Carl von Ossietzky
Universität
Oldenburg

meta
rep

# Multiverse Analysis: Theory and Conceptualisation

Cassie Short, Daniel Kristanto, Micha Burkhardt, Andrea Hildebrandt

Pre-Conference-Workshop, PuG 2025

17.-18.06.2025

# The schedule for the afternoon

- 13:30-15:00: Multiverse Analysis Theory and Conceptualisation
- 15:00-15:30: Coffee break
- 15:30-17:00: Multiverse Analysis in Neuroimaging
- 17:00-18:00: Conceptualisation activity

# Content

Why Multiverse Analysis

- The replication problem
- The five sources of uncertainty in empirical research
- Traditional empirical research approach
- On the multiplicity of analysis strategies
- Multiverse analysis
- Principled multiverse analysis
- Systematic multiverse analysis construction
- Summary and conclusions

# The replication problem

# Three main reasons for the replication problem

- Fragmented theoretical landscape and bold hypotheses
- Modest sample size resulting in low statistical power
- Undisclosed flexibility in data analyses

e.g., Paul, Short et al. (2022)

# Three main reasons for the replication problem

- Fragmented theoretical landscape and bold hypotheses
- Modest sample size resulting in low statistical power
- Undisclosed flexibility in data analyses

e.g., Paul, Short et al. (2022)

# Collaborative forking path analyses (cFPA, Wacker, 2017)



Paul, Short et al. (2022)

# Content

1. **Theoretical input**
   - The replication problem
   - The five sources of uncertainty in empirical research
   - Traditional empirical research approach
   - On the multiplicity of analysis strategies
   - Multiverse analysis
   - Principled multiverse analysis
   - Systematic multiverse analysis construction
   - Summary and conclusions

# Introduction - The goal of quantitative empirical research

- The aim of most empirical research is to test the (causal) relationship between an outcome ($Y$) and explanatory variables ($X_1$, $X_2$,..., $X_n$)
- Estimate a function to describe bi- or multivariate relationships

# Introduction - The goal of quantitative empirical research

- The aim of most empirical research is to test the (causal) relationship between an outcome ($Y$) and explanatory variables ($X_1$, $X_2$,..., $X_n$)
- Estimate a function to describe bi- or multivariate relationships
- In many cases in psychology we model multiple outcomes simultaneously, but for simplicity in this workshop we will only talk about models with a single outcome variable.

# Introduction - The goal of quantitative empirical research

- The aim of most empirical research is to test the (causal) relationship between an outcome ($Y$) and explanatory variables ($X_1$, $X_2$,..., $X_n$)
- Estimate a function to describe bi- or multivariate relationships

Mathematical function for relating an outcome variable to explanatory variable(s) and confounder(s)

$$Y = f(X_{1-n}, Z_{1-n}) \tag{1}$$

# Sources of uncertainty in empirical research - Introduction

- The aim of most empirical research is to test the (causal) relationship between an outcome ($Y$) and explanatory variables ($X_1$, $X_2$,..., $X_n$)
- Estimate a function to describe bi- or multivariate relationships

Mathematical function for relating an outcome variable to explanatory variable(s) and confounder(s)

$$Y = f(X_{1-n}, Z_{1-n}) \tag{2}$$

- Let's discuss an example

# Example

- Personality psychology
  - Is there an association between educational success ($ES$) and the big five personality traits ($O$, $C$, $E$, $A$, $N$) ?

# Example

- Personality psychology
  - Is there an association between educational success ($ES$) and the big five personality traits ($O$, $C$, $E$, $A$, $N$) ?

## Mathematical function

$$ES = f(O, C, E, A, N) \tag{3}$$

# Example

- Personality psychology
  - Is there an association between educational success ($ES$) and the big five personality traits ($O, C, E, A, N$) ?

- The variables $ES, O, C, E, A, N$ have to be operationalised
  1. Is $ES$ a binary success variable or a quantitative variable?
  2. Are the variables $O, C, E, A, N$ observed or latent?
- Are there confounders to be considered? For example motivation.
- The functional form of the association and the statistical method to estimate the model parameters need be specified, and these decisions will depend on the type of variables, on sample size, on sample characteristics, etc.

