

Statistics & Scientific reporting: Can we do better?

Abhraneel Sarma
School of Information
University of Michigan



**Scientists & researchers
can be really bad at
interpreting results of
statistical analysis.**

A simple independent means t-test
comparing the means of your control and
experimental groups ($n = 20$ each):

$t = 2.7$, d.f. = 18, $p < 0.01$

True or False?

The probability there is no difference between treatment and control is less

$t = 2.7, df = 18, p < 0.01$

[Oakes (1986): Statistical inference]

[Haller and Krauss (2002): Misinterpretations of significance: A problem students share with their teachers]

False!

The probability there is no difference
between treatment and control is less
than 1%

If there was no difference between the means of the two conditions, there is a less than 1% probability of obtaining the result

[Oakes (1986): Statistical inference]

[Haller and Krauss (2002): Misinterpretations of significance: A problem students share with their teachers]

**~ 90% people answered at least
one such questions incorrectly**

[Oakes (1986): Statistical inference]

[Haller and Krauss (2002): Misinterpretations of significance: A problem students share with their teachers]

**Misinterpretation may
lead to overestimation
of certainty**

1. Adopt Bayesian statistics
2. Uncertainty representations

1. Adopt Bayesian statistics
2. Uncertainty representations

A mixed-design ANOVA with sex of face (male, female) as a within-subjects factor and self-rated attractiveness (low, average, high) and oral contraceptive use (true, false) as between-subjects factors revealed a main effect of sex of face, $F(1, 1276) = 1372$, $p < .001$, $\eta_p^2 = .52$. This was qualified by interactions between sex of face and SRA, $F(2, 1276) = 6.90$, $p = .001$, $\eta_p^2 = .011$, and between sex of face and oral contraceptive use, $F(1, 1276) = 5.02$, $p = .025$, $\eta_p^2 = .004$. The predicted interaction among sex of face, SRA and oral contraceptive use was not significant, $F(2, 1276) = 0.06$, $p = .94$, $\eta_p^2 < .001$. All other main effects and interactions were non-significant and irrelevant to our hypotheses, all $F \leq 0.94$, $p \geq .39$, $\eta_p^2 \leq .001$.

A mixed-design ANOVA with sex of face (male, female) as within-subjects factor and self-rated attractiveness (low, average, high) and oral contraceptive use (true, false) as between-subjects factors revealed a main effect of sex of face, $F(1, 1276) = 1372, p < .01, \eta_p^2 = .52$. This was qualified by interactions between sex of face and SRA, $F(2, 1276) = 6.90, p = .001, \eta_p^2 = .011$, and between sex of face and oral contraceptive use, $F(1, 1276) = 5.02, p = .025, \eta_p^2 = .004$. The predicted interaction among sex of face, SRA and oral contraceptive use was not significant, $F(2, 1276) = 0.06, p = .94, \eta_p^2 < .001$. All other main effects and interactions were non-significant and irrelevant to our hypotheses, all $F \leq 0.94, p \geq .39, \eta_p^2 \leq .001$.

Alternatives...

Table 7
**Stevens et al. 2006, table 2: Determinants
of authoritarian aggression**

Variable	Coefficient (Standard Error)
Constant	.41 (.93)
Countries	
Argentina	1.31 (.33)**B,M
Chile	.93 (.32)**B,M
Colombia	1.46 (.32)**B,M
Mexico	.07 (.32) ^{A,CH,CO,V}
Venezuela	.96 (.37)**B,M
Threat	
Retrospective egocentric economic perceptions	.20 (.13)
Prospective egocentric economic perceptions	.22 (.12) [#]
Retrospective sociotropic economic perceptions	-.21 (.12) [#]
Prospective sociotropic economic perceptions	-.32 (.12)*
Ideological distance from president	-.27 (.07)**
Ideology	
Ideology	.23 (.07)**
Individual Differences	
Age	.00 (.01)
Female	-.03 (.21)
Education	.13 (.14)
Academic Sector	.15 (.29)
Business Sector	.31 (.25)
Government Sector	-.10 (.27)
R ²	.15
Adjusted R ²	.12
N	500

**p < .01, *p < .05, [#]p < .10 (twotailed)

[Kastellac and Leoni (2007): Using Graphs
Instead of Tables in Political Science]

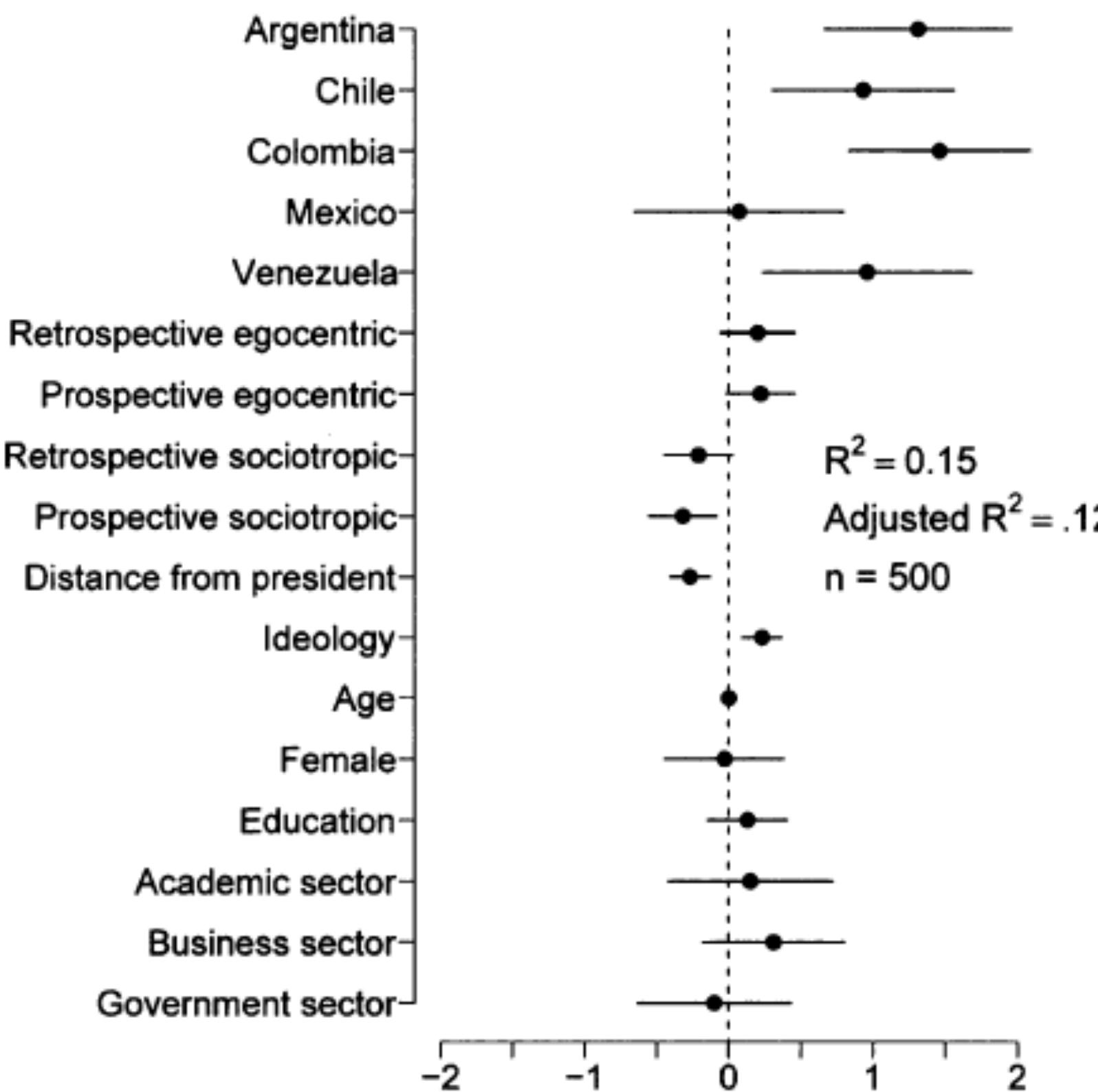
Alternatives...

