p_1$ | $1$ $q_2(0)$ | $1$ $p_2(0,1)$ |
| | | $0$ $1 - p_2(0,1)$ |
| | $0$ $1 - q_2(0)$ | $\mathbf{1}$ $\boldsymbol{p_2(0,0)}$ |
| | | $\mathbf{0}$ $\mathbf{1 - p_2(0,0)}$ |

# Observed and Unobserved Data

| $S_1$ | $R_2$ | $S_2$ | Proportion | |
|---|---|---|---|---|
| 1 $p_1$ | 1 $q_2(1)$ | 1 $p_2(1,1)$ | $f_{111}$ | |
| | | 0 $1-p_2(1,1)$ | $f_{110}$ | |
| | 0 $1-q_2(1)$ | **1** $\mathbf{p_2(1,0)}$ | $\boldsymbol{f_{101}}$ | $f_{10?}$ |
| | | **0** $\mathbf{1-p_2(1,0)}$ | $\boldsymbol{f_{100}}$ | |
| 0 $1-p_1$ | 1 $q_2(0)$ | 1 $p_2(0,1)$ | $f_{011}$ | |
| | | 0 $1-p_2(0,1)$ | $f_{010}$ | |
| | 0 $1-q_2(0)$ | **1** $\mathbf{p_2(0,0)}$ | $\boldsymbol{f_{001}}$ | $f_{00?}$ |
| | | **0** $\mathbf{1-p_2(0,0)}$ | $\boldsymbol{f_{000}}$ | |

## Distribution of Observed Data

- $p_1 = P[S_1 = 1]$

- $q_2(1) = P[R_2 = 1 | S_1 = 1]$

- $q_2(0) = P[R_2 = 1 | S_1 = 0]$

- $p_2(1, 1) = P[S_2 = 1 | S_1 = 1, R_2 = 1]$

- $p_2(1, 0) = P[S_2 = 1 | S_1 = 0, R_2 = 1]$

- $f_{111} = P[S_1 = 1, R_2 = 1, S_2 = 1] = p_1 q_2(1) p_2(1, 1)$

- $f_{110} = P[S_1 = 1, R_2 = 1, S_2 = 0] = p_1 q_2(1)\{1 - p_2(1, 1)\}$

- $f_{10?} = P[S_1 = 1, R_2 = 0, S_2 =?] = p_1\{1 - q_2(1)\}$

- $f_{011} = P[S_1 = 0, R_2 = 1, S_2 = 1] = \{1 - p_1\} q_2(1) p_2(1, 1)$

- $f_{010} = P[S_1 = 0, R_2 = 1, S_2 = 0] = \{1 - p_1\} q_2(1)\{1 - p_2(1, 1)\}$

- $f_{00?} = P[S_1 = 0, R_2 = 0, S_2 =?] = \{1 - p_1\}\{1 - q_2(1)\}$

## Distribution of Unobserved Data

- $\mathbf{p_2(1, 0)} = P[S_2 = 1 | S_1 = 1, R_2 = 0]$
- $\mathbf{p_2(0, 0)} = P[S_2 = 1 | S_1 = 0, R_2 = 0]$

<br>

- $\mathbf{f_{101}} = P[S_1 = 1, R_2 = 0, S_2 = 1] = p_1 \{1 - q_2(1)\} \mathbf{p_2(1, 0)}$
- $\mathbf{f_{100}} = P[S_1 = 1, R_2 = 0, S_2 = 0] = p_1 \{1 - q_2(1)\} \{1 - \mathbf{p_2(1, 0)}\}$
- $\mathbf{f_{001}} = P[S_1 = 0, R_2 = 0, S_2 = 1] = \{1 - p_1\} \{1 - q_2(1)\} \mathbf{p_2(1, 0)}$
- $\mathbf{f_{000}} = P[S_1 = 0, R_2 = 0, S_2 = 0] = \{1 - p_1\} \{1 - q_2(1)\} \{1 - \mathbf{p_2(1, 0)}\}$