# The five sources of uncertainty in empirical research - see Hoffmann et al. (2021)

- Measurement uncertainty arises from randomness in the operationalisation or measurement of input and output variables
- Data processing uncertainty arises from the multiplicity of choices in selecting the data to be analysed and in defining, cleaning and transforming the input and output variables
- Model uncertainty arises from the multiplicity of choices in the specification of the model structure to describe the phenomenon of interest
- Method uncertainty arises from the multiplicity of potential decisions in the choice of statistical methods to estimate model parameters
- Sampling uncertainty arises from randomness in the sampling from the population of interest

# The five sources of uncertainty in empirical research - see Hoffmann et al. (2021)

- **Measurement uncertainty** arises from randomness in the operationalisation or measurement of input and output variables
- **Data processing uncertainty** arises from the multiplicity of choices in selecting the data to be analysed and in defining, cleaning and transforming the input and output variables
- **Model uncertainty** arises from the multiplicity of choices in the specification of the model structure to describe the phenomenon of interest
- **Method uncertainty** arises from the multiplicity of potential decisions in the choice of statistical methods to estimate model parameters
- **Sampling uncertainty** arises from randomness in the sampling from the population of interest

Note that in frequentist statistics, sampling uncertainty is the only one routinely taken into account, e.g., in null hypothesis significance testing.

# This workshop

...is concerned with

- measurement uncertainty
- data processing uncertainty
- model uncertainty
- method uncertainty

# Back to our example

- Personality psychology
  - Is there an association between educational success (*ES*) and the big five personality traits (*O*, *C*, *E*, *A*, *N*) ?

If *ES* is taken as quantitative and the personality traits are operationalised and quantified as sum scores on the NEO-FFI questionnaire, the sample is representative of the population and the predictors are all normally distributed...

- Which model could be estimated to address the above question?
- How could the model parameters be estimated?

# Back to our example

- Personality psychology
  - ▸ Is there an association between educational success (*ES*) and the big five personality traits (*O*, *C*, *E*, *A*, *N*) ?

If *ES* is taken as quantitative and the personality traits are operationalised and quantified as sum scores on the NEO-FFI questionnaire, the sample is representative of the population and the predictors are all normally distributed:

---

**Mathematical function (option 1)**

$$ES = b_0 + b_1 \cdot O + b_2 \cdot C + b_3 \cdot E + b_4 \cdot A + b_5 \cdot N + \epsilon \tag{4}$$

---

The model parameters can be estimated with the Ordinary Least Squares (OLS) method or the Maximum Likelihood Estimator (MLE).

# Back to our example

- Personality psychology
  - ▸ Is there an association between educational success (*ES*) and the big five personality traits (*O*, *C*, *E*, *A*, *N*) ?

If *ES* is taken as binary and the personality traits are operationalised and quantified as sum scores on the NEO-FFI questionnaire, the sample is representative of the population and the predictors are all normally distributed:

---

Mathematical function (option 2)

$$P(ES = 1|O, C, E, A, N) = \frac{e^{b_0 + b_1 \cdot O + b_2 \cdot C + b_3 \cdot E + b_4 \cdot A + b_5 \cdot N}}{1 + e^{b_0 + b_1 \cdot O + b_2 \cdot C + b_3 \cdot E + b_4 \cdot A + b_5 \cdot N}} \quad (5)$$

---

The model parameters can be estimated with the Maximum Likelihood Estimator (MLE) or alternatives.

# Content

Theoretical input

- The replication problem
- The five sources of uncertainty in empirical research
- Traditional empirical research approach
- On the multiplicity of analysis strategies
- Multiverse analysis
- Principled multiverse analysis
- Systematic multiverse analysis construction
- Summary and conclusions

# Traditional empirical research approach

### Researcher A

1. Research question: Is there an association between educational success ($ES$) and the big five personality traits ($O$, $C$, $E$, $A$, $N$)?

2. Measurement: Use years of education as outcome and five NEO-FFI scale sum scores as predictors

3. Data processing: Exclude observations with missing data, exclude observations with $M \pm 3SD$ identified as outliers

4. Model: Apply an additive multiple regression model with the link function for the Gaussian distribution

5. Estimation method: Estimate the parameters by using the OLS method

6. Sampling: Null hypothesis statistical significance testing at $\alpha = 0.05$

# Traditional empirical research approach

## Researcher B

1. Research question: Is there an association between educational success ($ES$) and the big five personality traits ($O$, $C$, $E$, $A$, $N$)?

