



**MICHELLE BYRNE**

Senior Lecturer  
Monash University

[michelle.byrne@monash.edu](mailto:michelle.byrne@monash.edu)

# COVID Hot Mess & Best Practices for Missing Data in Longitudinal Studies

**ABCD Workshop on  
Brain Development  
and Mental Health**



**Sponsored by NIMH**

**June 22nd - July 22nd,  
2021**

# A lot of this material comes from...

- 1) Matta, T.H., Flournoy, J.C., Byrne, M.L., 2018. Making an unknown unknown a known unknown: Missing data in longitudinal neuroimaging studies, *DCN*

<https://pubmed.ncbi.nlm.nih.gov/29129673/>

(actually check out our whole special issue: Methodological Challenges in Developmental Neuroimaging: Contemporary Approaches and Solutions

<https://www.sciencedirect.com/journal/developmental-cognitive-neuroscience/vol/33/suppl/C> )

- 2) Dr. Josh Wiley at Monash University

[Joshuwiley.com](http://Joshuwiley.com)

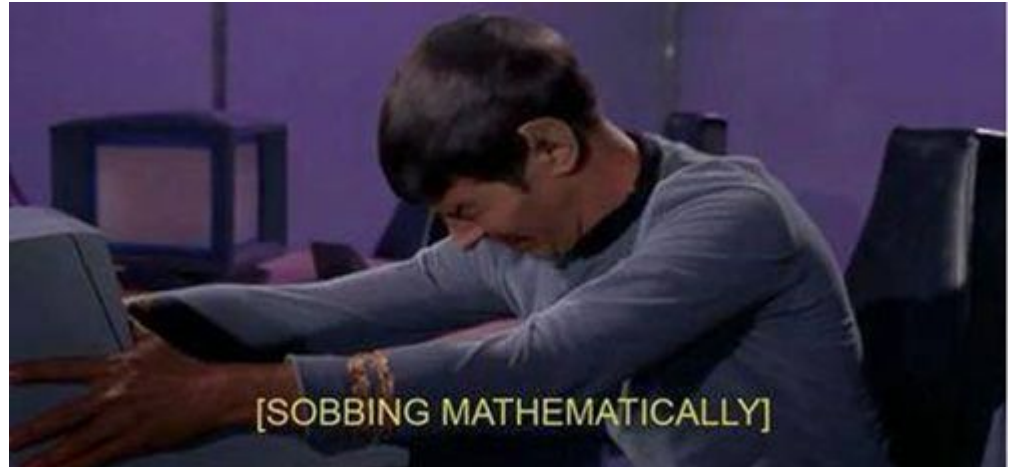
<https://github.com/jwiley> (check out his JWileyMisc R package for a lot of useful functions)

- 3) A perspectives piece yet to be published (stay tuned) on more COVID issues in developmental studies - with Jessica Flannery, Kate McLaughlin, Bonnie Nagel, Jennifer Pfeifer, Eva Telzer, John Flournoy, and Valérie Courchesne

# Missing Data Makes Me Sad

If we only use complete cases (i.e., listwise deletion):

1. Missing data cause a loss of efficiency and makes everyone sad
2. Results from the non-missing data may be biased and that's a waste of time and also sad



# Why it matters

T.H. Matta, et al.

