



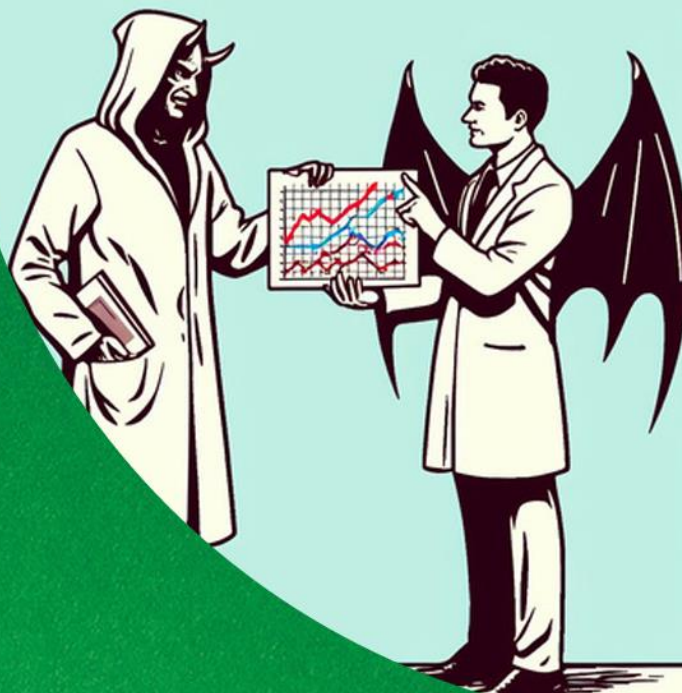
Palacký University  
Olomouc

FTK UP pořádá odbornou diskusi

# Replikační krize ve vědě: 7 (ne)smrtných hříchů moderní statistiky

Změna termínu:

25. března 2025 13:30-15:00 NA 2.10



Fakulta  
tělesné kultury  
Univerzita Palackého  
v Olomouci



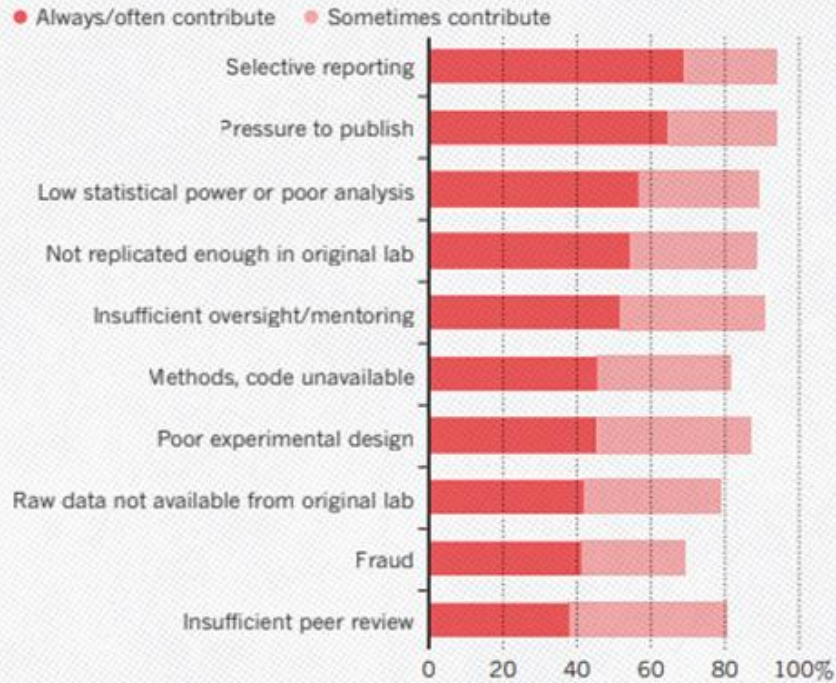
## Předpokládaný časový harmonogram

- 13:30-13:50 přivítání a úvod: Ladislav Baloun (Proč řešit 7 nesmrtelných hříchů moderní statistiky?).
- 13:50-14:05 tematická přednáška: Tomáš Fürst (13 lekcí o vědě, které by měl slyšet každý prvák na UP).
- 14:05-14:20 tematická přednáška: Ondřej Vencálek (Trable s p-values).
- 14:20-15:00 otevřená diskuse.