Table 7
Stevens et al. 2006, table 2: Determinants
of authoritarian aggression

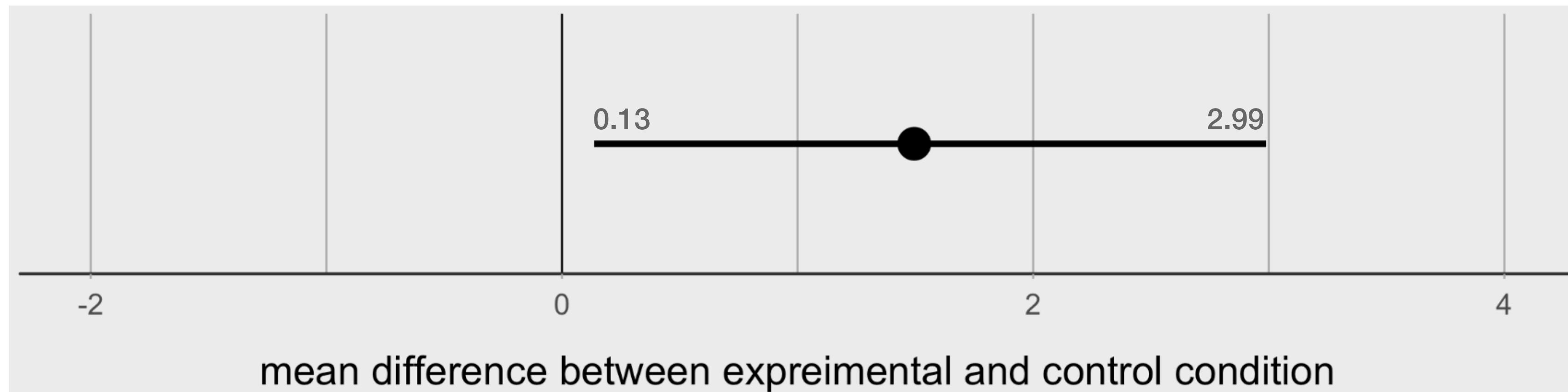
Variable	Coefficient (Standard Error)
Constant	.41 (.93)
Countries	
Argentina	1.31 (.33)**B,M
Chile	.93 (.32)**B,M
Colombia	1.46 (.32)**B,M
Mexico	.07 (.32) ^{A,CH,CO,V}
Venezuela	.96 (.37)**B,M
Threat	
Retrospective egocentric economic perceptions	.20 (.13)
Prospective egocentric economic perceptions	.22 (.12) [#]
Retrospective sociotropic economic perceptions	-.21 (.12) [#]
Prospective sociotropic economic perceptions	-.32 (.12)*
Ideological distance from president	-.27 (.07)**
Ideology	
Ideology	.23 (.07)**
Individual Differences	
Age	.00 (.01)
Female	-.03 (.21)
Education	.13 (.14)
Academic Sector	.15 (.29)
Business Sector	.31 (.25)
Government Sector	-.10 (.27)
R ²	.15
Adjusted R ²	.12
N	500

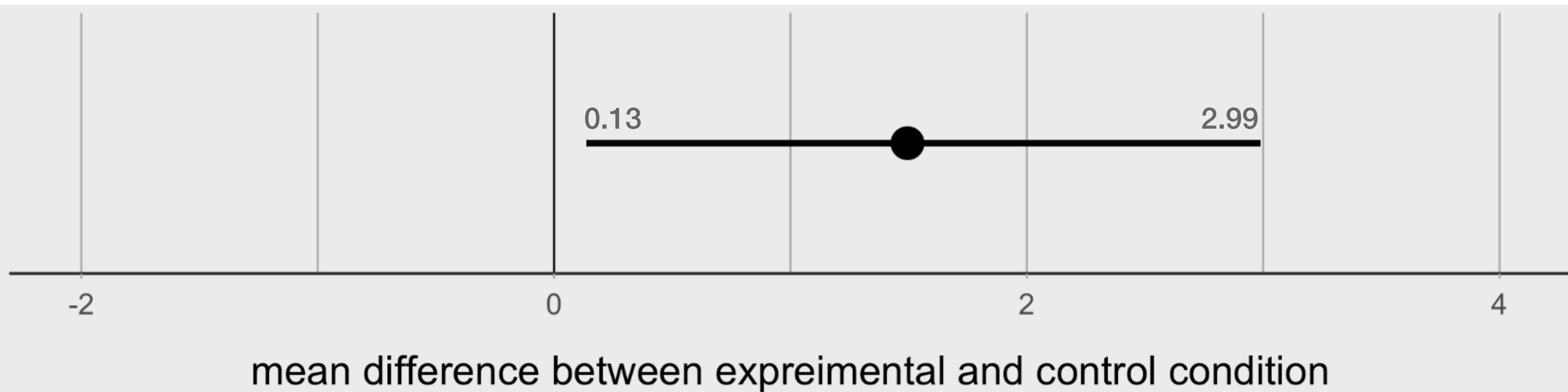
**p < .01, *p < .05, [#]p < .10 (twotailed)

Figure 6
Presenting a single regression model using
a dot plot with error bars.