2. Measurement: Use university degree (yes vs. no) as outcome and five NEO-FFI principal component scores as predictors

3. Data processing: Use mean imputation replace missing data on the predictors, exclude observations with $M \pm 2.5SD$ identified as outliers

4. Model: Apply a logistic regression model with the link function for the Binomial distribution

5. Estimation method: Estimate the parameters by using the MLE

6. Sampling: Null hypothesis statistical significance testing at $\alpha = 0.01$

# Content

Theoretical input

- The replication problem
- The five sources of uncertainty in empirical research
- Traditional empirical research approach
- On the multiplicity of analysis strategies
- Multiverse analysis
- Principled multiverse analysis
- Systematic multiverse analysis construction
- Summary and conclusions

# On the multiplicity of analysis strategies (Gelman & Loken, 2014)

- The combination of many possible choices throughout the research process results in a theoretical multiverse of statistical outcomes, known as the *Garden of Forking Paths*.

# The multiplicity of analysis strategies

# Sankey Diagram

- The combination of many possible choices throughout the research process results in a theoretical multiverse of statistical outcomes, known as the *Garden of Forking Paths*.



Magnitude of the multiverse = Cartesian product of all methodological decisions: $2 \cdot 3 \cdot 11 \cdot 8 = 528$ forking paths.

# How to address the multiplicity of analyses strategies?

# How to address the multiplicity of analyses strategies? - Hoffmann et al., 2021

- Reduce uncertainty
    - integrate existing knowledge
    - improve measurements
    - formulate more precise theories
    - increase sample size

# How to address the multiplicity of analyses strategies? - Hoffmann et al., 2021

1. Accept uncertainty
   - aim for multiple lines of evidence
   - conduct replication studies
   - acknowledge constraints on generalizability
   - conduct meta-analysis

# How to address the multiplicity of analyses strategies? - Hoffmann et al., 2021

1. Integrate uncertainty
   - Bayesian model averaging
   - Bayesian deep learning
   - Probabilistic sensitivity analysis

# How to address the multiplicity of analyses strategies? - Hoffmann et al., 2021

- Report uncertainty
    - sensitivity analysis
    - robustness analysis
    - multiverse analysis / vibration of effects / specification curves

# How to address the multiplicity of analysis strategies? - Our focus

1. Reduce uncertainty
2. Accept uncertainty
3. Integrate uncertainty
4. Report uncertainty

# Content

Theoretical input

- The replication problem
- The five sources of uncertainty in empirical research
- Traditional empirical research approach
- On the multiplicity of analysis strategies
- Multiverse analysis
- Principled multiverse analysis
- Systematic multiverse analysis construction
- Summary and conclusions

# Multiverse analysis

Alternative terms for same or very similar approaches: "Vibration of effects" in epidemiology, "specification curve analysis" in psychology, "measure of robustness to misspecification" in economics, "multimodel analysis" and "computational robustness analysis" in sociology

- Primary goal of multiverse analysis is to enhance research transparency and uncover various sources of uncertainties

# Multiverse analysis

- Primary goal of multiverse analysis is to enhance research transparency and uncover various sources of uncertainties
- Multiverse analysis investigates the effect of arbitrary decisions at the level of data processing and statistical modeling given the raw data file, research questions, and hypotheses

# Multiverse analysis

- Primary goal of multiverse analysis is to enhance research transparency and uncover various sources of uncertainties
- Multiverse analysis investigates the effect of arbitrary decisions at the level of data processing and statistical modeling given the raw data file, research questions, and hypotheses
- Steps
  1. Identify decision nodes and arbitrary alternatives concerning preparation of the data and the statistical approach to construct all reasonable combinations of decisions
  2. Perform statistical analysis across all combinations specified in the previous step to obtain a set of results for each combination
  3. Examine variability in results by graphical representations
  4. Potentially: Conduct joint inference across the multiverse of findings

# Back to our example and identify decision nodes

Research question: Is there an association between educational success (*ES*) and the big five personality traits (*O*, *C*, *E*, *A*, *N*)?