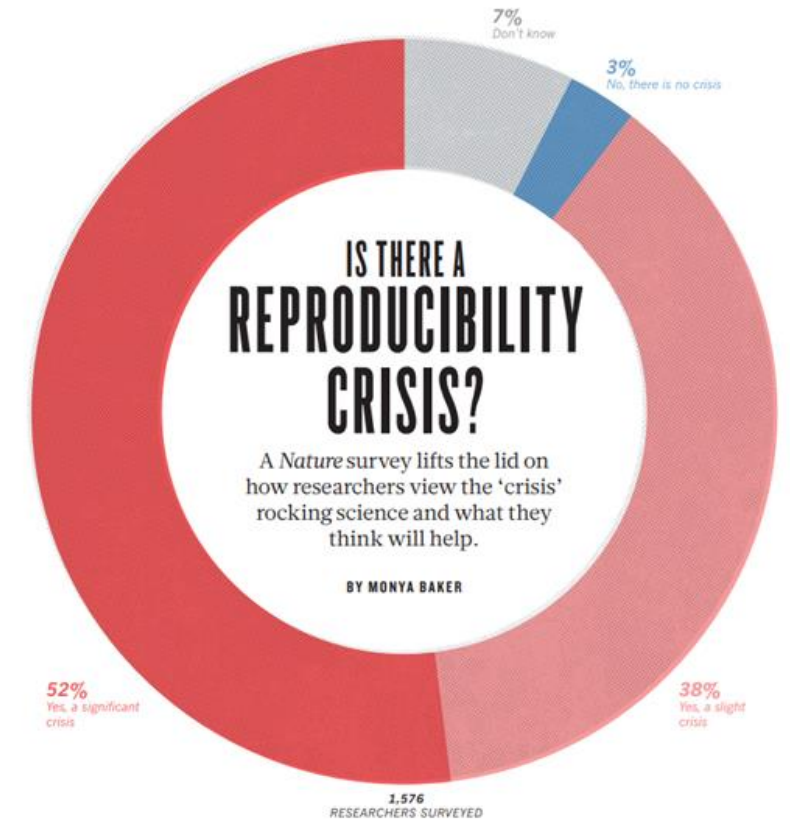
## WHAT FACTORS CONTRIBUTE TO IRREPRODUCIBLE RESEARCH?

Many top-rated factors relate to intense competition and time pressure.



## WHAT FACTORS COULD BOOST REPRODUCIBILITY?

Respondents were positive about most proposed improvements but emphasized training in particular.



# Proč sedm (ne)smrtebných hřichů moderní statistiky?

Ladislav Baloun, FTK UP

ladislav.baloun@upol.cz

Baker, M. (2016). Reproducibility crisis. *Nature*, 533(26), 353-66



## 7 hříchů statistické inference

1. Málo časové dotace pro výuku statistiky
2. Špatné představy o tom jak data vznikají
3. Smal sample size
4. Předpoklad objektivních výsledků: subjektivní analýza objektivních dat.
5. Dichotomizace výsledků – obsese pro pozitivní výsledky
6. P hacking
7. Harking
8. Selektivní reportování



## Předpoklad objektivních výsledků: subjektivní analýza objektivních dat.

- ....data analysts should recognize that subjectivity and potential bias are inherent in all data analysis, exploratory or otherwise.
- ...datoví analytici by si měli uvědomit, že subjektivita a potenciální zkreslení jsou vlastní veškeré analýze dat, průzkumné i jiné.
- One great danger in overmathematizing data analysis is believing that the reliability and precision of mathematics itself imbue reliability and precision to the data and the data analysis.
- Jedním velkým nebezpečím při přematematizaci analýzy dat je věřit, že spolehlivost a přesnost matematiky sama o sobě dodává spolehlivost a přesnost datům a analýze dat.