[Kastellac and Leoni (2007): Using Graphs
Instead of Tables in Political Science]

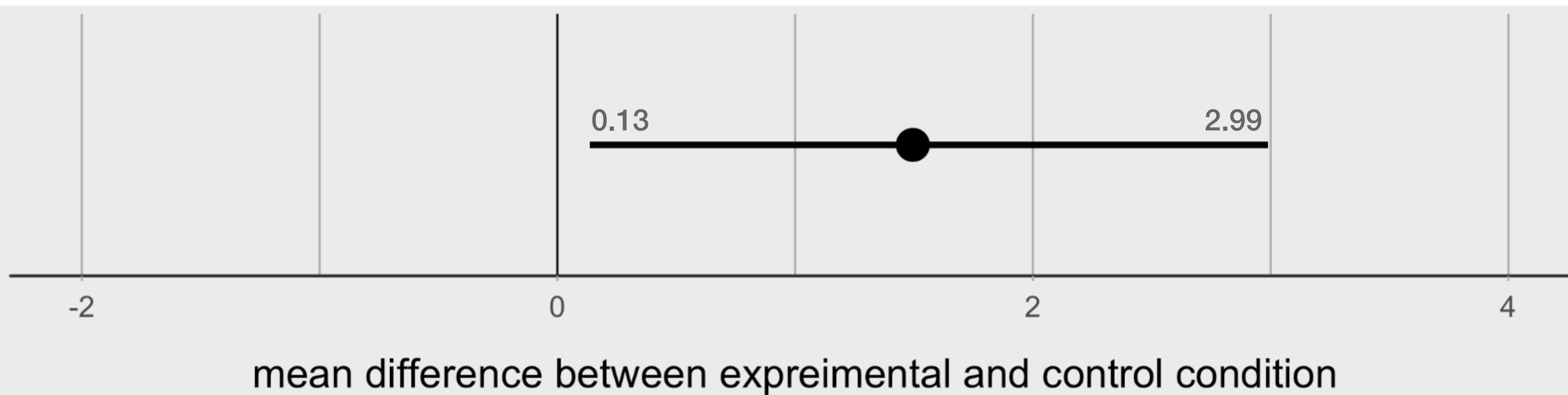




There is a 95% probability that the mean difference between the experimental and control conditions lie in the interval [0.13, 2.99]

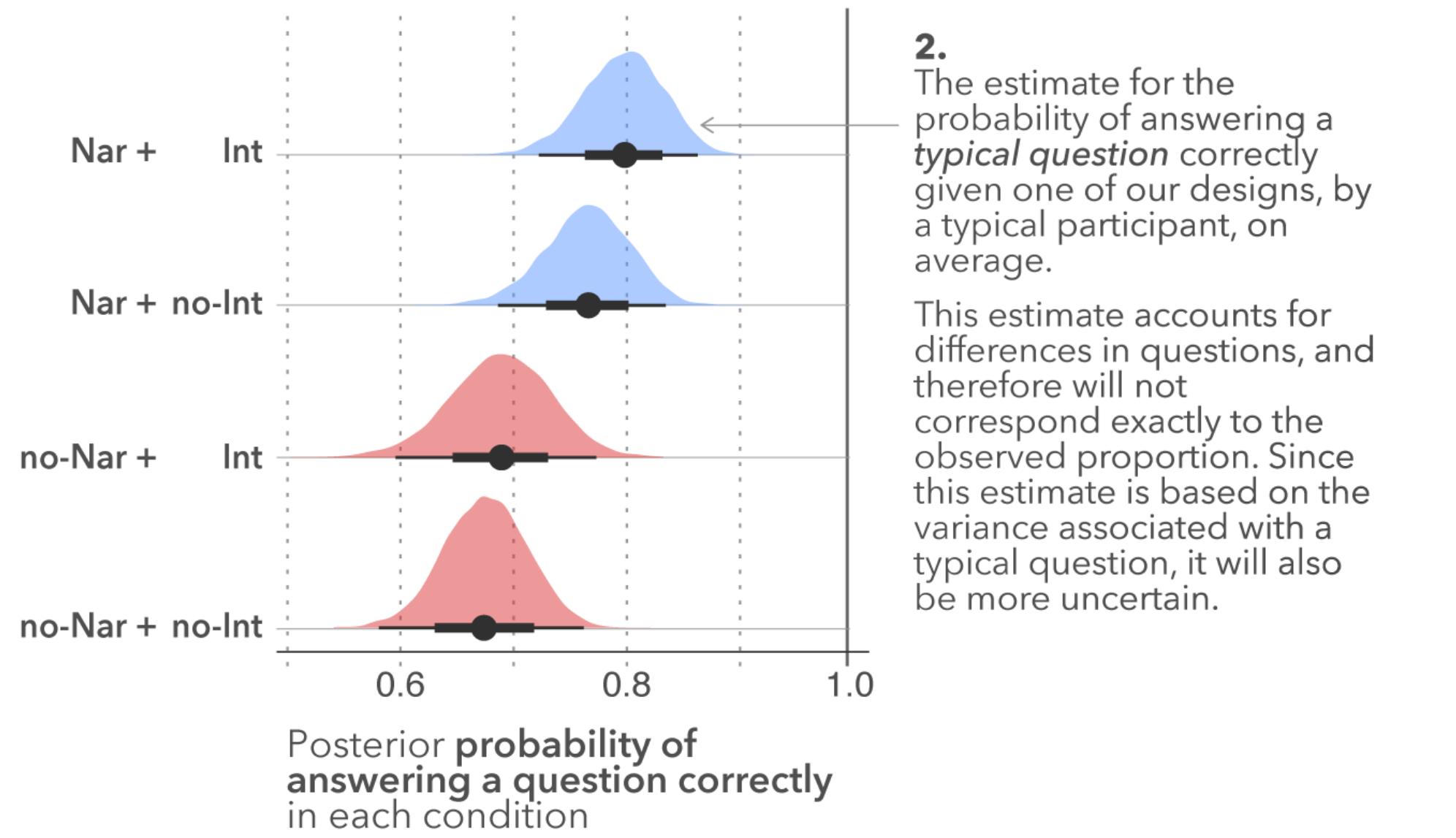
False!