| $S_1$ | $R_2$ | $S_2$ | Proportion | |
|---|---|---|---|---|
| 1 $p_1$ | 1 $q_2(1)$ | 1 $p_2(1,1)$ | $f_{111}$ | |
| | | 0 $1-p_2(1,1)$ | $f_{110}$ | |
| | 0 $1-q_2(1)$ | **1** $\mathbf{p_2(1,0)}$ | $\mathbf{f_{101}}$ | $f_{10?}$ |
| | | **0** $\mathbf{1-p_2(1,0)}$ | $\mathbf{f_{100}}$ | |
| 0 $1-p_1$ | 1 $q_2(0)$ | 1 $p_2(0,1)$ | $f_{011}$ | |
| | | 0 $1-p_2(0,1)$ | $f_{010}$ | |
| | 0 $1-q_2(0)$ | **1** $\mathbf{p_2(0,0)}$ | $\mathbf{f_{001}}$ | $f_{00?}$ |
| | | **0** $\mathbf{1-p_2(0,0)}$ | $\mathbf{f_{000}}$ | |

# Fundamental Problem

- Even with infinite data, we cannot learn about the probability of having symptoms at V2.
- We don't know the probability of have symptoms for individuals who have dropped out prior to V2.
- **Need to make assumptions**!
- With assumptions, we can compute $P[S_2 = 1]$

Worst Case

- If $R_2 = 0$ then $S_2 = 1$

# Worst Case

| $S_1$ | $R_2$ | $S_2$ | Proportion |
|---|---|---|---|
| 1 $p_1$ | 1 $q_2(1)$ | 1 $p_2(1,1)$ | $f_{111}$ |
| | | 0 $1 - p_2(1,1)$ | $f_{110}$ |
| | 0 $1 - q_2(1)$ | 1 $p_2(1,0) = 1$ | $f_{101} = f_{10?}$ |
| | | 0 $1 - p_2(1,0) = 0$ | $f_{100} = 0$ |
| 0 $1 - p_1$ | 1 $q_2(0)$ | 1 $p_2(0,1)$ | $f_{011}$ |
| | | 0 $1 - p_2(0,1)$ | $f_{010}$ |
| | 0 $1 - q_2(0)$ | 1 $p_2(0,0) = 1$ | $f_{001} = f_{00?}$ |
| | | 0 $1 - p_2(0,0) = 0$ | $f_{000} = 0$ |

$f_{10?}$

$f_{00?}$

Best Case

- If $R_2 = 0$ then $S_2 = 0$

# Best Case



| $S_1$ | $R_2$ | $S_2$ | Proportion | |
|---|---|---|---|---|
| 1 $p_1$ | 1 $q_2(1)$ | 1 $p_2(1,1)$ | $f_{111}$ | |
| | | 0 $1 - p_2(1,1)$ | $f_{110}$ | |
| | 0 $1 - q_2(1)$ | 1 $p_2(1,0) = 0$ | $f_{101} = 0$ | $f_{10?}$ |
| | | 0 $1 - p_2(1,0) = 1$ | $f_{100} = f_{10?}$ | |
| 0 $1 - p_1$ | 1 $q_2(0)$ | 1 $p_2(0,1)$ | $f_{011}$ | |
| | | 0 $1 - p_2(0,1)$ | $f_{010}$ | |
| | 0 $1 - q_2(0)$ | 1 $p_2(0,0) = 0$ | $f_{001} = 0$ | $f_{00?}$ |
| | | 0 $1 - p_2(0,0) = 1$ | $f_{000} = f_{00?}$ | |