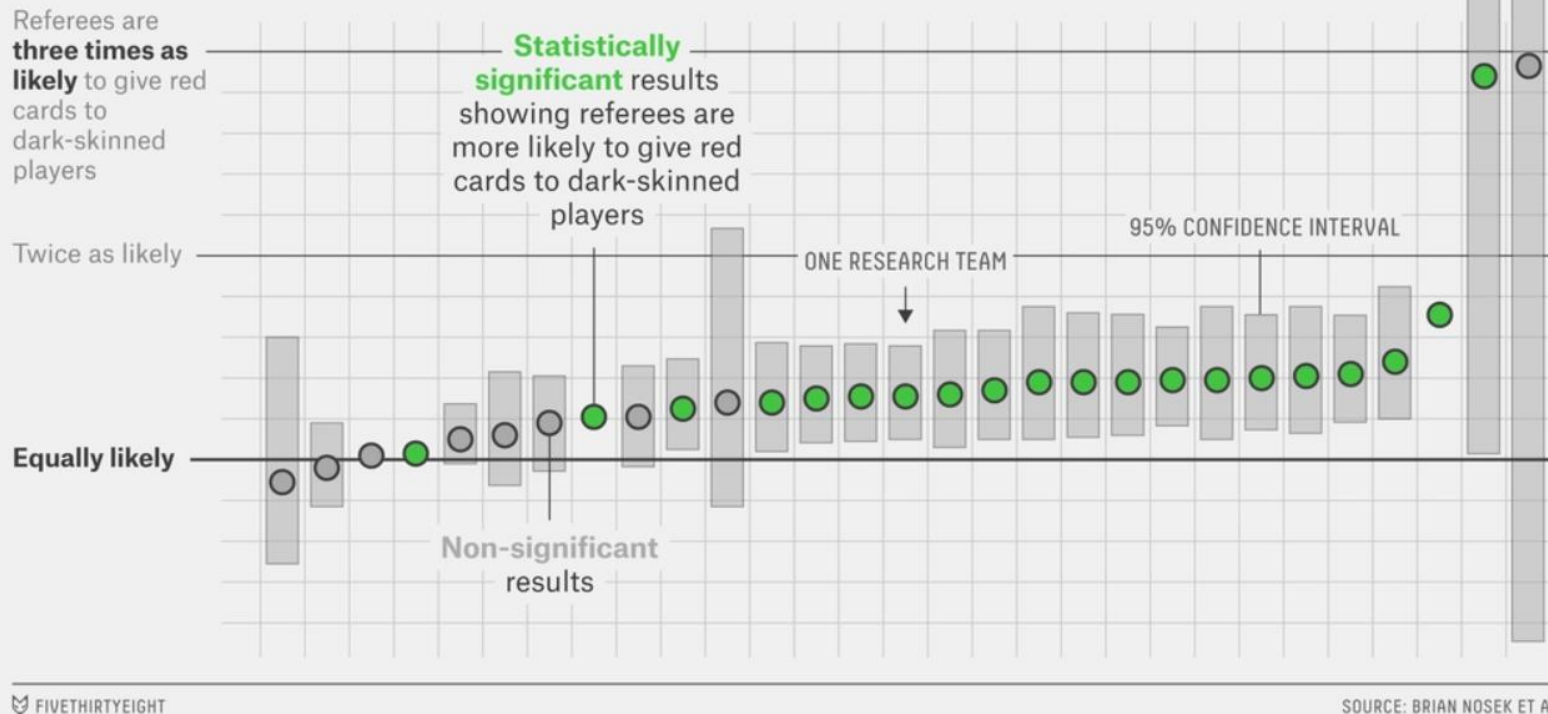


# Reproducibility - How do we work with data?

Silberzahn et al. (2018): different approaches, different type of analysis.

## Same Data, Different Conclusions

Twenty-nine research teams were given the same set of soccer data and asked to determine if referees are more likely to give red cards to dark-skinned players. Each team used a different statistical method, and each found a different relationship between skin color and red cards.



Many hands  
make tight work

# How do we work with data?

RESEARCH ARTICLE | SOCIAL SCIENCES | 8



## Observing many researchers using the same data and hypothesis reveals a hidden universe of uncertainty

Nate Breznau , Eike Mark Rinke , Alexander Wuttke , , , and Tomasz Żółtak  [Authors Info & Affiliations](#)

Edited by Douglas Massey, Princeton University, Princeton, NJ; received March 6, 2022; accepted August 22, 2022

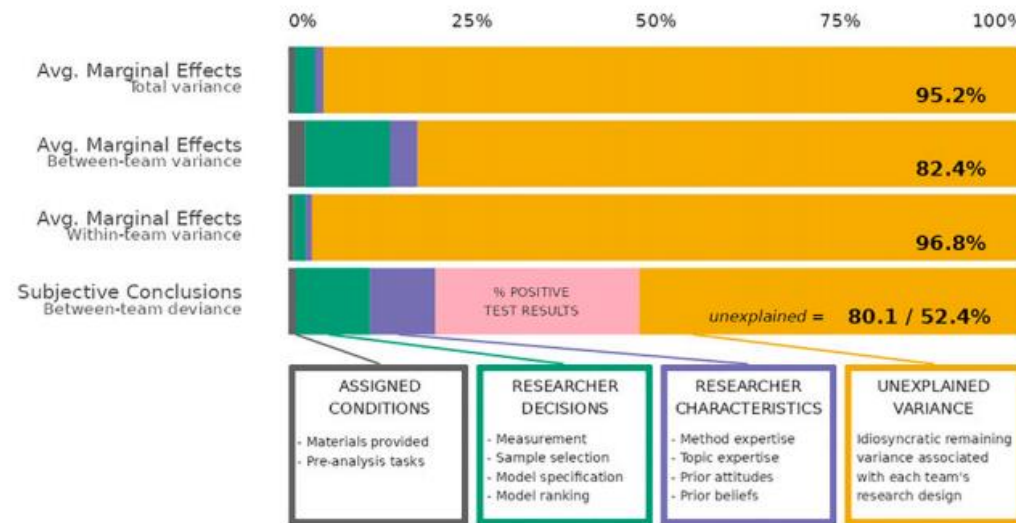
October 28, 2022 | 119 (44) e2203150119 | <https://doi.org/10.1073/pnas.2203150119>

[VIEW RELATED CONTENT](#) +

There are two primary explanations for variation in forking decisions. The competency hypothesis posits that researchers may make different analytical decisions because of varying levels of statistical and subject expertise that lead to different judgments as to what constitutes the “ideal” analysis in a given research situation. The confirmation bias hypothesis holds that researchers may make reliably different analytical choices because of differences in preexisting beliefs and attitudes, which may lead to justification of analytical approaches favoring certain outcomes post hoc. However, many other covert or idiosyncratic influences, large and small, may also lead to unreliable and unexplainable variation in analytical decision pathways (10). Sometimes even the tiniest of these differences may add up and interact to produce widely varying outcomes.