*Developmental Cognitive Neuroscience* 33 (2018) 83–98

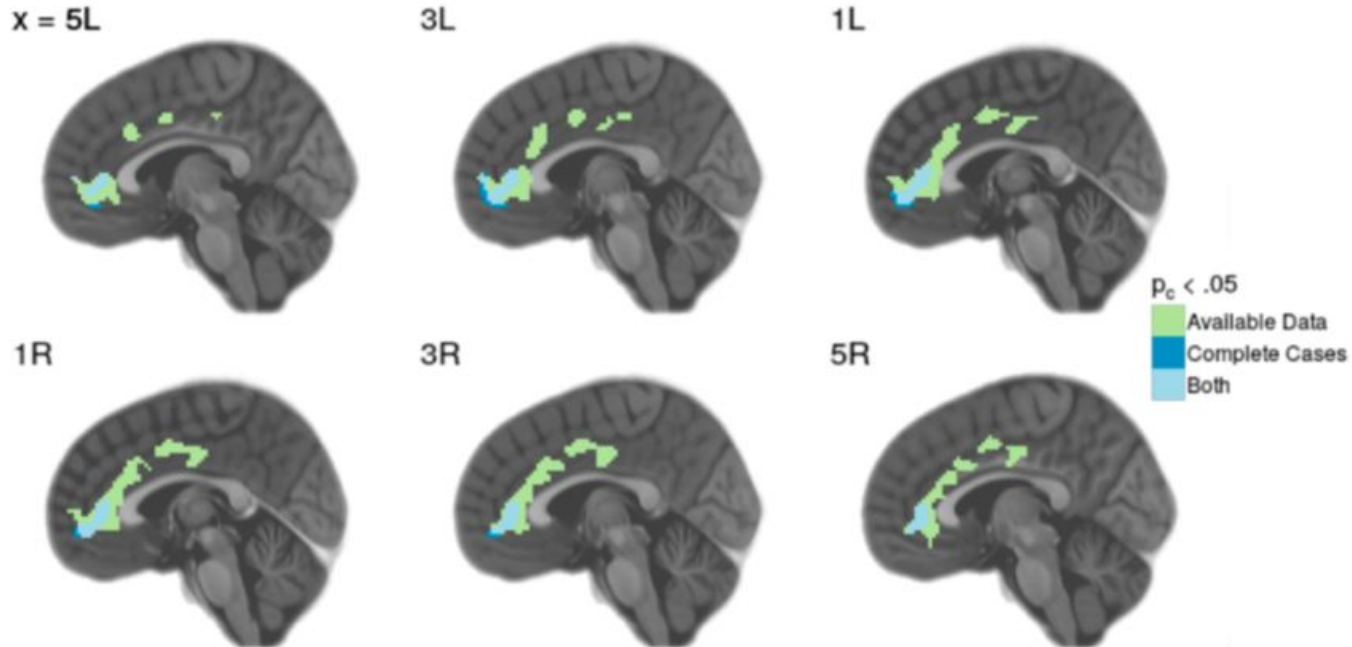


Fig. 3. Significant clusters identified in both available data and complete case analysis is indicated in blue, while significant clusters identified in the available data analysis only are indicated in green. Slice labels indicate the MNI coordinate along which the slice was acquired. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of the article.)

# Some reasons to care about it in ABCD

## Baseline brain function in the preadolescents of the ABCD Study

B. Charani, S. Hahn, [...] the ABCD Consortium

*Nature Neuroscience* (2021)

<https://doi.org/10.1038/s41593-021-00867-9>

**Table 1 | Inclusion criteria in the task fMRI analyses and number of participants remaining after each step of exclusions**

Criteria\task	EN-back	SST	MID
Total number scanned participants	10,189	10,294	10,385
Two runs that passed MRI quality control	8,981	9,035	9,140
Data available excluding Philips scans	8,163	8,140	8,201
Mutual vertex and voxel data availability	7,969	7,288	7,427
Motion censoring (mean FD < 0.9 mm)	7,680	7,000	7,239
d.f. across runs >200	7,680	7,000	7,225
Beta weights outlier detection	6,666	6,995	7,214
Passed behavioral performance QC	6,085	5,116	6,753
No missing covariates	6,009	5,547	6,657

Covariates include age, sex, education, puberty, race, family and scanner ID. FD, framewise displacement; QC, quality control.