~~There is a 95% probability that the mean difference between the experimental and control conditions lie in the interval [0.13, 2.99]~~

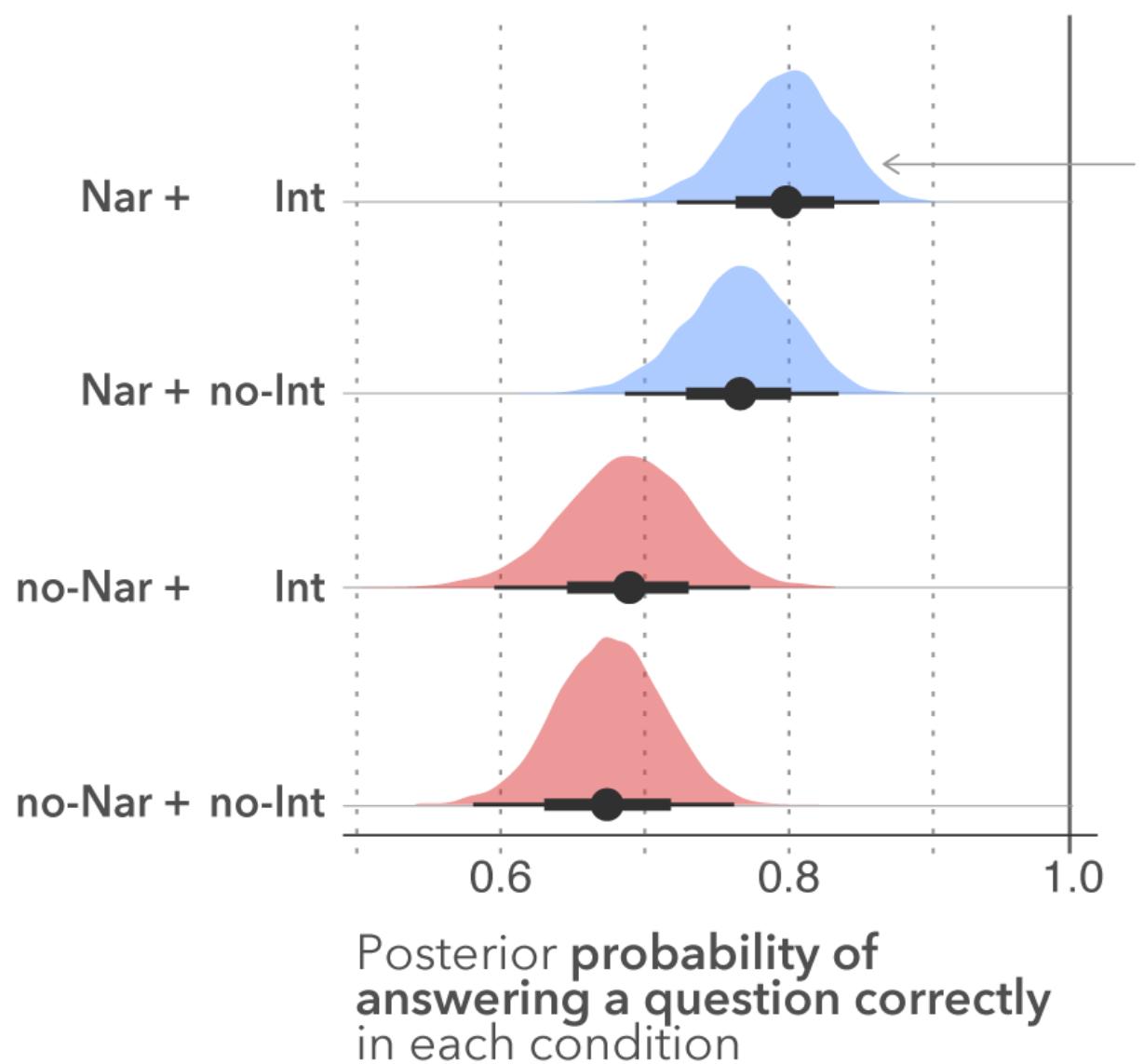


If you were to repeat the experiment over and over, then the fraction of calculated confidence intervals (which would differ for each sample) that encompass the true population parameter would tend towards 95%.

More alternatives...



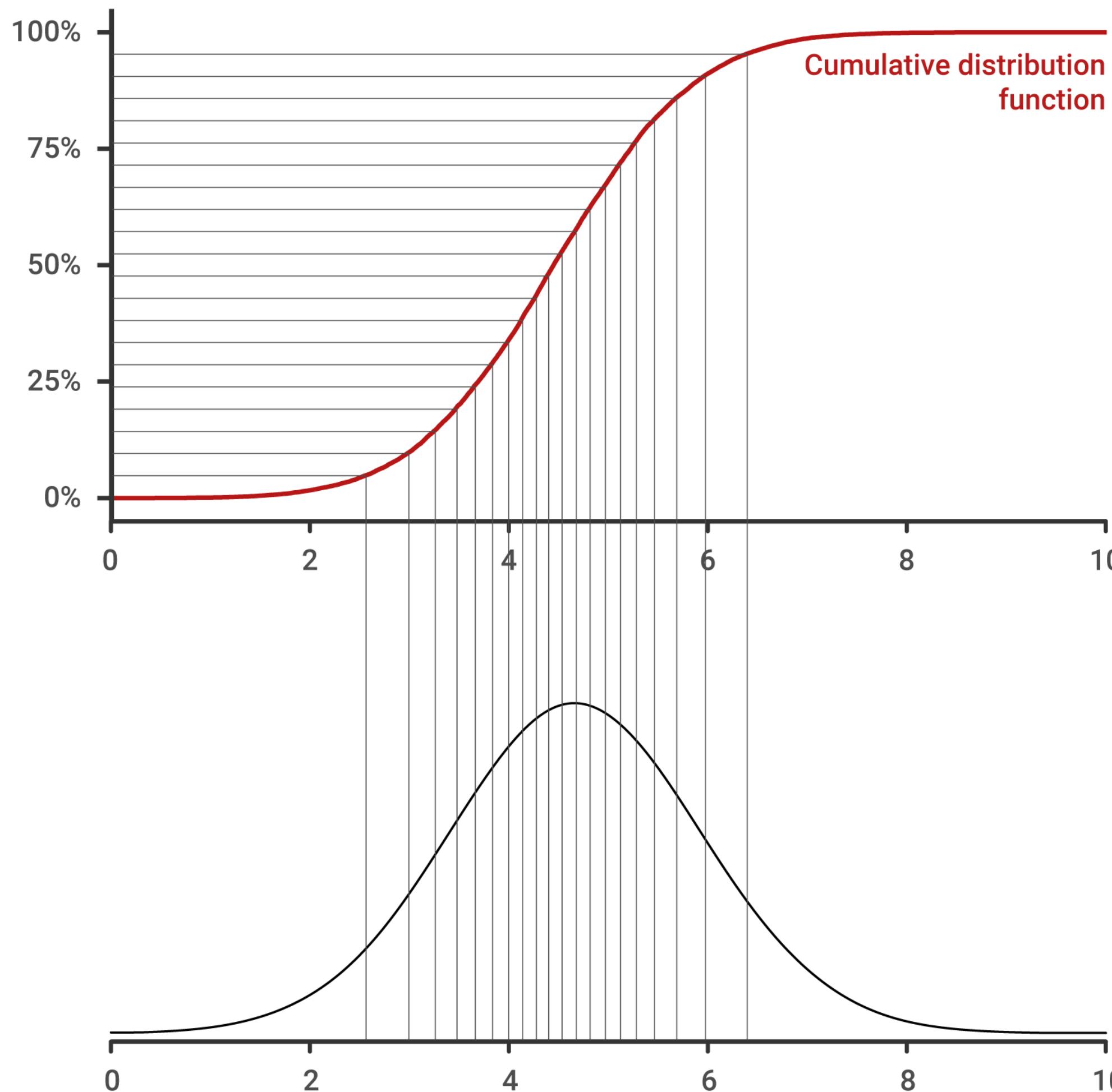
More alternatives...