73 výzkumných týmů (identická data z mezinárodního šetření)

hypotéza: větší imigrace sníží veřejnou podporu pro vládní poskytování sociálních politik



**Fig. 1.** Broad variation in the findings from 73 teams testing the same hypothesis with the same data. The distribution of estimated AMEs across all converged models ( $n = 1,253$ ) includes results that are negative (yellow; in the direction predicted by the given hypothesis the teams were testing), not different from zero (gray), or positive (blue) using a 95% CI. AME are  $xy$  standardized. The y axis contains two scaling breaks at  $\pm 0.05$ . Numbers inside circles represent the percentages of the distribution of each outcome inversely weighted by the number of models per team.

a statistical results and substantive conclusions between and within teams is mostly unexplained by conditions, researcher





# How do we work with data?

## Variability in the analysis of a single neuroimaging dataset by many teams

<https://doi.org/10.1038/s41586-020-2314-9>

Received: 14 November 2019

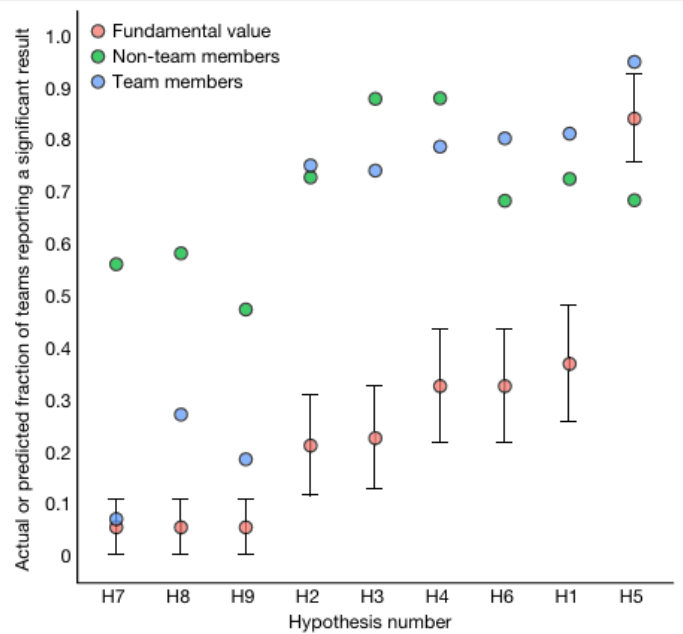
Accepted: 7 April 2020

Published online: 20 May 2020

Check for updates

A list of authors and affiliations appears in the online version of the paper.

Data analysis workflows in many scientific domains have become increasingly complex and flexible. Here we assess the effect of this flexibility on the results of functional magnetic resonance imaging by asking 70 independent teams to analyse the same dataset, testing the same 9 ex-ante hypotheses<sup>1</sup>. The flexibility of analytical approaches is exemplified by the fact that no two teams chose identical workflows to analyse the data. This flexibility resulted in sizeable variation in the results of hypothesis tests, even for teams whose statistical maps were highly correlated at intermediate stages of the analysis pipeline. Variation in reported results was related to several aspects of analysis methodology. Notably, a meta-analytical approach that aggregated information across teams yielded a significant consensus in activated regions. Furthermore, prediction markets of researchers in the field revealed an overestimation of the likelihood of significant findings, even by researchers with direct knowledge of the dataset<sup>2-5</sup>. Our findings show that analytical flexibility can have substantial effects on scientific conclusions, and identify factors that may be related to variability in the analysis of functional magnetic resonance imaging. The results emphasize the importance of validating and sharing complex analysis workflows, and demonstrate the need for performing and reporting multiple analyses of the same data. Potential approaches that could be used to mitigate issues related to analytical variability are discussed.



**Fig. 1 | Fraction of teams reporting a significant result and prediction market beliefs.** The observed fraction of teams reporting significant results (fundamental value, pink dots;  $n = 70$  analysis teams), as well as final market prices for the team members markets (blue dots;  $n = 83$  active traders) and the non-team members markets (green dots;  $n = 65$  active traders). The corresponding 95% confidence intervals are shown for each of the nine hypotheses (note that hypotheses are sorted on the basis of the fundamental value). Confidence intervals were constructed by assuming convergence of the binomial distribution towards the normal.