Missing data type	Assumptions	Conditions for unbiasedness
<p><b>Missing completely at random (MCAR):</b> missingness is entirely independent of our outcome of interest (<math>y</math>) and any covariates <math>y</math> depends on</p>	<p>missingness is</p> <ul style="list-style-type: none"> <li>(1) not dependent on observations of <math>Y</math> that we don't have</li> <li>(2) not dependent on observations of <math>Y</math> that we do have, and</li> <li>(3) not dependent on covariates <math>X</math>, that <math>Y</math> is dependent on.</li> </ul>	<p>Unbiased results using complete case analysis, or all available data in maximum likelihood, multiple imputation.</p>
<p><b>Covariate-dependent MCAR:</b> missingness is only dependent on the values of variables that affect <math>y</math></p>	<p>(1) and (2); dropping assumption 3.</p>	<p>Unbiased if covariate is included in the models for maximum likelihood analysis or imputation.</p>
<p><b>Missing at random (MAR):</b> missingness (for a particular participant or unit) is independent of the unobserved values of <math>y</math>, i.e., depends only on the values of <math>y</math> (for that unit) that we were able to collect</p>	<p>(1); dropping assumptions 2 &amp; 3.</p>	<p>Unbiased estimates only if all available data are used in a maximum likelihood or multiple imputation framework.</p>
<p><b>Missing not at random (MNAR):</b> missingness is dependent on the unobserved values of <math>y</math>, i.e., the values of the missing data depend on the outcome values that we were not able to collect.</p>	<p>No assumptions made.</p>	<p>Biased estimates (sensitivity analyses should be performed).</p>
<p><b>Take-away:</b> It's best to choose an estimation method that allows you to assume the data are MAR and do sensitivity analyses under the assumption that the data are actually MNAR.</p>		

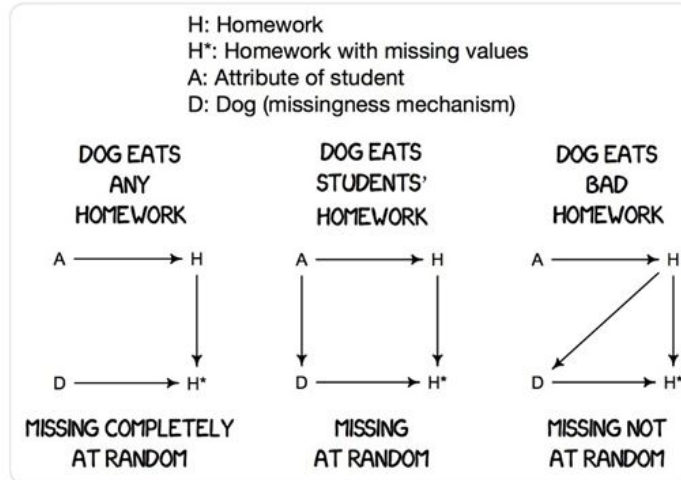
# Missing Data Mechanisms



Richard McElreath  
@rjmcElreath

Follow

In today's lecture, I tried to redefine missing data types (MCAR, MAR, MNAR) as different reasons a dog might eat your homework. This needs more work, but audience seemed to appreciate it.

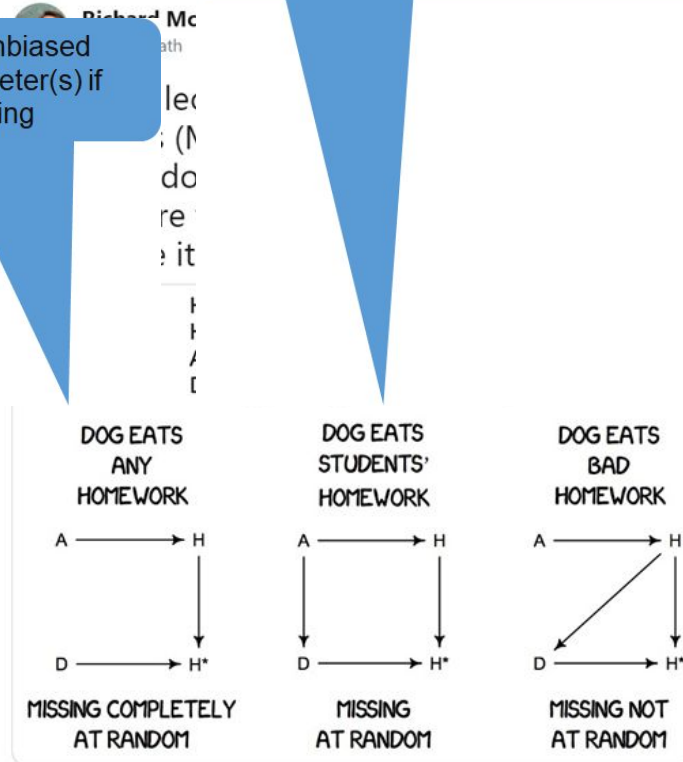


# Missing Data Mechanism

listwise deletion will yield unbiased estimates of the true parameter(s) if the data had not been missing

possible to recover unbiased estimates if the right other variables are present.