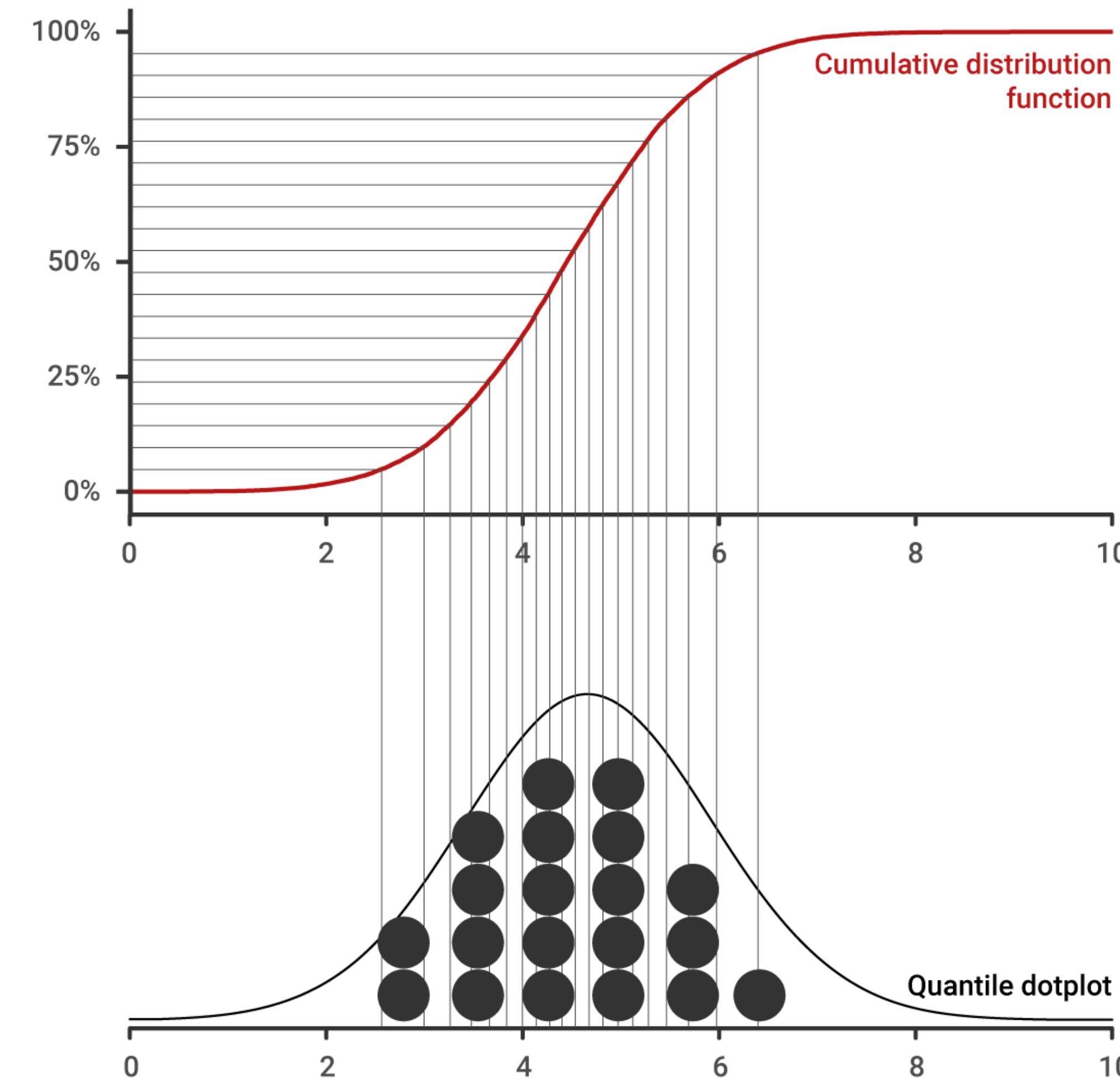
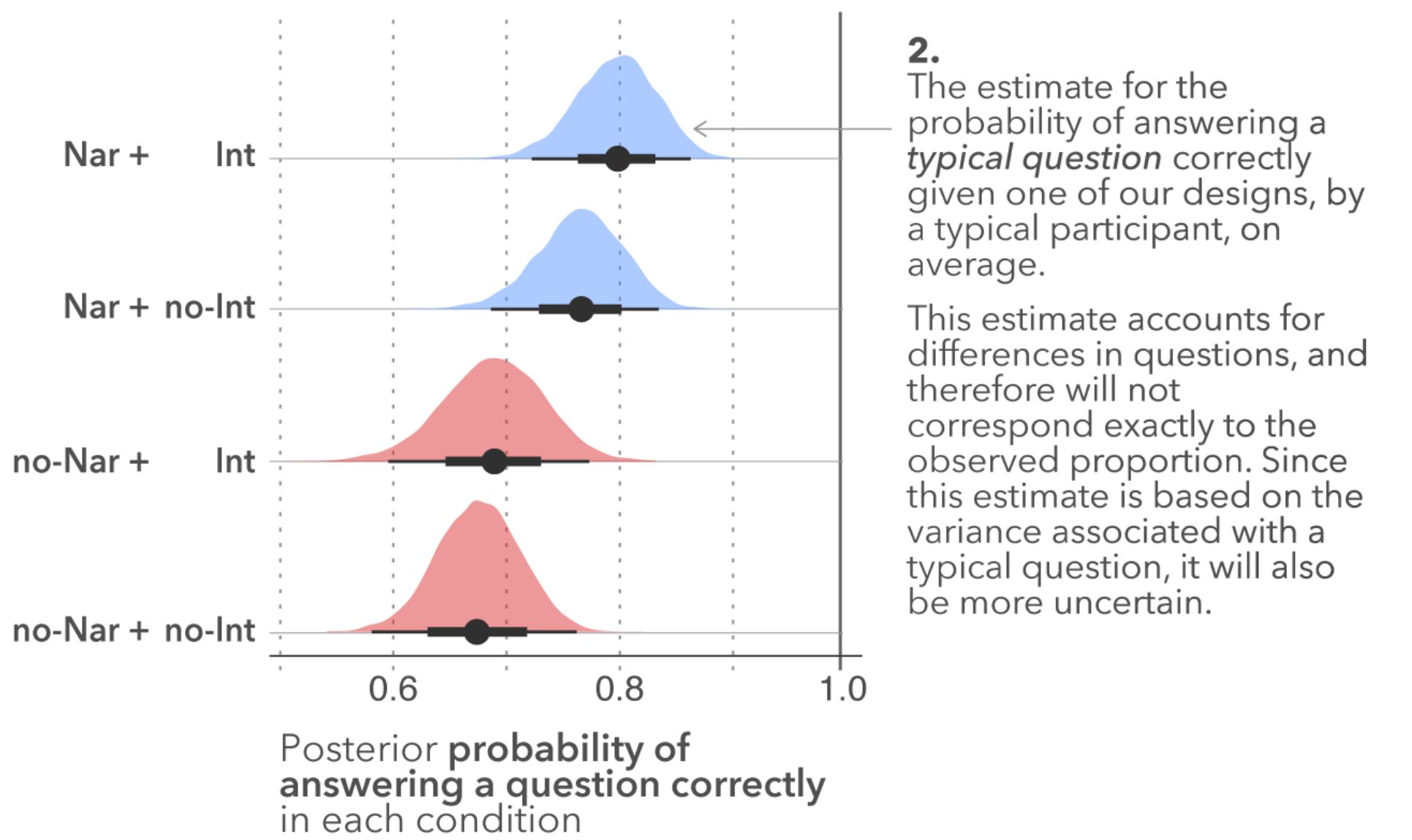
2.