fMRI data from 108 individuals  
70 teams  
9 hypotheses

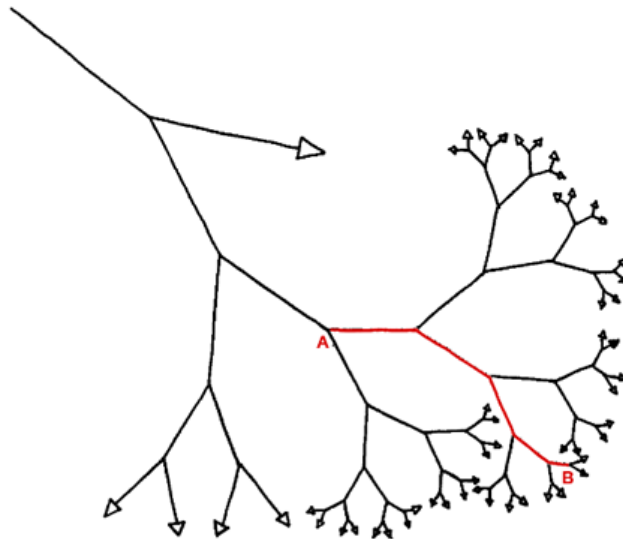




# Zahrada rozvětvených cestiček

Researcher degrees of freedom (Stupně volnosti výzkumníka)

## The Garden of Forking Paths by Jorge Luis Borges



**A dataset can be analyzed in so many different ways** (with the choices being not just what statistical test to perform but also decisions on what data to exclude or include, what measures to study, what interactions to consider, etc.),

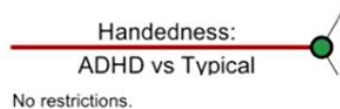


Palacký University  
Olomouc

# Zahrada rozvětvených cestiček

## Researcher degrees of freedom (Stupně volnosti výzkumníka)

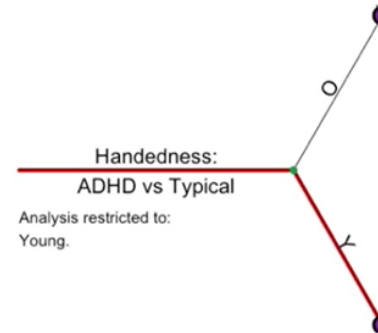
Large population  
database used to explore  
link between ADHD and  
handedness



1 contrast

Probability of a  
'significant' p-value  
< .05 = .05

Large population  
database used to explore  
link between ADHD and  
handedness



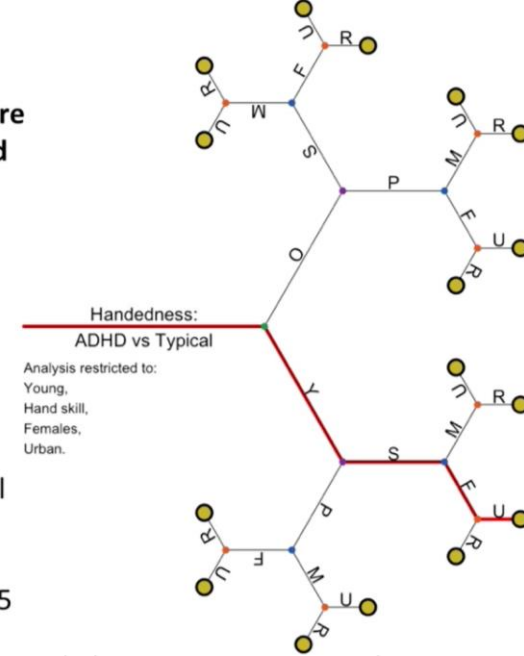
Focus just on Young  
subgroup:  
2 contrasts at this level

Probability of a  
'significant' p-value < .05  
= .10

Large population  
database used to explore  
link between ADHD and  
handedness

Focus just on Young,  
Urban, Females on  
measure of hand skill:  
16 contrasts at this level

Probability of a  
'significant' p-value < .05  
= .56





Palacký University  
Olomouc

# Zahrada rozvětvených cestiček

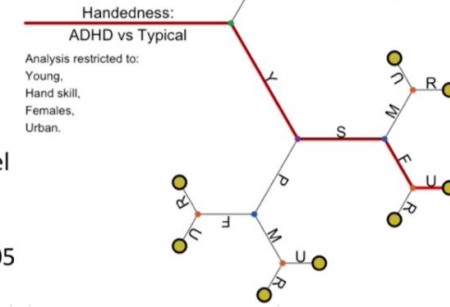
Researcher degrees of freedom (Stupně volnosti výzkumníka)

- Výzkumníci nezkouší více testů, aby zjistili, který má nejlepší p-hodnotu; spíše používají svůj vědecký zdravý rozum, aby formulovali své hypotézy rozumným způsobem, s ohledem na data, která mají.
- Chyba je v domněnání, že pokud konkrétní cesta, která byla zvolena, poskytuje statistickou významnost, že je to silný důkaz ve prospěch hypotézy.

