

Arbitrary choices, arbitrary results: A multiverse analysis of L2 reaction time data

Ryo Maie

Michigan State University

American Association for Applied Linguistics 2021

Warm-up

Suppose you have a research question:

Do L2 English speakers process lexical items more slowly than L1 speakers?

- You set up a lexical decision task where participants judged whether the word presented on a screen was a word or a nonword.
- You collected the response time as the dependent variable.

Warm-up

Prior to statistical analyses, you have to make the following decisions:

- ① Should all responses be included or only those that are correct?
 - **all or only correct**
- ② Should participants be excluded if they have an unacceptably low accuracy rate. What criterion must they reach?
 - **80%, 85%, or 90%**

Warm-up

Prior to statistical analyses, you have to make the following decisions:

- ③ Should data points be excluded if they are abnormally faster or slower than normally expected? At what point should you consider them to be unacceptable?
 - **200ms, 250ms, or 300ms, for the lower bound**
 - **1500ms, 2000ms, or 2500ms for the upper bound**
- ④ Should data points be excluded if they deviate from the overall trend of individual participants? What range should you adopt as an acceptable range of variation?
 - **2SD, 2.5SD, or 3SD**

Flexibility in data processing

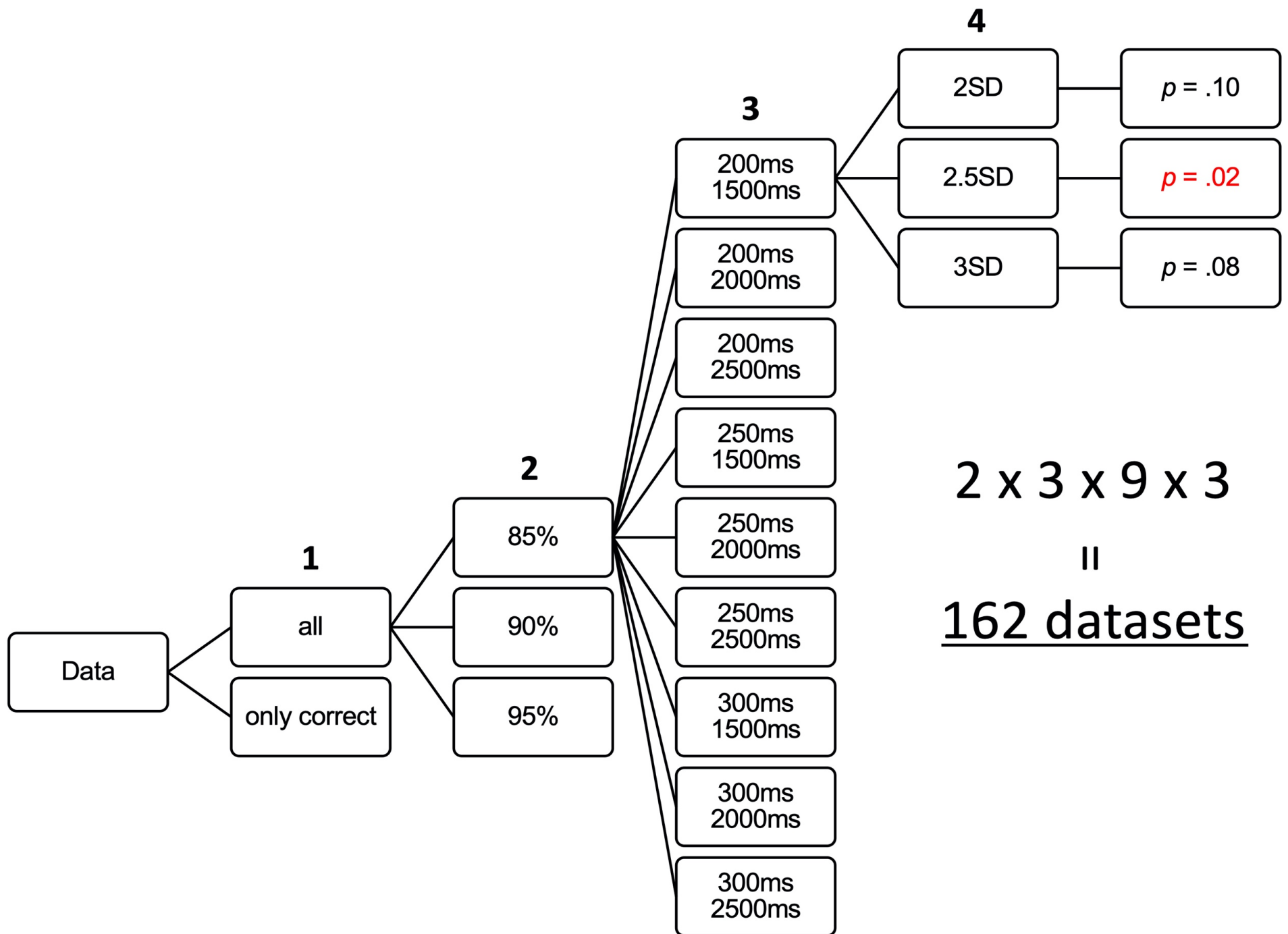
Out of these four decisions, how many were you able to justify with any substantive reasoning?