Maintained Response After Dropout

- If $R_2 = 0$, $S_2 = S_1$

| $S_1$ | $R_2$ | $S_2$ | Proportion | |
|---|---|---|---|---|
| 1<br><br>$p_1$ | 1<br><br>$q_2(1)$ | 1<br>$p_2(1,1)$ | $f_{111}$ | |
| | | 0<br>$1 - p_2(1,1)$ | $f_{110}$ | |
| | 0<br><br>$1 - q_2(1)$ | 1<br>$p_2(1,0) = 1$ | $f_{101} = f_{10?}$ | $f_{10?}$ |
| | | 0<br>$1 - p_2(1,0) = 0$ | $f_{100} = 0$ | |
| 0<br><br>$1 - p_1$ | 1<br><br>$q_2(0)$ | 1<br>$p_2(0,1)$ | $f_{011}$ | |
| | | 0<br>$1 - p_2(0,1)$ | $f_{010}$ | |
| | 0<br><br>$1 - q_2(0)$ | 1<br>$p_2(0,0) = 0$ | $f_{001} = 0$ | $f_{00?}$ |
| | | 0<br>$1 - p_2(0,0) = 1$ | $f_{000} = f_{00?}$ | |

Missing at Random (MAR)

$$R_2 \text{ independent of } S_2 \text{ given } S_1$$

$\mathbf{p_2(1,0)} = P[S_2 = 1|S_1 = 1, R_2 = 0] = P[S_2 = 1|S_1 = 1, R_2 = 1] = p_2(1,1)$

$\mathbf{p_2(0,0)} = P[S_2 = 1|S_1 = 0, R_2 = 0] = P[S_2 = 1|S_1 = 0, R_2 = 1] = p_2(0,1)$

# Missing At Random



| $S_1$ | $R_2$ | $S_2$ | Proportion | |
|---|---|---|---|---|
| 1 $p_1$ | 1 $q_2(1)$ | 1 $p_2(1,1)$ | $f_{111}$ | |
| | | 0 $1-p_2(1,1)$ | $f_{110}$ | |
| | 0 $1-q_2(1)$ | 1 $p_2(1,0)=p_2(1,1)$ | $f_{101}$ | $f_{10?}$ |
| | | 0 $1-p_2(1,0)$ | $f_{100}$ | |
| 0 $1-p_1$ | 1 $q_2(0)$ | 1 $p_2(0,1)$ | $f_{011}$ | |
| | | 0 $1-p_2(0,1)$ | $f_{010}$ | |
| | 0 $1-q_2(0)$ | 1 $p_2(0,0)=p_2(0,1)$ | $f_{001}$ | $f_{00?}$ |
| | | 0 $1-p_2(0,0)$ | $f_{000}$ | |

# Missing Not at Random (MNAR)

- Missing at Random doesn't hold
- Best/Worst Case and Maintained Response After Drop-out are MNAR assumptions

# Missing Not at Random (MNAR)

$$\overbrace{P[S_2 = 1 | S_1 = 1, R_2 = 0]}^{\mathbf{p_2(1,0)}}$$
$$\propto \underbrace{P[S_2 = 1 | S_1 = 1, R_2 = 1]}_{p_2(1,1)} \exp(\alpha)$$

$$\overbrace{P[S_2 = 1 | S_1 = 0, R_2 = 0]}^{\mathbf{p_2(0,0)}}$$
$$\propto \underbrace{P[S_2 = 1 | S_1 = 0, R_2 = 1]}_{p_2(0,1)} \exp(\alpha)$$

- Exponential Tilting
- $\alpha$ is a sensitivity analysis parameter
- $\alpha = 0$ corresponds to MAR

# Missing Not at Random (MNAR)

| $\alpha$ | $p_2(1,1)$ $P[S_2 = 1 \mid S_1 = 1, R_2 = 1]$ | $\mathbf{p_2(1,0)}$ $P[S_2 = 1 \mid S_1 = 1, R_2 = 0]$ |
|---|---|---|
| -1 | 0.2 | 0.084 |
| -0.5 | 0.2 | 0.132 |
| 0 | 0.2 | 0.200 |
| 0.5 | 0.2 | 0.292 |
| 1 | 0.2 | 0.405 |