Large population database used to explore link between ADHD and handedness

Focus just on Young, Urban, Females on measure of hand skill: 16 contrasts at this level

Probability of a 'significant' p-value < .05 = .56







# Zahrada rozvětvených cestiček

## Researcher degrees of freedom (Stupně volnosti výzkumníka)

- Chceme objasnit, že vícenásobné srovnání může být velkým problémem, aniž by to znamenalo, že dotyční výzkumníci podvádějí nebo jsou hloupí nebo se snaží systém manipulovat.
- Když jsme tyto druhy výzkumu popsali jako rybářské výpravy, udělali jsme chybu.
- Fishing a p-hacking znamenají aktivní snahu o statistickou významnost, zatímco to, o co by se zde mohlo stát, je soubor možností analýzy dat, které by mohly být rozumné, nebýt problémů s malou velikostí vzorku a chybou měření, kvůli kterým jsou výsledky hlučnější, než si lidé uvědomují.

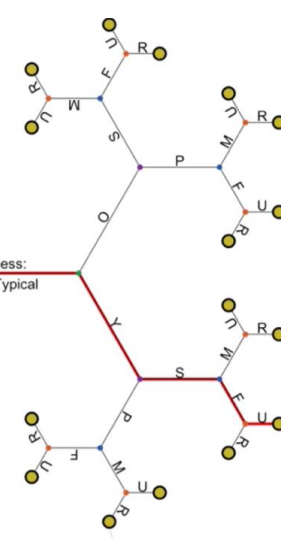
Large population database used to explore link between ADHD and handedness

Focus just on Young, Urban, Females on measure of hand skill: 16 contrasts at this level

Probability of a 'significant' p-value < .05 = .56

Handedness:  
ADHD vs Typical

Analysis restricted to:  
Young,  
Hand skill,  
Females,  
Urban.





Palacký University  
Olomouc

# Multiverse Analysis

## Increasing Transparency Through a Multiverse Analysis

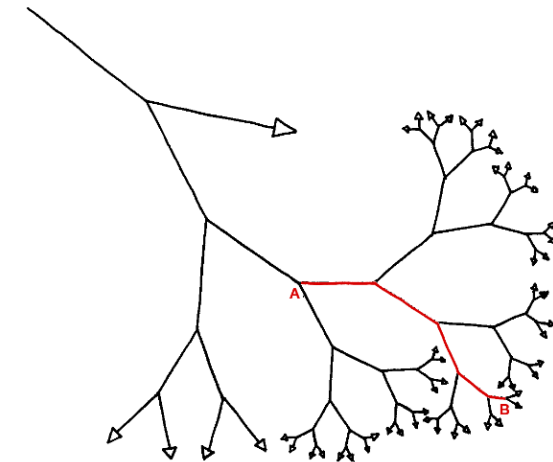
Sara Steegen<sup>1</sup>, Francis Tuerlinckx<sup>1</sup>, Andrew Gelman<sup>2</sup>, and Wolf Vanpaemel<sup>1</sup>

<sup>1</sup>KU Leuven, University of Leuven and <sup>2</sup>Columbia University

### Abstract

Empirical research inevitably includes constructing a data set by processing raw data into a form ready for statistical analysis. Data processing often involves choices among several reasonable options for excluding, transforming, and coding data. We suggest that instead of performing only one analysis, researchers could perform a multiverse analysis, which involves performing all analyses across the whole set of alternatively processed data sets corresponding to a large set of reasonable scenarios. Using an example focusing on the effect of fertility on religiosity and political attitudes, we show that analyzing a single data set can be misleading and propose a multiverse analysis as an alternative practice. A multiverse analysis offers an idea of how much the conclusions change because of arbitrary choices in data construction and gives pointers as to which choices are most consequential in the fragility of the result.

## The Garden of Forking Paths by Jorge Luis Borges



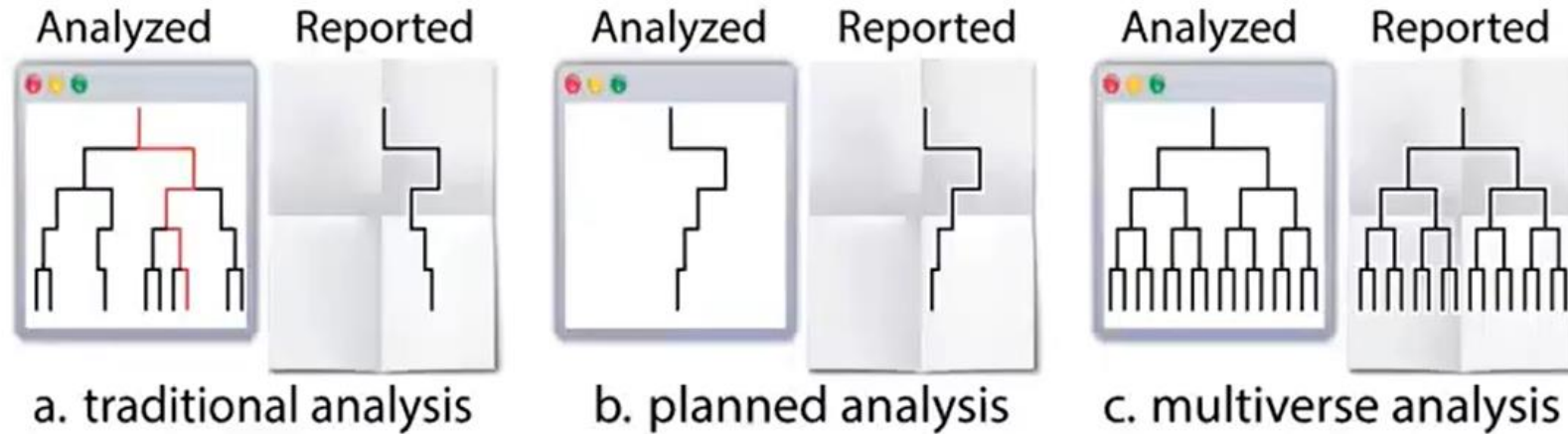
- A multiverse analysis starts from the observation that data used in an analysis are usually not just passively recorded in an experiment or an observational study.
- Rather, data are to a certain extent actively constructed.
- Data construction occurs when the raw data are converted into a form ready for analysis.
- When preparing their data for analysis, researchers often take several processing steps, such as discretization of variables into categories, combination of variables, transformation of variables, data exclusion, and so on.
- These processing steps typically come with many researcher degrees of freedom.