- In two groups

# Content

Theoretical input

- The replication problem
- The five sources of uncertainty in empirical research
- Traditional empirical research approach
- On the multiplicity of analysis strategies
- Multiverse analysis
- Principled multiverse analysis
- Systematic multiverse analysis construction
- Summary and conclusions

# Principled multiverse analysis

- Garbage in - garbage out
- Del Giudice & Gangestad (2021) - "The multiverse is a dangerous place"

# Principled multiverse analysis

- Del Giudice & Gangestad (2021) - "The multiverse is a dangerous place"
  - "In principle, multiverse-style analyses can be highly instructive. At the same time, analyses that explore multiverse spaces that are not homogeneous can produce misleading results and interpretations, lead scholars to dismiss the robustness of theoretically important findings that do exist, and discourage them from following fruitful avenues of research. This can hinder scientific progress just as much as the proliferation of false, unreplicable findings does."
  - "The main danger of multiverse-style methods lies in their potential for combinatorial explosion. Just a few decisions incorrectly treated as arbitrary can quickly explode the size of the multiverse, drowning reasonable effect estimates in a sea of unjustified alternatives."

# Principled multiverse analysis

- Del Giudice & Gangestad (2021) - "The multiverse is a dangerous place"
  - "In principle, multiverse-style analyses can be highly instructive. At the same time, analyses that explore multiverse spaces that are not homogeneous can produce misleading results and interpretations, lead scholars to dismiss the robustness of theoretically important findings that do exist, and discourage them from following fruitful avenues of research. This can hinder scientific progress just as much as the proliferation of false, unreplicable findings does."
  - "The main danger of multiverse-style methods lies in their potential for combinatorial explosion. Just a few decisions incorrectly treated as arbitrary can quickly explode the size of the multiverse, drowning reasonable effect estimates in a sea of unjustified alternatives."

A principled multiverse analysis is necessary!

# Principled multiverse analysis (Del Giudice & Gangestad, 2021)

- Multiverse analysis in a principled way is based on the a priori assessment of the equivalence of alternatives at each decision node

# Principled multiverse analysis (Del Giudice & Gangestad, 2021)

- Multiverse analysis in a principled way is based on the a priori assessment of the equivalence of alternatives at each decision node
- Type E decisions - principled equivalence: Specifications are equivalent and effectively arbitrary, e.g., alternatives have comparable validity, examine the same effect, or estimate the effect with comparable precision
- Type N decisions - principled equivalence: Specification are nonequivalent, i.e., some of the alternatives are more justified than others as a means of estimating the effect of interest

# Principled multiverse analysis (Del Giudice & Gangestad, 2021)

- Multiverse analysis in a principled way is based on the a priori assessment of the equivalence of alternatives at each decision node

- Type E decisions - principled equivalence: Specifications are equivalent and effectively arbitrary, e.g., alternatives have comparable validity, examine the same effect, or estimate the effect with comparable precision

- Type N decisions - principled equivalence: Specification are nonequivalent, i.e., some of the alternatives are more justified than others as a means of estimating the effect of interest

- Type E decisions should be selected for multiverse analysis exploring robustness, whereas Type N decisions should not be selected for multiverse analysis

# Equivalence assessment of alternatives

- The assessment of the equivalence of alternatives at each decision node can be based on three kinds of nonequivalence

# Equivalence assessment of alternatives

- The assessment of the equivalence of alternatives at each decision node can be based on three kinds of nonequivalence
  1. Measurement Nonequivalence: Alternative measurement choices yield systematic differences in validity and reliability
     - ★ Example?

# Equivalence assessment of alternatives

- The assessment of the equivalence of alternatives at each decision node can be based on three kinds of nonequivalence
  1. Measurement Nonequivalence: Alternative measurement choices yield systematic differences in validity and reliability
  2. Effect Nonequivalence: Alternative specification investigates a different effect not in line with the effect of interest
     - ★ Example?