- Some may be empirically justifiable: word recognition < 300ms is unlikely
- Others are equally (un)justifiable and you can choose any of them
- **There is always an article to cite to defend any of your decisions!!!**

Flexibility in data processing

How do you decide on:

- 2SD, 2.5SD, or 3SD for outliers?
- 1500ms, 2000ms, or 3000ms for unusual responses?



What flexibility causes

- Flexibility in data processing -> Researcher degree of freedom
- When exploited, researcher *df* can lead to:
 - arbitrary variation in statistical results between studies
 - questionable research practices (such as *p*-hacking)
 - What if a researcher goes through many possibilities and only report the one he likes?

Simmons et al. (2011) demonstrated that:

- ① even **one** binary decision inflates the rate of false-positive discoveries to the point that is unacceptable (i.e., larger than 0.05)
- ② with four degrees of freedom, it can become as high as an astonishing rate of **61%**

A series of Earths are arranged in a curved line across the top half of the image, starting from the left and curving towards the right. The Earths increase in size as they move along the curve. The background is a deep purple and black space filled with numerous white stars of varying brightness. The text "Somewhere in the Multiverse...." is written in a bold, white, sans-serif font in the lower-left quadrant of the image.

**Somewhere in the
Multiverse....**

Multiverse analysis

Steegen, Tuerlinckx, Gelman and Vanpaemel (2016)

- Data used in a statistical analysis are usually not person-free. They are “to a certain extent **actively constructed**” (p. 702)
- An analyst constructs a given dataset from raw data, but a single analysis on one dataset is not sufficient because it is the only one of many equally reasonable possibilities, that is, a many worlds or **multiverse** of datasets.
- The concept of a multiverse of datasets is critical because the arbitrariness in reaching a particular dataset is inevitably inherited by the statistical result, and thus:
 - “the data multiverse directly implies a multiverse of statistical results”

Multiverse analysis

In multiverse analysis, a researcher:

- ① identifies every possible method of data processing (e.g., coding, cleaning, and transforming data) that are equally justifiable
- ② performs the same set of analysis across the whole data multiverse
- There is also, however, **a multiverse of statistical models** that can interact with and multiply the multiverse of datasets (we do not consider it here)

A multiverse analysis of L2 RT data: Study

Maie and DeKeyser (2020)

- Adult L1 speakers of English learned a semiartificial language called, *Japlish* (i.e., English lexicon and Japanese syntax and case-marking system), from incidental exposure
- The researchers investigated whether the exposure results in explicit and/or implicit knowledge
- Three structures:
 - Simple: O-S-V and O-S-I-V
 - Complex: O-S-[S-V]-V and O-S-[S-I-V]-V
 - Case: *-ga*, *-o*, and *-ni*

A multiverse analysis of L2 RT data: Analysis

- Variables to consider
 - ① **Complexity**: Simple (-1), Complex (0 or 1), and Case (0 or 1)
 - Effect coded
 - ② **Grammaticality**: Grammatical (0) and Ungrammatical (1)
 - Dummy coded
- The original study used a combination of (M)ANOVAs
- For the current demo, a mixed-effects model:
 - $RT \sim \text{Complex} * \text{Grammaticality} + \text{Case} * \text{Grammaticality} + (1 | \text{Subject}) + (1 | \text{Item})$
- Here, we only consider the regression coefficient, $\text{Case} : \text{Grammaticality}$, which estimates a difference between the grammatical and ungrammatical sentences on Case
 - The original study found an effect here!