In the light of this problem of selective reporting, we propose to use a multiverse analysis as an alternative to a single data set analysis.



P  
C

# Multiverse Terms



nature  
human behaviour

RESOURCE

<https://doi.org/10.1038/s41562-020-0912-z>

Check for updates

## Specification curve analysis

Uri Simonsohn<sup>1</sup>, Joseph P. Simmons<sup>2</sup> and Leif D. Nelson<sup>3</sup>

**Webcast Lecture Highlight - Multiverse Analysis with Dr. Aaron Hill**

<https://www.youtube.com/watch?v=Q9bbCi6bV8Q>

**Multiverse analysis**

<https://www.youtube.com/watch?v=QfvjZL7jY24>





## Děkuji za pozornost

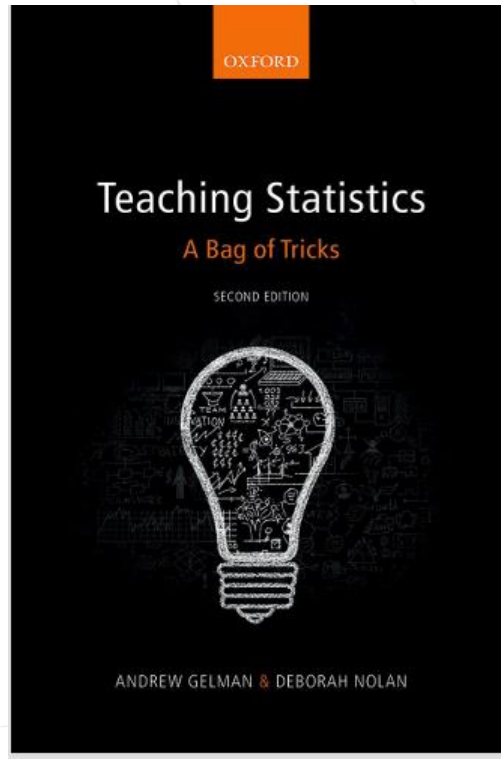
- The numbers have no way of speaking for themselves. We speak for them. We imbue them with meaning.
- Čísla nemohou mluvit sama za sebe. Mluvíme za ně. Dáváme jim význam.

Nate Silver, *The Signal and the Noise*



# Prvotní hřích – málo času

## Více důrazu a časové dotace na statistiku a analýzu dat



### 120 PROBABILITY

00111000110010000100	01000101001100010100
00100010001000000001	11101001100011110100
00110010101100001111	01110100011000110111
11001100010101100100	10001001011011011100
10001000000011111001	01100100010010000100

**Fig. 8.1** Two binary sequences produced by students in an eighth grade class for the demonstration of Section 8.3.2. Can you figure out which is the actual sequence of 100 coin flips and which is the fake? The answer appears on page 366.



# iPod Shuffle Problems: How Random is the iPod Shuffle?

Humans are not good at identifying randomness: our minds naturally look for patterns, even when there are none.

**Furthermore, we are poor at creating random data. Famously, as a result of listener complaints, the first iPod 'shuffle function' had to be changed to make it less random, but appear more random to the human ear**







# Meet Yoshitaka Fujii, the most prolific fraudster in modern science

By Joseph Stromberg | May 21, 2015, 1:10pm EDT

- Those wishing to invent data have a hard task. They must ensure that all the data satisfy several layers of statistical cross-examination  
....It is therefore always so much easier actually to do the experiment than to invent its results.

## The new retraction record holder is a German anesthesiologist, with 184

The German anesthesiologist Joachim Boldt has lost 20 more papers since January 2023, earning him the top spot in our leaderboard, with 184 retractions.