# Missing Not At Random



| $S_1$ | $R_2$ | $S_2$ | Proportion | |
|---|---|---|---|---|
| 1<br>$p_1$ | 1<br>$q_2(1)$ | 1<br>$p_2(1,1)$ | $f_{111}$ | |
| | | 0<br>$1 - p_2(1,1)$ | $f_{110}$ | |
| | 0<br>$1 - q_2(1)$ | 1<br>$p_2(1,0) \propto p_2(1,1)e^\alpha$ | $f_{101}$ | $f_{10?}$ |
| | | 0<br>$1 - p_2(1,0)$ | $f_{100}$ | |
| 0<br>$1 - p_1$ | 1<br>$q_2(0)$ | 1<br>$p_2(0,1)$ | $f_{011}$ | |
| | | 0<br>$1 - p_2(0,1)$ | $f_{010}$ | |
| | 0<br>$1 - q_2(0)$ | 1<br>$p_2(0,0) \propto p_2(0,1)e^\alpha$ | $f_{001}$ | $f_{00?}$ |
| | | 0<br>$1 - p_2(0,0)$ | $f_{000}$ | |

# Inference in Finite Samples

- Under the above assumptions, $P[S_2 = 1]$ depends on the distribution of the observed data.
- Estimate $P[S_2 = 1]$ by plugging-in the estimated distribution of the observed data.
- Standard errors and confidence intervals: Re-sampling methods such as jackknife and bootstrap.

## Case Study

- Women were enrolled in a randomized trial to evaluate two doses (100 and 150 mg) of the contraceptive DMPA.
- 4 doses (administered via injection) were scheduled to be given at 90 day intervals with the first dose at randomization.
- Women were asked to fill out a daily diary recording bleeding/spotting.
- A women was coded as having "amenorrhea" at an injection visit if she did not have bleeding/spotting for 80 consecutive days since the previous injection.
- The analysis population is restricted to the 1151 women who were randomized and returned their first diary.
- We focus on the analysis of the first two diaries.

# Low Dose (Tx 0)

| $S_1$ | $R_2$ | $S_2$ | Proportion | |
|---|---|---|---|---|
| 1 $\frac{107}{576}$ = 19% | 1 $\frac{84}{107}$ = 79% | 1 $\frac{54}{84}$ = 64% | 9.4% | |
| | | 0 $\frac{30}{84}$ = 36% | 5.2% | |
| | 0 $\frac{23}{107}$ = 21% | **1** $p_2(1,0)$ | $f_{101}$ | 4.0% |
| | | **0** $1 - p_2(1,0)$ | $f_{100}$ | |
| 0 $\frac{469}{576}$ = 81% | 1 $\frac{393}{469}$ = 84% | 1 $\frac{71}{393}$ = 18% | 12.3% | |
| | | 0 $\frac{322}{393}$ = 82% | 55.9% | |
| | 0 $\frac{76}{469}$ = 16% | **1** $p_2(0,0)$ | $f_{001}$ | 13.2% |
| | | **0** $1 - p_2(0,0)$ | $f_{000}$ | |

# High Dose (Tx 1)

| $S_1$ | $R_2$ | $S_2$ | Proportion | |
|---|---|---|---|---|
| 1 <br><br> $\dfrac{118}{575} = 21\%$ | 1 <br><br> $\dfrac{87}{118} = 74\%$ | 1 <br><br> $\dfrac{56}{87} = 64\%$ | 9.7% | |
| | | 0 <br><br> $\dfrac{31}{87} = 36\%$ | 5.4% | |
| | 0 <br><br> $\dfrac{31}{118} = 26\%$ | **1** <br><br> $\boldsymbol{p_2(1,0)}$ | $\boldsymbol{f_{101}}$ | 5.4% |
| | | **0** <br><br> $\boldsymbol{1 - p_2(1,0)}$ | $\boldsymbol{f_{100}}$ | |
| 0 <br><br> $\dfrac{457}{575} = 79\%$ | 1 <br><br> $\dfrac{389}{457} = 85\%$ | 1 <br><br> $\dfrac{104}{389} = 27\%$ | 18.1% | |
| | | 0 <br><br> $\dfrac{285}{389} = 73\%$ | 49.6% | |
| | 0 <br><br> $\dfrac{68}{457} = 15\%$ | **1** <br><br> $\boldsymbol{p_2(0,0)}$ | $\boldsymbol{f_{001}}$ | 11.8% |
| | | **0** <br><br> $\boldsymbol{1 - p_2(0,0)}$ | $\boldsymbol{f_{000}}$ | |