The estimate for the probability of answering a *typical question* correctly given one of our designs, by a typical participant, on average.

This estimate accounts for differences in questions, and therefore will not correspond exactly to the observed proportion. Since this estimate is based on the variance associated with a typical question, it will also be more uncertain.



More alternatives...

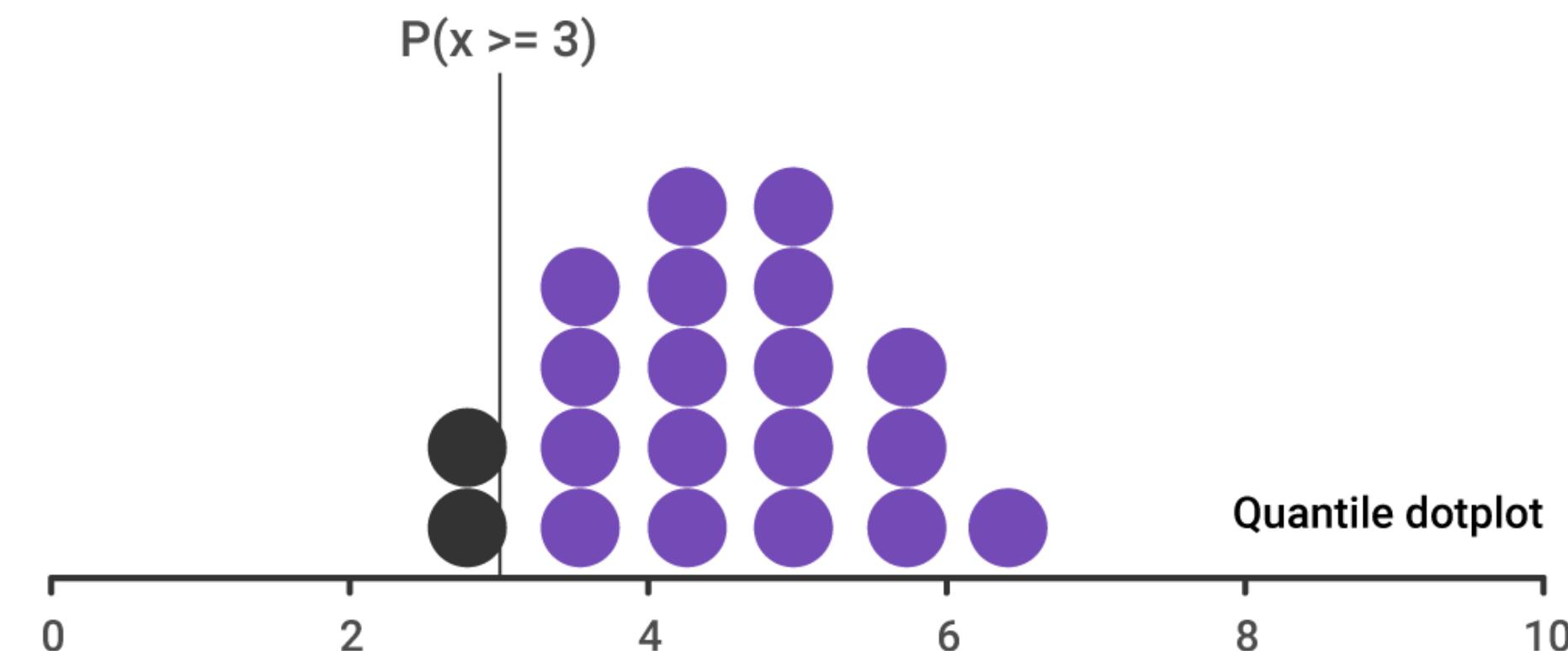


[Kay, Kola, Hullman and Munson (2016): When (ish) is my bus?: User-centered visualizations of uncertainty in everyday, mobile predictive systems]

More alternatives...

What is the probability of $x \geq 3$?