cannot recover unbiased estimates





# COVID hot mess

Situation	Missingness mechanism	What do
<p>Beginning of the pandemic</p> <ul style="list-style-type: none"><li>- Research universally shut down</li><li>- Participants missing an entire wave/MRI</li></ul>	<p>Likely <b>MCAR</b> because the temporal order in which participants were assessed should be random with respect to y</p>	<p>missingness depends only on variables that do not affect the outcome; Obtain unbiased estimates using either complete cases or all available data (ML estimation or multiple imputation)</p>
<p>Later in the pandemic</p> <ul style="list-style-type: none"><li>- Some regions (or sites!) opened up later</li><li>- In-person maybe conditional on region</li></ul>	<p>If region was dependent on initial observed anxiety, missingness on an outcome of anxiety would be considered <b>MAR</b></p>	<p>missingness dependent on observed levels of outcome variable; estimate unbiased coefficients using all available data</p>
<ul style="list-style-type: none"><li>- The outcome may be influenced by the missingness itself</li></ul>	<p>If missingness is caused by anxiety specific to the missed observation, the data are <b>MNAR</b></p>	<p>the missed observation of anxiety depends on its unobserved level; obtain biased estimates regardless of method - sensitivity analysis</p>
<ul style="list-style-type: none"><li>- Tasks shifted online that are harder for younger children (neurocog tasks)</li></ul>	<p>Differences in <i>age</i> cause missing an online assessment, and age also causes neurocog performance: <b>covariate- dependent MCAR</b></p>	<p>Obtain unbiased estimates by including age in the model</p>