# Analysis

**Treatment 0**

**Treatment 1**

- Imagine a study in which eligible individuals are to receive a new drug to relieve symptoms.
- Individuals are expected to return for three post-enrollment visits (V1-V3) at which the presence (1) or absence (0) of symptoms is recorded.
- The goal is to learn about probability of having symptoms at V3.
- Assume all individuals show up at V1 and some individuals drop out of the study before V3.
- To start, imagine that we conduct this study in "infinite" population so that there is no sampling variability.

# Observed and Unobserved Data

| $S_1$ | $R_2$ | $S_2$ | $R_3$ | $S_3$ | Proportion | |
|---|---|---|---|---|---|---|
| | | | 1 $q_2(1,1,1)$ | 1 $p_3(1,1,1,1)$ | $f_{11111}$ | |
| | | | | 0 $1 - p_3(1,1,1,1)$ | $f_{11110}$ | |
| | | 1 $p_2(1,1)$ | 0 $1 - q_2(1,1,1)$ | **1** $p_3(\mathbf{1,1,1,0})$ | $f_{11101}$ | $f_{1110?}$ |
| | | | | **0** $1 - p_3(\mathbf{1,1,1,0})$ | $f_{11100}$ | |
| | 1 $q_2(1)$ | | 1 $q_2(1,1,0)$ | 1 $p_3(1,1,0,1)$ | $f_{11011}$ | |
| | | | | 0 $1 - p_3(1,1,0,1)$ | $f_{11010}$ | |
| 1 $p_1$ | | 0 $1 - p_2(1,1)$ | 0 $1 - q_2(1,1,0)$ | **1** $p_3(\mathbf{1,1,0,0})$ | $f_{11001}$ | $f_{1100?}$ |
| | | | | **0** $1 - p_3(\mathbf{1,1,0,0})$ | $f_{11000}$ | |
| | | **1** $p_2(\mathbf{1,0})$ | 0 $q_2(1,0,1) = 1$ | **1** $p_3(\mathbf{1,0,1,0})$ | $f_{10101}$ | $f_{1010?}$ |
| | | | | **0** $1 - p_3(\mathbf{1,0,1,0})$ | $f_{10100}$ | |
| | 0 $1 - q_2(1)$ | | 0 $q_2(1,0,0) = 1$ | **1** $p_3(\mathbf{1,0,0,0})$ | $f_{10001}$ | |
| | | **0** $1 - p_2(\mathbf{1,0})$ | | **0** $1 - p_3(\mathbf{1,0,0,0})$ | $f_{10000}$ | |

## Assumptions

- Worst Case
- Best Case
- Maintained Response after Dropout
- Missing at Random

$$R_2 \text{ independent } (S_2, S_3) \text{ given } S_1$$

$$R_3 \text{ independent } S_3 \text{ given } R_2 = 1, S_2, S_1$$

- Missing Not at Random: Exponential Tilting

$$R_2 \text{ independent } S_3 \text{ given } S_2, S_1$$

$$P[S_2 = 1 | R_2 = 0, S_1 = s_1] \propto P[S_2 | R_2 = 1, S_1 = s_1] \exp(\alpha)$$