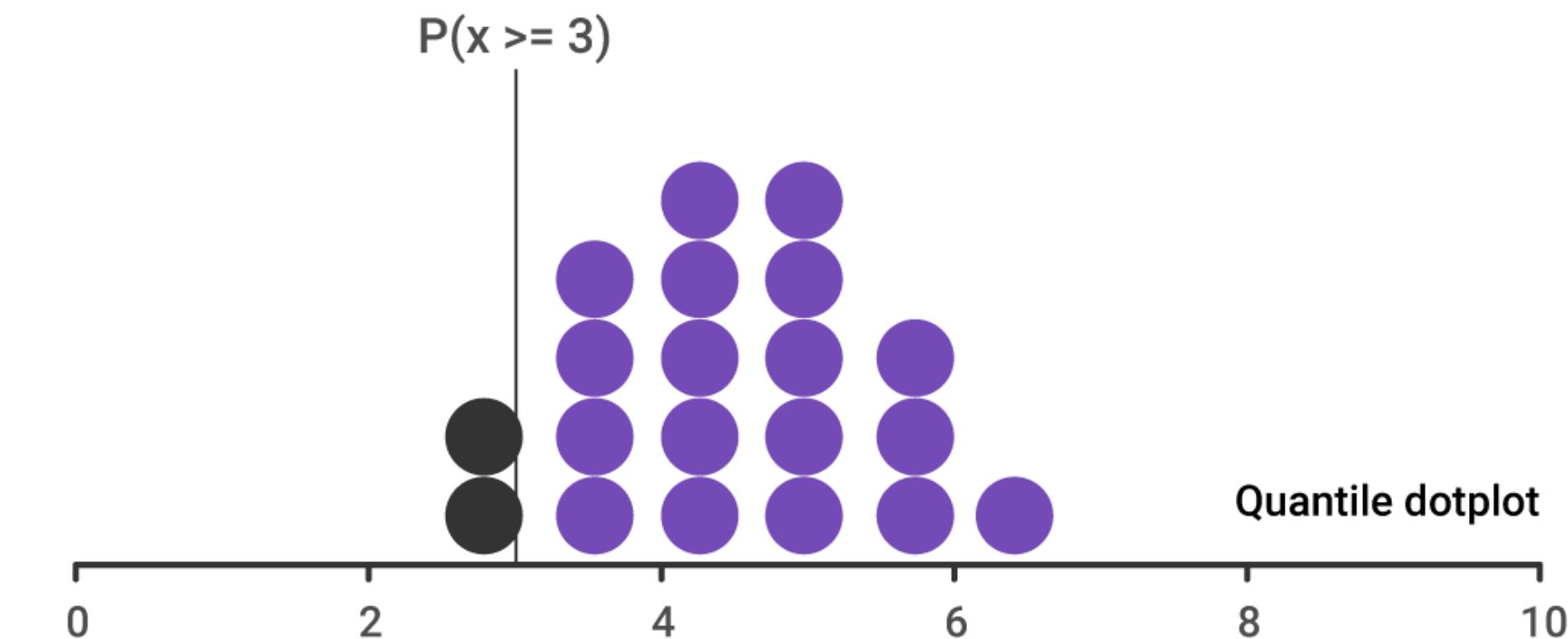
~ 90%



[Kay, Kola, Hullman and Munson (2016): When (ish) is my bus?: User-centered visualizations of uncertainty in everyday, mobile predictive systems]

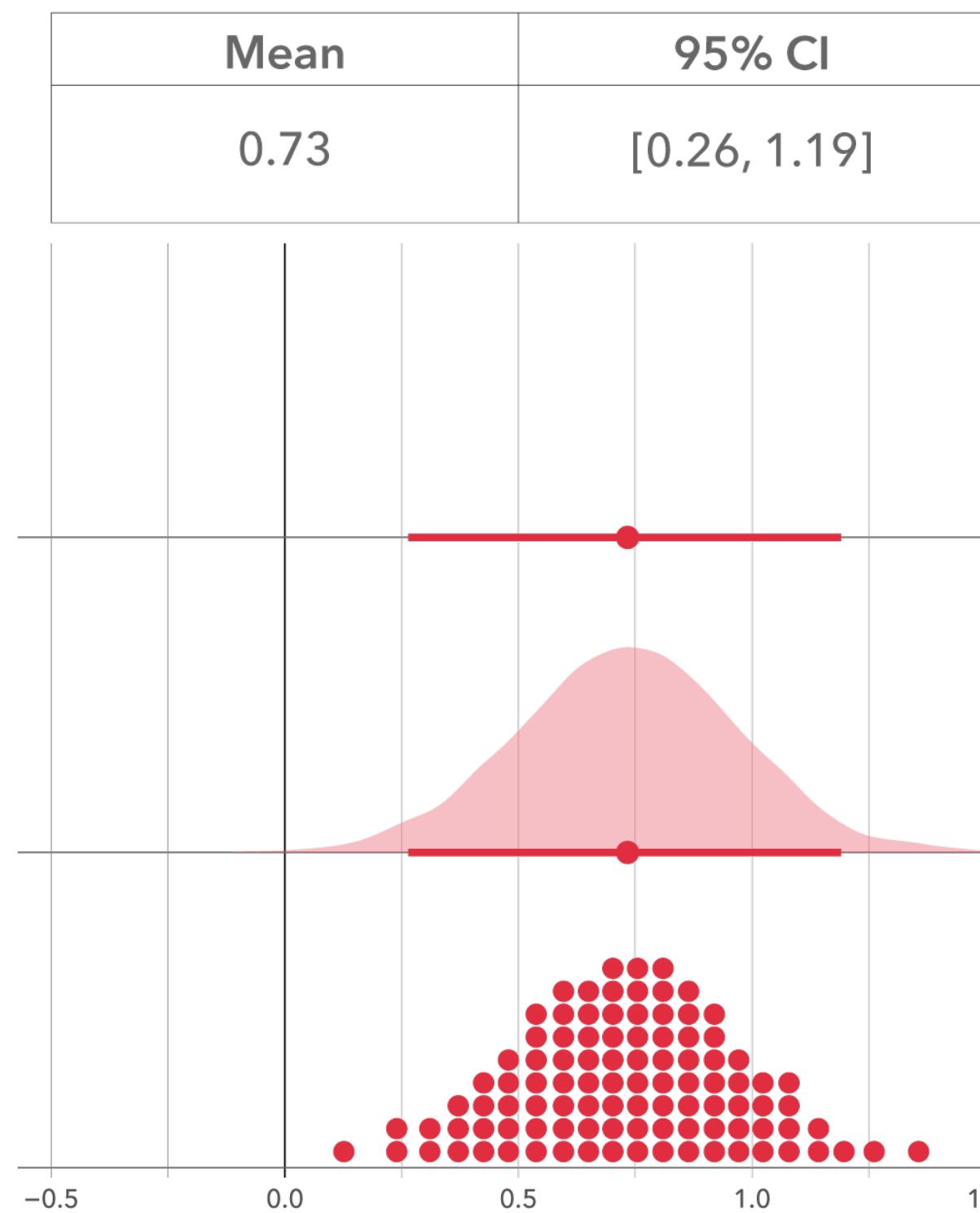
More alternatives...