Creating data multiverse

Literature Review on:

- ① Previous studies that used a word-monitoring task as a measure of implicit knowledge (e.g., Godfroid 2016; Granena, 2012; Suzuki, 2015)
- ② An often-cited book that discussed technical aspects of conducting L2 RT research (Jiang, 2012)
- I identified **five steps** in processing raw RT data (from a WMT) to construct datasets

Creating data multiverse

- ① Excluding participants who do not reach a certain criterion level of accuracy in comprehension questions
 - 63% (Jiang, 2004)
 - 70% (Jiang, 2012)
 - 75% (Granena, 2012; Suzuki, 2015)

Creating data multiverse

- ② Excluding outlier data points that are larger or smaller than some absolute cutoff points
 - Lower bound
 - 100ms (Suzuki, 2015)
 - 120ms (Jiang, Novokshanova, Masuda, & Wang, 2011)
 - Upper bound
 - 1500ms (Jiang, 2012)
 - 2000ms (Jiang et al, 2011)
 - 2500ms (Suzuki, 2015)

Creating data multiverse

- ③ Excluding data points that are outside of the normal range of a given participant
 - 2SD (Jiang, 2012)
 - 2.5SD (Godfroid, 2016)
 - 3SD (Granena, 2012; Suzuki, 2015)

Creating data multiverse

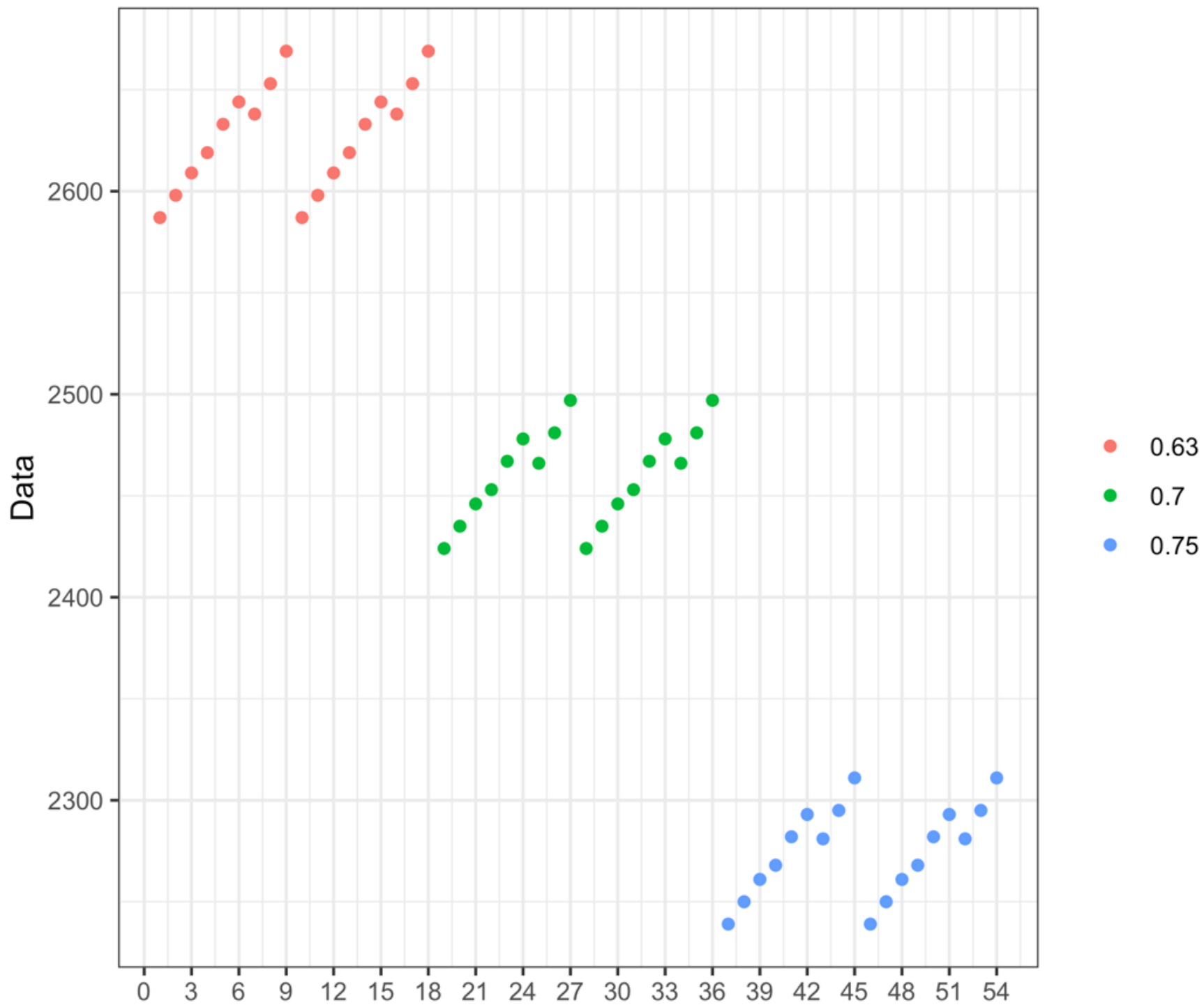
- ④ Excluding participants whose mean word monitoring latency is larger or smaller than the range of the whole group
 - 2SD
 - 2.5SD
 - 3SD