$$P[S_3 = 1 | R_3 = 0, R_2 = 1, S_2 = s_2, S_1 = s_1]$$
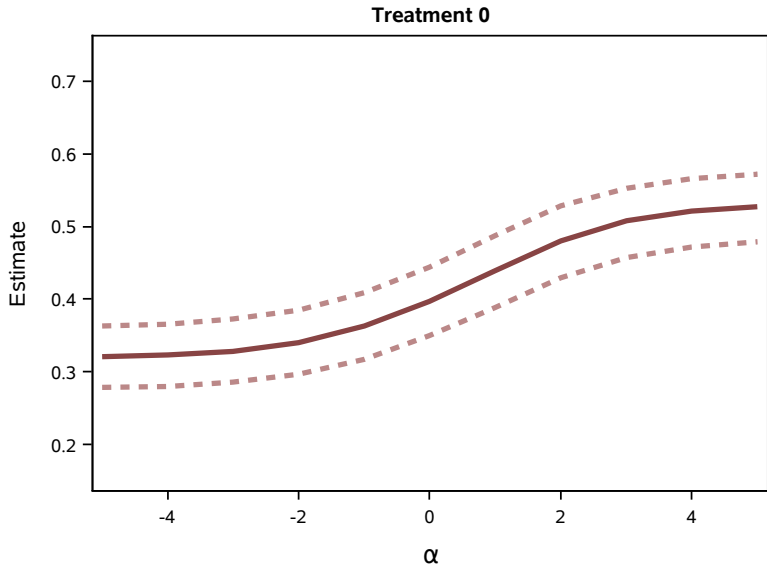$$\propto P[S_3 = 1 | R_3 = 1, S_2 = s_2, S_1 = s_1] \exp(\alpha)$$

# Missing at Random

# Missing Not at Random

- What are the treatment-specific probabilities of symptoms at V3?
- How do these probabilities compare?

**Treatment 0**

Treatment 1
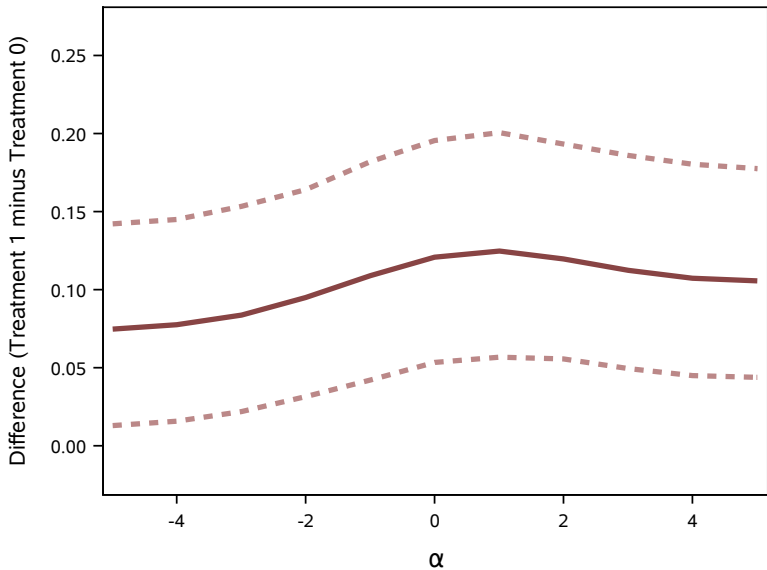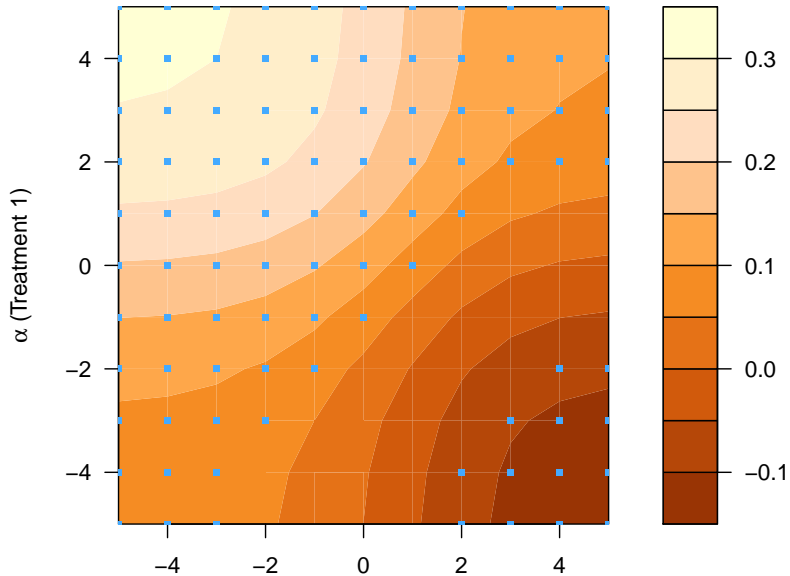
## Other Approaches