# Equivalence assessment of alternatives

- The assessment of the equivalence of alternatives at each decision node can be based on three kinds of nonequivalence

  1. Measurement Nonequivalence: Alternative measurement choices yield systematic differences in validity and reliability
  2. Effect Nonequivalence: Alternative specification investigates a different effect not in line with the effect of interest
  3. Power/Precision Nonequivalence: Alternative specification results in a lower precision of estimating an effect and lower statistical power to detect an effect
     - Example?

# Content

Theoretical input

- The replication problem
- The five sources of uncertainty in empirical research
- Traditional empirical research approach
- On the multiplicity of analysis strategies
- Multiverse analysis
- Principled multiverse analysis
- Systematic multiverse analysis construction
- Summary and conclusions

# Systematic multiverse analysis construction

- **Community guidelines** for transparent, comprehensive and systematic multiverse analysis construction
  - ▸ interdisciplinary guidance on key procedural considerations

- **Systematic Multiverse Analysis Registration Tool (SMART)**
  - ▸ increased transparency, systematicity and reproducibility
  - ▸ reduced uncertainty and potential for QRPs

# Community Guidelines



Multi-curious: A Multi-Disciplinary Guide to Multiverse Analysis

AUTHORS
Cassie Short, Nate Breznau, Maria Bruntsch, Micha Burkhardt, Niko Busch, Elena Cesnaite, Maximilian Frank, Carsten Gießing, Daniel Krähmer, Daniel Kristanto, and 8 more

# SMART App

# Content

Theoretical input

- The replication problem
- The five sources of uncertainty in empirical research
- Traditional empirical research approach
- On the multiplicity of analysis strategies
- Multiverse analysis
- Principled multiverse analysis
- Systematic multiverse analysis construction
- Summary and conclusions

# Summary

- Multiverse analysis is an open science practice for dealing with uncertainties in empirical research
- There are multiple sources of uncertainty (measurement, data processing, model, method, and sampling), but only sampling uncertainty is routinely taken into account
- Multiplicity of analysis strategies resulting in a multiverse of results can be explored using multiverse analysis, where decision nodes and arbitrary alternatives are identified
- Importantly, a priori assessment of the equivalence of alternatives at each decision node is needed to only include Type E (equivalence) decisions in the multiverse analysis
- Note that multiverse analysis aims at assessing robustness to arbitrary changes in data analysis, but is not intended to assess substantive robustness

Uncertainty in research has in most cases an epistemic source, resulting from a lack of knowledge (Hoffmann et al., 2021). Multiverse analysis helps to increase knowledge and reduce uncertainty in the longer term.



https://news.uchicago.edu/

# The schedule for the afternoon

- 13:30-15:00: Multiverse Analysis Theory and Conceptualisation
- 15:00-15:30: Coffee break
- 15:30-17:00: Multiverse Analysis in Neuroimaging
- 17:00-18:00: Conceptualisation activity

# The schedule for the afternoon

- 13:30-15:00: Multiverse Analysis Theory and Conceptualisation
- 15:00-15:30: Coffee break
- 15:30-17:00: Multiverse Analysis in Neuroimaging
- 17:00-18:00: Conceptualisation activity

# Conceptualisation Activity

- EEG Multiverse Analysis
  - Construct a multiverse analysis for one of your EEG projects, or
  - Construct a multiverse analysis for our example (next slides)
- fMRI Multiverse Analysis
  - Construct a multiverse analysis for one of your fMRI projects, or
  - Construct a multiverse analysis for our example (next slides)

**We will discuss your garden of forking paths together.**

# Conceptualisation Activity

## EEG multiverse anlaysis example

**Model:**

Extraversion ~ *f*(happiness LPP - neutral, anger LPP - neutral, fear LPP - neutral, surprise LPP - neutral, sadness LPP - neutral, disgust LPP – neutral)

---

**Data:**

**Sample:** 98 healthy adults ($M_{age}$ = 26.64, $SD_{age}$ = 4.82)

**Extraversion:** NEO Personality Inventory Revised (*M* = 2.26, *SD* = 0.43)

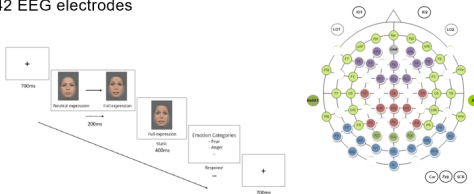**Emotion recognition task with EEG recording:**

6 dynamic emotional expressions

- *happiness, sadness, anger, fear, disgust, surprise*

1 dynamic neutral expression

- *either chewing or blinking*

42 EEG electrodes

# References

Botvinik-Nezer et al. (2020). Variability in the analysis of a single neuroimaging dataset by many teams. *Nature, 582(7810)*:84-88

Del Giudice, M., Gangestad, S. W. (2021). A traveler's guide to the multiverse: Promises, pitfalls, and a framework for the evaluation of analytic decisions. *Advances in Methods and Practices in Psychological Science, 4*.