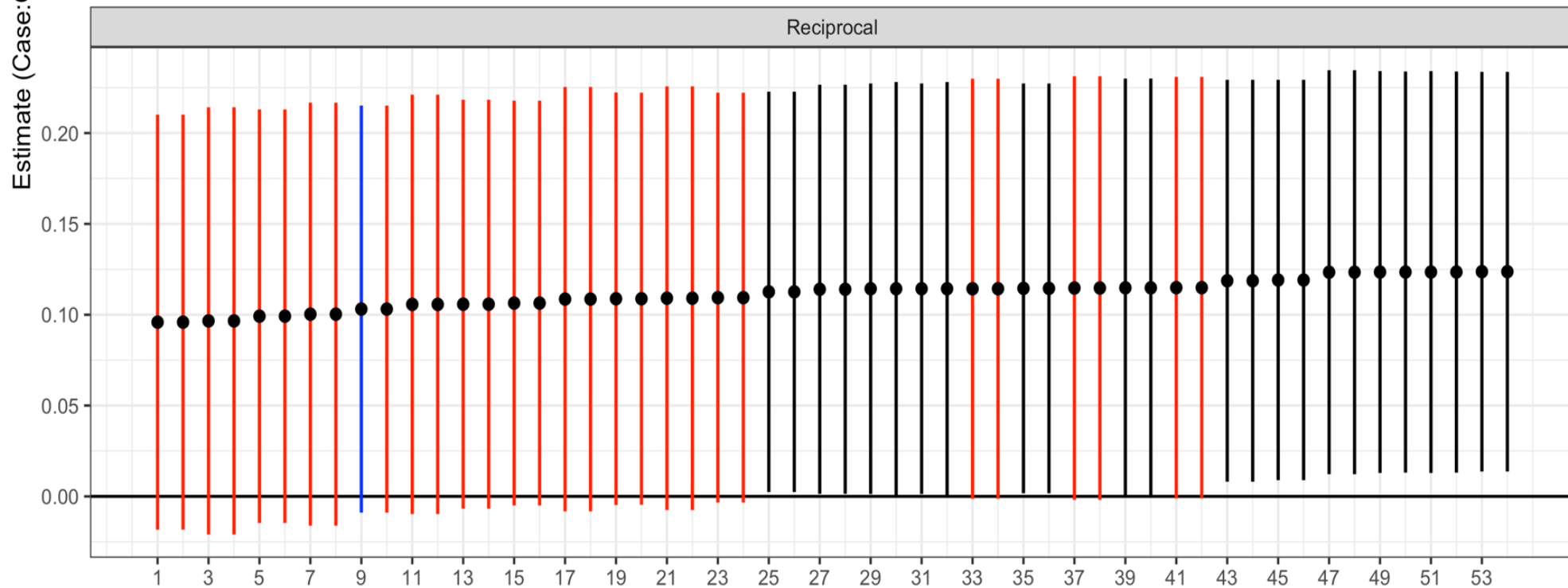
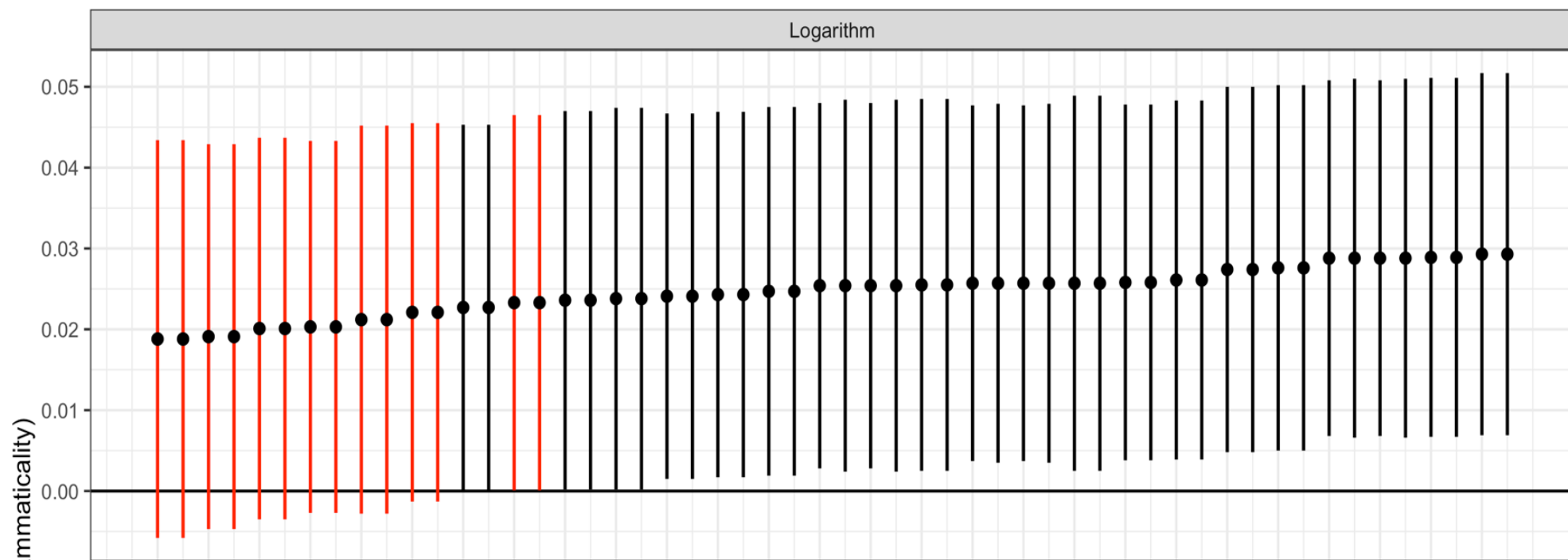
Creating data multiverse