All require assumptions!!

- Multiple imputation
  - For each individual, draw from the predictive distribution of the missing outcomes given observed outcomes
  - Perform multiple times to generate a series of datasets with complete data
  - Analyze each complete dataset using standard methods
  - Combine results
- Likelihood-based
  - Mixed models
  - Pattern-mixture models
- Estimating equations
  - Inverse-weighted estimators
  - Doubly-robust estimators

The set of possible assumptions about the missing data mechanism is very large and cannot be fully explored. There are different approaches to sensitivity analysis:

- Ad-hoc
- Local
- Global

# Ad-hoc Sensitivity Analysis

- Analyzing data using a few different analytic methods and evaluate whether the resulting inferences are consistent.
- The problem with this approach is that the assumptions that underlie these methods are very strong and for many of these methods unreasonable.
- More importantly, just because the inferences are consistent does not mean that there are no other reasonable assumptions under which the inference about the treatment effect is different.

- Specify a reasonable benchmark assumption (e.g., missing at random) and evaluate the robustness of the results within a small neighborhood of this assumption.
- What if there are assumptions outside the local neighborhood which are plausible?

# Global Sensitivity Analysis

- Evaluate robustness of results across a much broader range of assumptions that include a reasonable benchmark assumption and a collection of additional assumptions that trend toward best and worst case assumptions.
- Emphasized in Chapter 5 of the NRC report.
- This approach is substantially more informative because it operates like "stress testing" in reliability engineering, where a product is systematically subjected to increasingly exaggerated forces/conditions in order to determine its breaking point.

# Global Sensitivity Analysis

- In the missing data setting, global sensitivity analysis allows one to see how far one needs to deviate from the benchmark assumption in order for inferences to change.
- "Tipping point" analysis
- If the assumptions under which the inferences change are judged to be sufficiently far from the benchmark assumption, then greater credibility is lent to the benchmark analysis; if not, the benchmark analysis can be considered to be fragile.

# PCORI Standards

- Properly account for statistical uncertainty
- Single imputation (e.g., last observation carried forward) should not be the primary analytic approach
- Examine sensitivity to assumptions

# Properly account for statistical uncertainty

- Statistical inference of intervention effects or measures of association should account for statistical uncertainty attributable to missing data.
- This means that methods used for imputing missing data should have valid type I error rates and that confidence intervals have the nominal coverage properties.
- This standard applies to all study designs for any type of research question.

# Single imputation should not be the primary analytic approach

- Single imputation methods like last observation carried forward and baseline observation carried forward generally should not be used as the primary approach for handling missing data in the analysis.
- This standard applies to all study designs for any type of research question.

# Examine sensitivity to assumptions

- Examining sensitivity to the assumptions about the missing data mechanism (i.e., sensitivity analysis) should be a mandatory component of the study protocol, analysis, and reporting.
- This standard applies to all study designs for any type of research question.