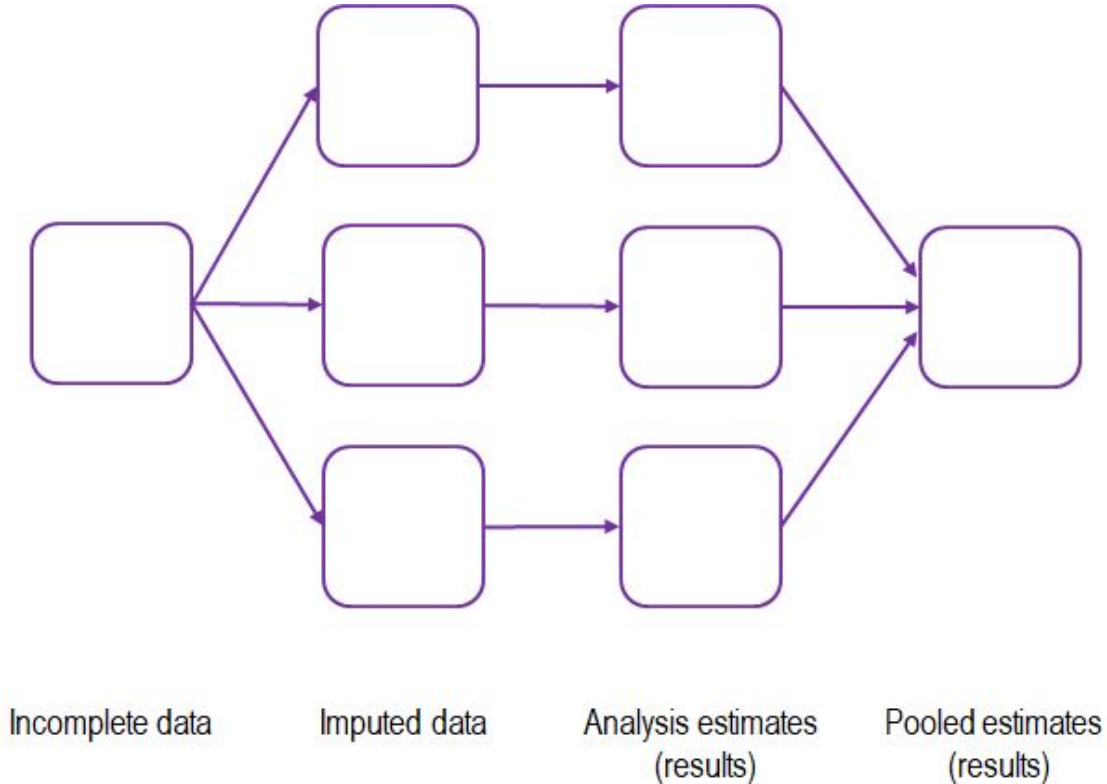
# COVID in ABCD

- Offering assessments virtually (over video chat or phone)
- Eliminated some components of the protocol
- Some sites but not others paused in-person visits completely and have only been conducting assessments remotely.
- “Scheduling flexibility” - so some time points might be “off time”.
- Some MRIs will not be completed until the site determines it is safe to reopen
- COVID-19 supplemental data release - more electronic surveys specific to coping during the pandemic and other situational variables (May, June, Aug, Dec 2020).

# Recommendations:

1. Always include any variables upon which  $y$  theoretically (not empirically) depends.
2. Carefully consider the COVID-related missing data mechanism. Is the missing outcome variable conditionally dependent on some other variable that *did* get collected (MAR or covariate-dependent MCAR), or is it assumed that the value of the outcome variable (if we did know it) is the reason for the missingness (MNAR)?
3. It is impossible to empirically rule out MNAR. Sensitivity analyses can be performed to assess the effect of missing observations had they been collected. Tutorials: [Coertjens et al., 2017](#); [Leurent et al., 2018](#); [Resseguier et al., 2011](#).
4. Always conduct analyses with all available data using maximum likelihood or multiple imputation methods. In rare cases, ML methods may not be available in which case multiple-imputation may be helpful.