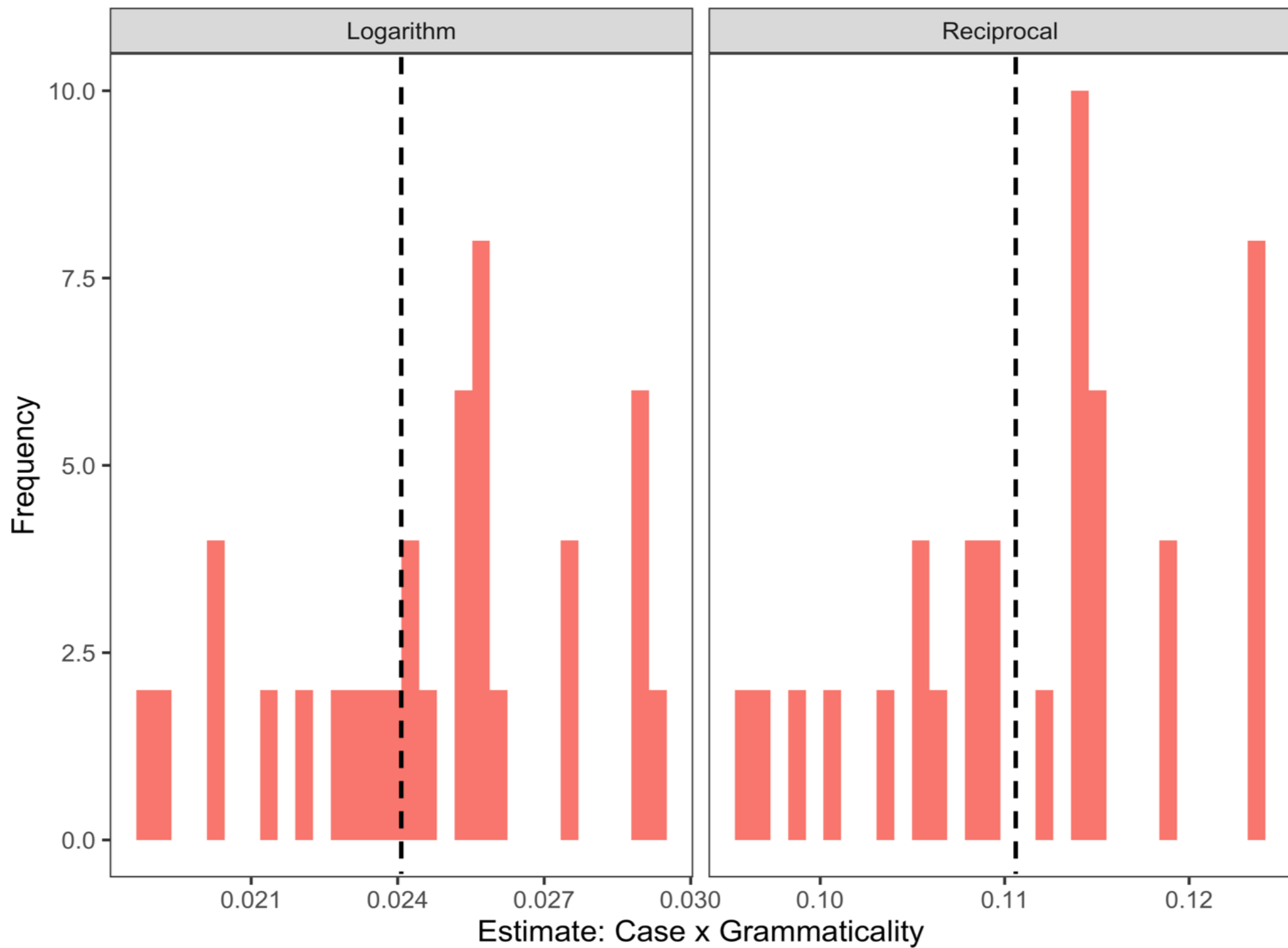
- ⑤ Transforming data for statistical analysis
 - Logarithm (with base = 10)
 - Reciprocal