Boldt, readers may recall, was once one of the leading international figures in perioperative medicine. His work, particularly studies involving



Ludwigshafen Hospital, via Wikimedia

<https://www.vox.com/2015/5/21/8636569/retraction-yoshitaka-fujii>

<https://retractionwatch.com/2023/07/12/the-new-retraction-record-holder-is-a-german-anesthesiologist-with-184/#more-127491>

Pandit, J. J. (2012). On statistical methods to test if sampling in trials is genuinely random. *Anaesthesia*, 67(5), 456-462.

<https://nautil.us/how-the-biggest-fabricator-in-science-got-caught-235421/>



Univerzita Palackého  
v Olomouci

## Vtípek na závěr Jak učíme statistiku?

### HOW TO DRAW AN OWL



1. Draw some circles

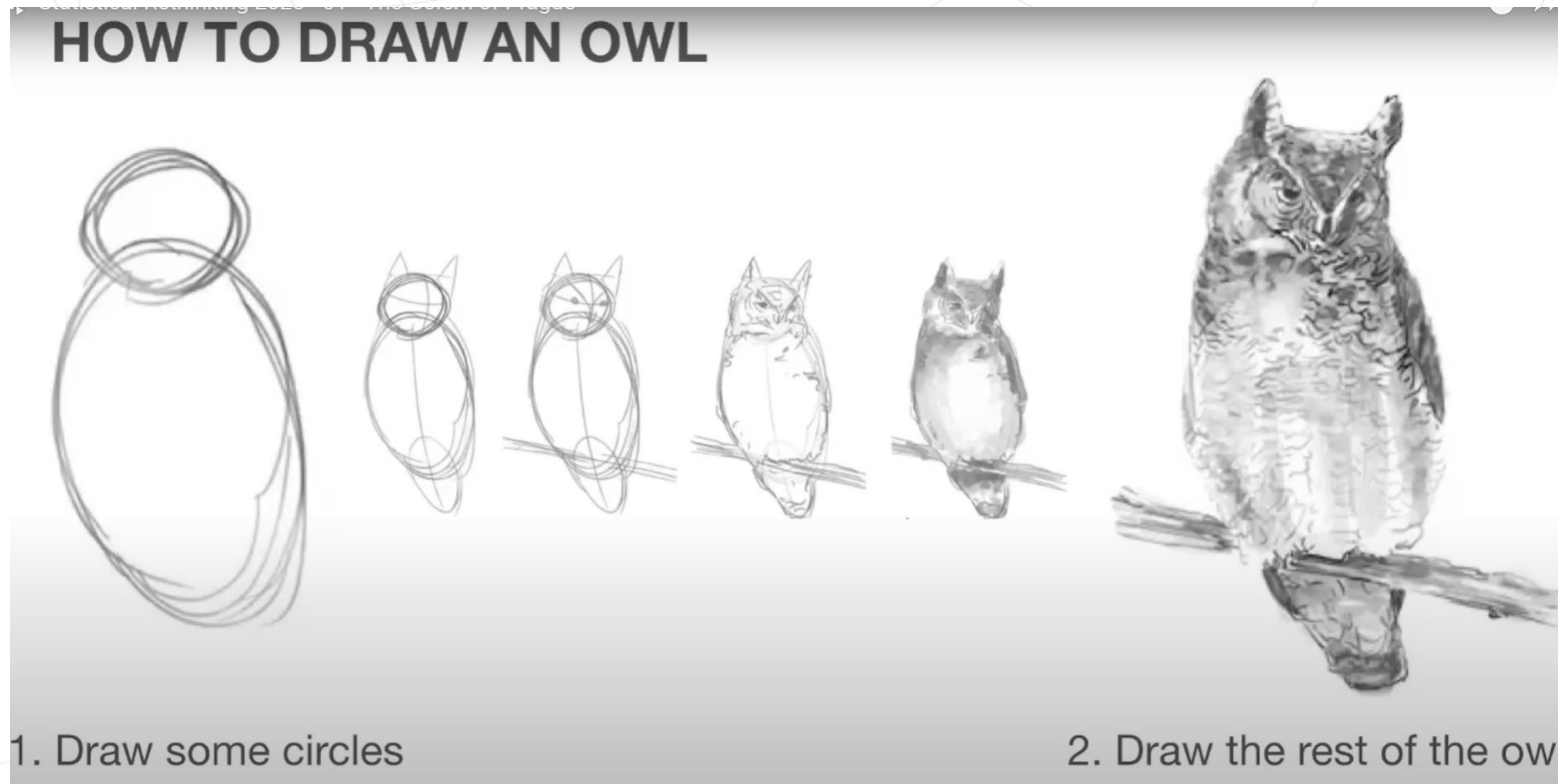


2. Draw the rest of the owl



Univerzita Palackého  
v Olomouci

## Vtípek na závěr Jak učíme statistiku?





Palacký University  
Olomouc

# Prostor na vaše dotazy