# A very crash course in Multiple Imputation (MI)

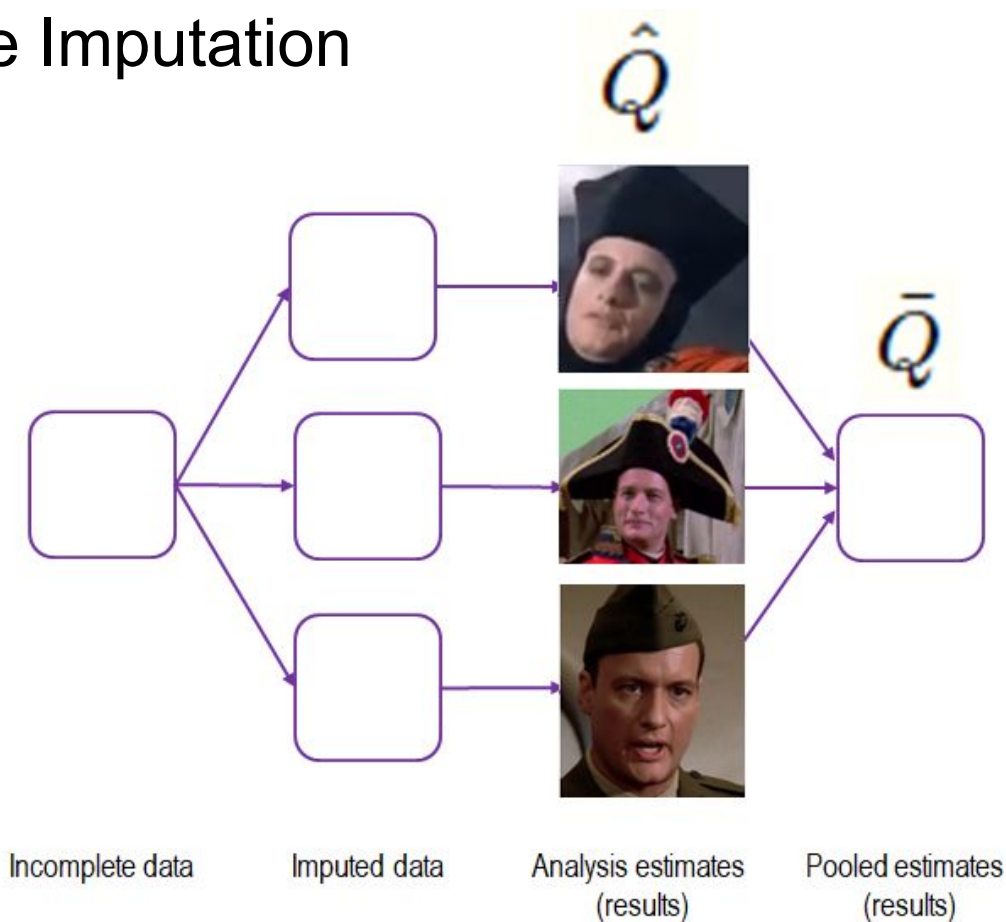


# Multiple Imputation



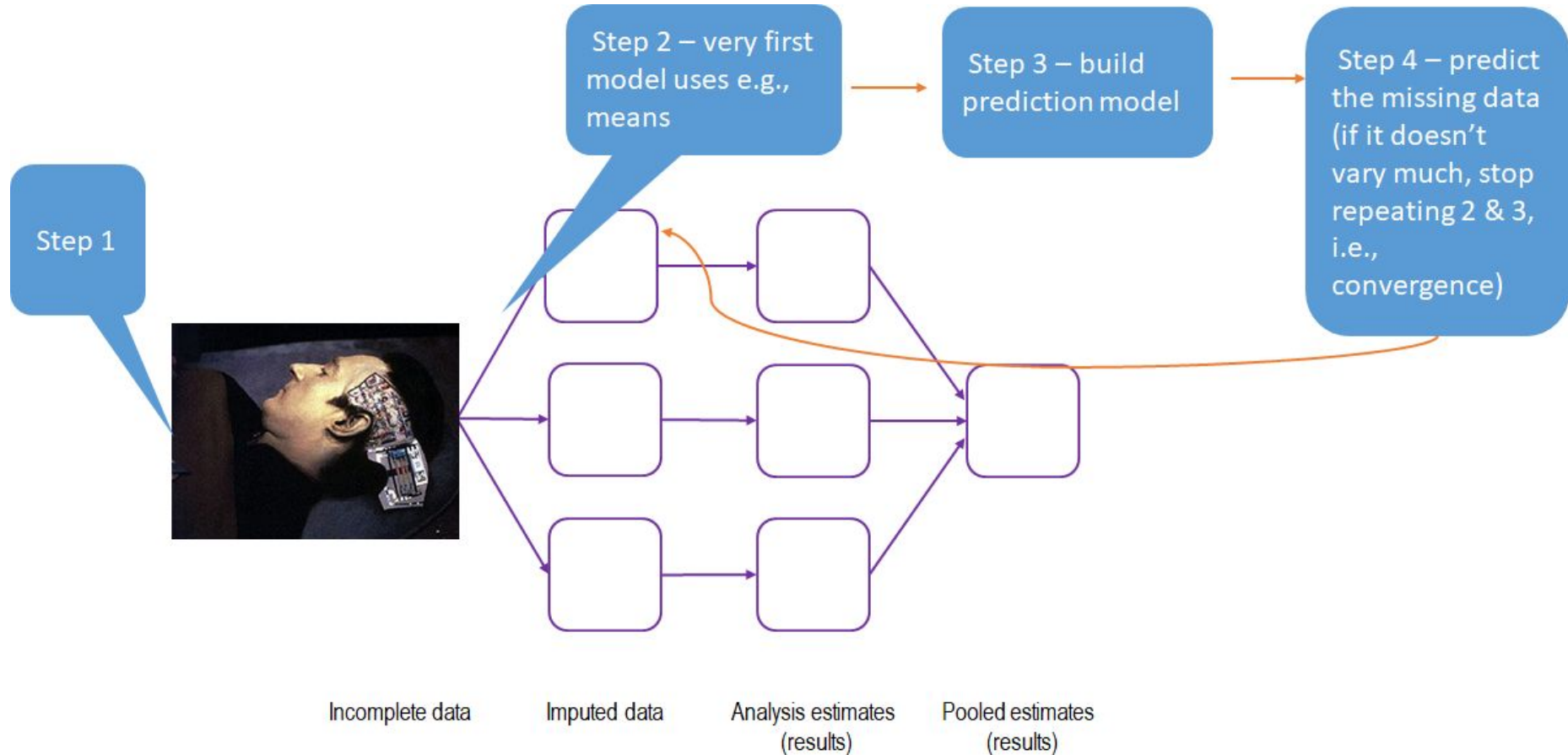
- Let  $Q$  be some population value (e.g., a mean, a regression coefficient).
- Let  $\hat{Q}$  be an estimate of  $Q$  with some estimate of uncertainty due to sampling variation, calculated typically in each imputed dataset.
- Let  $\bar{Q}$  be the average of a set of estimates,  $\hat{Q}$  across different imputed datasets, with some estimate of uncertainty both due to sampling variation impacting  $\hat{Q}$  and missing data uncertainty (causing variation in  $\hat{Q}$  from one imputed dataset to the next).