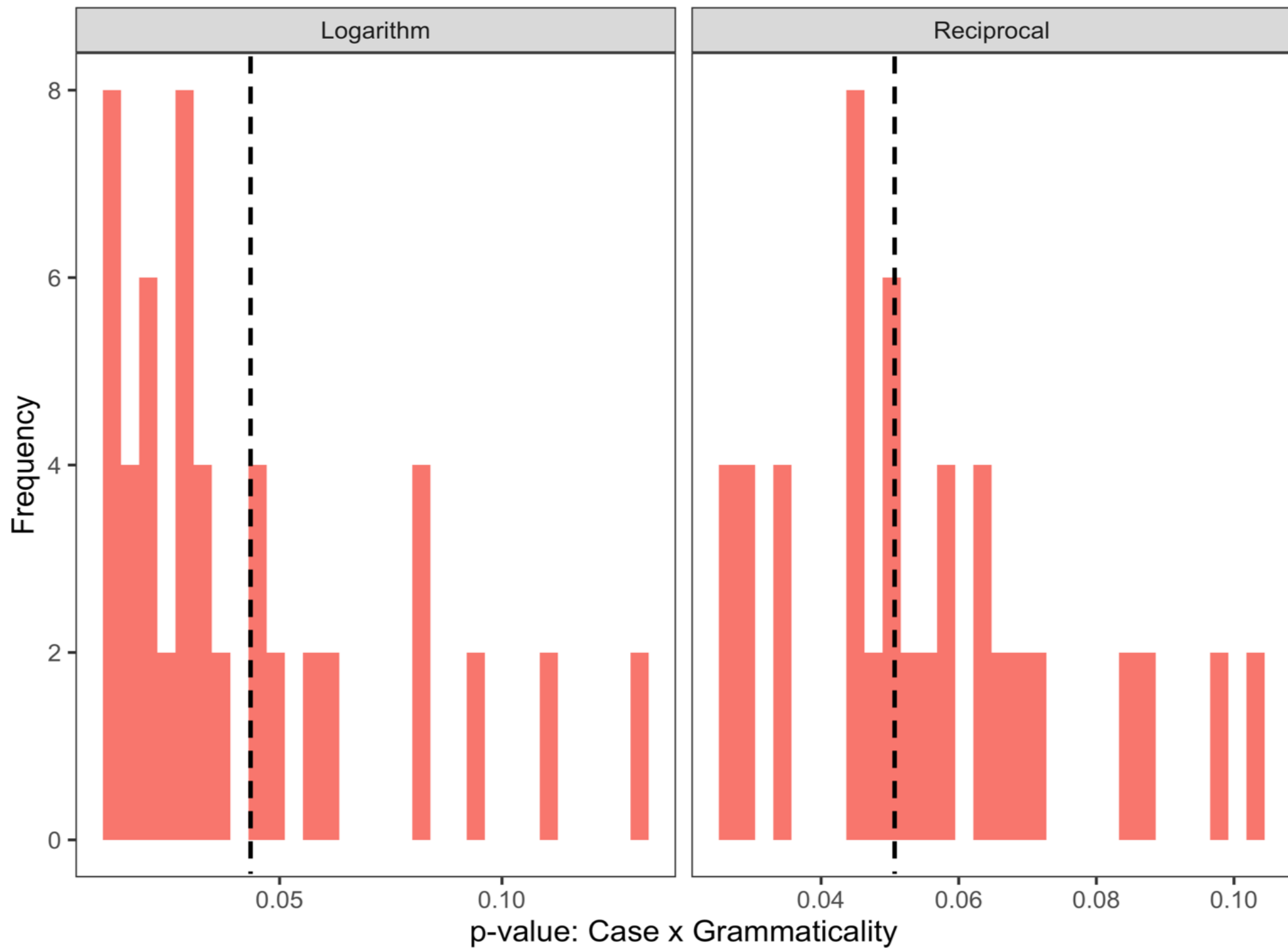
Creating data multiverse

- ① Excluding participants not reaching a critical level of accuracy in comprehension
- ② Excluding outlier data points larger or smaller than absolute cutoff points
- ③ Excluding data points outside of the normal range of a participant
- ④ Excluding participants with mean word monitoring RT larger or smaller than the group
- ⑤ Transforming data for statistical analysis
 - **$3 \times 2 \times 3 \times 3 \times 2 = 108$ datasets**









Step 4

AA [100, 1500] BA [120, 1500]
 AB [100, 2000] BB [120, 2000]
 AC [100, 2500] BC [120, 2500]

Step 2

Log

Rec

A	A	A	B	A	C	B	A	B	B	B	C
0.010	0.010	0.010	0.010	0.034	0.010	0.010	0.010	0.010	0.010	0.034	0.034
0.011	0.017	0.017	0.017	0.057	0.011	0.017	0.017	0.017	0.017	0.057	0.057
0.011	0.017	0.017	0.017	0.080	0.011	0.017	0.017	0.017	0.017	0.080	0.080

0.021	0.022	0.022	0.022	0.047	0.021	0.022	0.022	0.022	0.022	0.047	0.047
0.023	0.036	0.036	0.036	0.081	0.023	0.036	0.036	0.036	0.036	0.081	0.081
0.022	0.035	0.035	0.035	0.112	0.022	0.035	0.035	0.035	0.035	0.112	0.112

0.028	0.029	0.029	0.029	0.062	0.028	0.029	0.029	0.029	0.029	0.062	0.062
0.030	0.047	0.047	0.047	0.093	0.030	0.047	0.047	0.047	0.047	0.093	0.093
0.029	0.046	0.046	0.046	0.129	0.029	0.046	0.046	0.046	0.046	0.129	0.129

A	A	A	B	A	C	B	A	B	B	B	C
0.028	0.027	0.027	0.027	0.044	0.028	0.027	0.027	0.027	0.027	0.044	0.044
0.029	0.034	0.034	0.034	0.059	0.029	0.034	0.034	0.034	0.034	0.059	0.059
0.028	0.033	0.033	0.033	<u>0.069</u>	0.028	0.033	0.033	0.033	0.033	0.069	0.069