Gelman, A., & Loken, E. (2014). The Statistical Crisis in Science. *American Scientist, 102*, 460-465.

Hoffmann, S., Schönbrodt, F., Elsas, R., Wilson, R., Strasser, U., & Boulesteix, A.-L. (2021). The multiplicity of analysis strategies jeopardizes replicability: Lessons learned across disciplines. *Royal Society Open Science, 8*, 4.

Paul, K., Short, C. A., Beauducel, A., Per Carsten H., Härpfer, K., Hennig, J., Hewig, J., Hildebrandt, A., Kührt, C. Mueller, E.M., Osinsky, R. Porth, E., Riesel, A., Rodrigues, J., Scheffel, C., Stahl, J., Strobel, A., & Wacker, J. (2022). The Methodology and Dataset of the CoScience EEG-Personality Project – A Large-Scale, Multi-Laboratory Project Grounded in Cooperative Forking Paths Analysis. *Personality Science, 3*, e7177.

Short, C. A., Breznau, N., Bruntsch, M., Burkhardt, M., Busch, N., Cesnaite, E., Frank, M., Gießing, C., Krähmer, D., Kristanto, D., Lonsdorf, T., Neuendorf, C., Nguyen, H., Rausch, M., Schmalz, X., Schneck, A., Tabakci, C., Hildebrandt, A. (2025). Multi-curious: A Multi-Disciplinary Guide to Multiverse Analysis. *MetaArXiV*.

Short, C. A., Inceler, Y. C., Frank, M., Hildebrandt, A. (2025). The Systematic Multiverse Analysis Registration Tool (SMART) for Defining Multiverse Analyses. *MetaArXiV*.

Simonsohn, U., Simmons, J. P., & Nelson, L. D. (2020). Specification curve analysis. *Nature Human Behavior, 4*, 1208-1214.

Wacker, J. (2017). Increasing the reproducibility of science through close cooperation and forking path analysis. *Frontiers in Psychology, 8*, 1332.

# Don't get lost in the Garden of Forking Paths



Think carefully about the analytical choices you can make!

Thank you for attending this workshop!

andrea.hildebrant@uol.de
cassie.short@uol.de
daniel.kristanto@uol.de
micha.burkhardt@uol.de

# Day 2: EEG Multiverse Practical

**Model:**

Extraversion ~ f(happiness LPP - neutral, anger LPP - neutral, fear LPP - neutral, surprise LPP - neutral, sadness LPP - neutral, disgust LPP – neutral)

---

**Data:**

**Sample:** 98 healthy adults ($M_{age}$ = 26.64, $SD_{age}$ = 4.82)

**Extraversion:** NEO Personality Inventory Revised ($M$ = 2.26, $SD$ = 0.43)

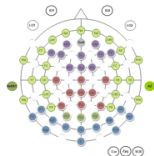**Emotion recognition task with EEG recording:**

6 dynamic emotional expressions

- *happiness, sadness, anger, fear, disgust, surprise*

1 dynamic neutral expression

- *either chewing or blinking*

42 EEG electrodes



**Multiverse Use Case:**

**Baseline correction (2)**

- -100m
- -200ms

**Reference (2)**

- Common average reference
- Linked mastoids

**Time Window for LPP quantification (4)**

- 400 – 600ms
- 500 – 700ms
- 450 – 750ms
- +/- 200 ms around subject average peak

**Electrode cluster for LPP quantification (4)**

- P3, P4, CP1, CP2
- P3, Pz, P4
- CP1, CP2
- Pz

Cartesian product: 2 * 2 * 4 * 4 = **64 pipelines**