# Multiple Imputation

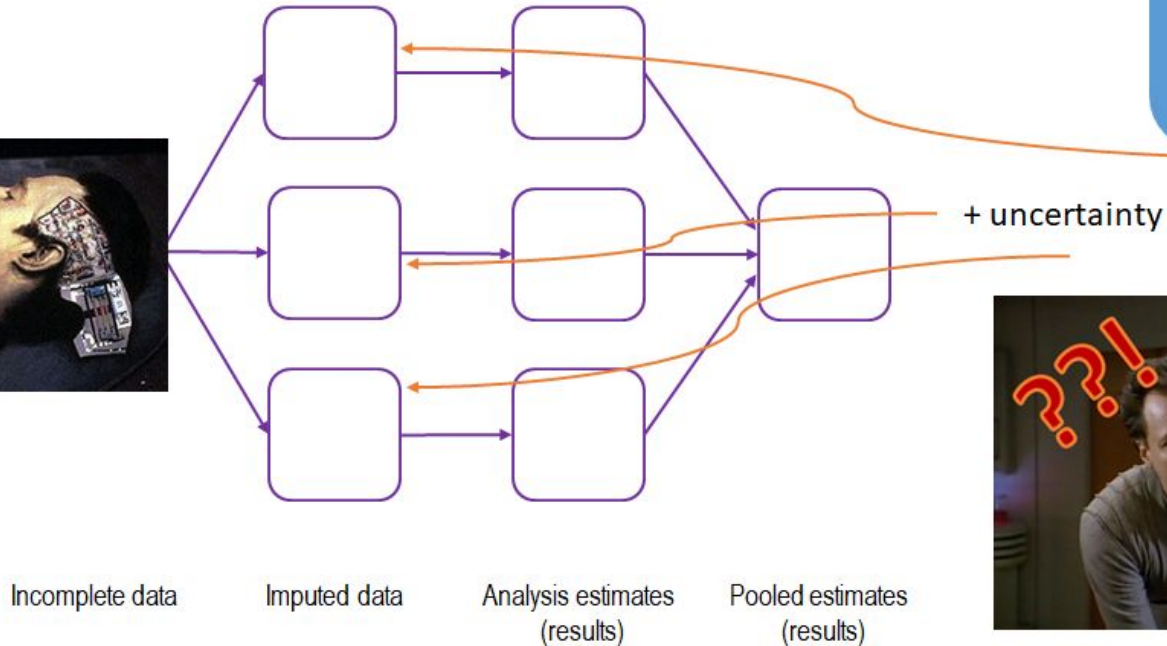
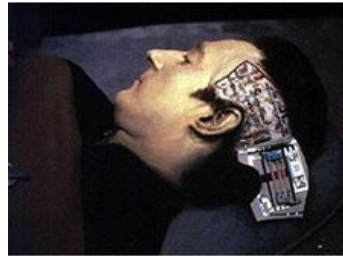


$$\bar{Q} = \frac{1}{m} \sum_{i=1}^m \hat{Q}_i$$

# Multiple imputation STEPS



# Multiple imputation STEPS



Step 4 – predict the missing data (if it doesn't vary much, stop repeating 2 & 3, i.e., convergence)





# Neat code from Josh Wiley to show how imputation works

[https://github.com/michellebyrne1/MonashHonoursStatistics/blob/8e16b6f141340869bb99296aa25940dafa6cdd7a/MissingData\\_Michelle.Rmd#L233](https://github.com/michellebyrne1/MonashHonoursStatistics/blob/8e16b6f141340869bb99296aa25940dafa6cdd7a/MissingData_Michelle.Rmd#L233)