0.044	0.046	0.046	0.046	0.046	0.046	0.046	0.046	0.046	0.046	0.063	0.063
0.059	0.047	0.047	0.047	0.056	0.047	0.058	0.058	0.058	0.058	0.084	0.084
0.069	0.045	0.045	0.045	0.063	0.045	0.056	0.056	0.056	0.056	0.097	0.097

0.049	0.051	0.051	0.051	0.070	0.049	0.051	0.051	0.051	0.051	0.070	0.070
0.052	0.066	0.066	0.066	0.089	0.052	0.066	0.066	0.066	0.066	0.089	0.089
0.051	0.064	0.064	0.064	0.103	0.051	0.064	0.064	0.064	0.064	0.103	0.103

2SD

2.5SD

3SD

Step 3

63%

2SD

2.5SD

3SD

70%

2SD

2.5SD

3SD

75%

Step 1

Discussion: What to do?

- Variation in statistical results due to arbitrary choices during data processing and **you can reach different conclusions depending on which datasets you analyzed!**
- This cannot be solved by pre-registering a study!
- Multiverse analysis is similar in concepts to sensitivity and outlier analysis, but it differs by **explicitly recognizing and incorporating uncertainty in statistical results induced by arbitrary choices.**

Discussion: What to do?

- One way to go around this problem is to take the average over the multiverse of statistical results
 - Estimate:
 - Mean = .024 (.018-.029) for Log and .114 (.095-.123) for Reciprocal
 - p -value:
 - Mean = .042 (.010-.129) for Log and .055 (.026-.106) for Reciprocal

Discussion: What to do?

And then make a statistical inference

- with an effect-based term:
 - 1 average over all the (multiverse of) regression models
 - conduct the same analysis for all other regression coefficients
 - 2 enter the value of independent variables (e.g., Gram vs Ungram) to linear predictors
 - 3 then, interpret the predicted value from the averaged model
 - What if the predicted difference between Gram and Ungram is 50ms?
- with uncertainty embraced and admitted
 - “**Statistics is the science of uncertainty and variation**” (Gelman, 2018, p. 41), but often, it is treated as a form of alchemy that coverts randomness into certainty - **uncertainty laundering**

Discussion: Moving forward

- Multiverse analysis is certainly recommended!
 - It does not have to be on the main document
- In Bayesian analysis, you can not only incorporate uncertainty of parameter estimation but also uncertainty that is caused by the multiverse of statistical results

Thank you for listening!

Acknowledgement

- Robert DeKeyser (University of Maryland)
- Michael Long (University of Maryland)
- Steve Ross (University of Maryland)
- Bronson Hui (Michigan State University)
- Masaki Eguchi (University of Oregon)

References - 1

Gelman, A. (2018). Ethics in statistical practice and communication: Five recommendations. *Royal Statistical Society*, 15, 40-43.

Godfroid, A. (2016). The effects of implicit instruction on implicit and explicit knowledge development. *Studies in Second Language Acquisition*, 38, 177-215.

Granena, G. (2012). *Age differences and cognitive aptitudes for implicit and explicit learning in ultimate second language attainment* (Unpublished doctoral dissertation). University of Maryland, College Park, College Park, MD.

Jiang, N. (2004). Morphological insensitivity in second language processing. *Applied Psycholinguistics*, 25, 603-634.

References - 2

Jiang, N. (2012). *Conducting reaction time research in second language studies*. London, UK: Routledge.

Jiang, N., Novokshanova, E., Masuda, K., & Wang, X. (2011). Morphological congruency and the acquisition of L2 morphemes. *Language Learning*, 61, 940-967.

Maie, R., & DeKeyser, R. M. (2020). Conflicting evidence of explicit and implicit knowledge from objective and subjective measures. *Studies in Second Language Acquisition*, 42, 359-382.

Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-positive psychology: Undisclosed flexibility in data collection and analysis allows presenting anything as significant. *Psychological Science*, 22, 1359-1366

References - 3

Steegen, S., Tuerlinckx, F., Gelman, A., & Vanpaemel, W. (2016). Increasing transparency through a multiverse analysis. *Perspectives on Psychological Science*, 11, 702-712.

Suzuki, Y. (2015). *Using new measures of implicit L2 knowledge to study the interface of explicit and implicit knowledge* (Unpublished doctoral dissertation). University of Maryland, College Park, College Park, MD.