On average,
quantile dotplots
with 50 outcomes
improve transit
decision making.



[Fernandes, Walls, Munson, Hullman, and Kay (2018):
Uncertainty Displays Using Quantile Dotplots or CDFs
Improve Transit Decision-Making]

Are these better at addressing misinterpretation?



A. Text tables with confidence interval (here, 95%)

Textual communication of point estimates and the corresponding 95% confidence interval is still commonly used in NHST which is advocated widely to be used in place of p-values (Cumming et. al.); if the null hypothesis is outside the interval, $p < 0.05$.

B. 95% interval

The point estimate and 95% confidence interval from (A) represented graphically. Graphical representation of statistical results have been commonly advocated as another way of emphasizing uncertainty.

C. Density + interval plot (here 95%)

The density shows the sampling distribution (or Bayesian posterior distribution). The addition of the point estimate and interval adds precision.

D. Quantile dotplot (here 100 dots)

A quantile dotplot allows precise estimation of many intervals and by providing a frequency based framing of the probability of the parameter, may improve understanding of uncertainty through hypothetical outcomes.

Are these better at addressing misinterpretation?



A. Text tables with confidence interval (here, 95%)
Textual communication of point estimates and the corresponding 95% confidence interval is still commonly used in NHST which is advocated widely to be used in place of p-values (Cumming et. al.); if the null hypothesis is outside the interval, $p < 0.05$.

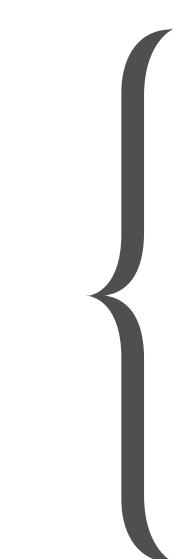
B. Point estimate + 95% confidence interval
The point estimate and 95% confidence interval from (A) represented graphically. Graphical representation of statistical results have been commonly advocated as another way of emphasizing uncertainty.

C. Density + interval plot (here 95%)
The density shows the sampling distribution (or Bayesian posterior distribution). The addition of the point estimate and interval adds precision.

D. Quantile dotplot (here 100 dots)
A quantile dotplot allows precise estimation of many intervals and by providing a frequency based framing of the probability of the parameter, may improve understanding of uncertainty through hypothetical outcomes.

A non exhaustive set of statements describing a statistical result

There exists

- 
- 
- strong evidence
 - weak evidence
 - inconclusive evidence
 - no evidence

that an effect exists

A non exhaustive set of statements describing a statistical result

There exists

- 
- { - strong evidence
- weak evidence
- inconclusive evidence
- no evidence }

that an effect exists

“The lure of incredible certitude”

Existing incentives make it tempting for researchers to maintain assumptions far stronger than they can persuasively defend, in order to draw strong conclusions.

**Let's step back from
strictly probabilistic
uncertainty.**

Garden of forking paths

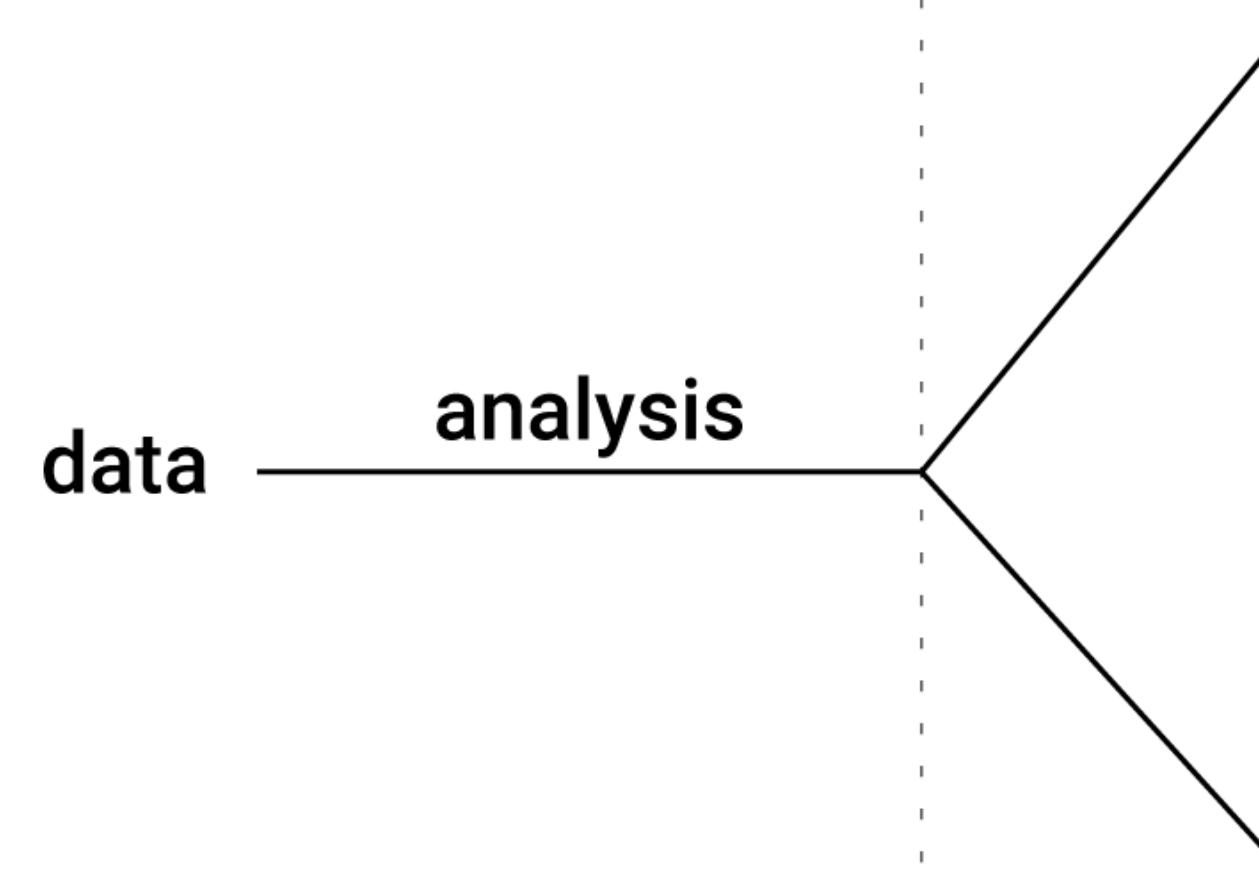
[Gelman and Loken (2016)]

data $\xrightarrow{\text{analysis}}$ $p < 0.05$

Garden of forking paths

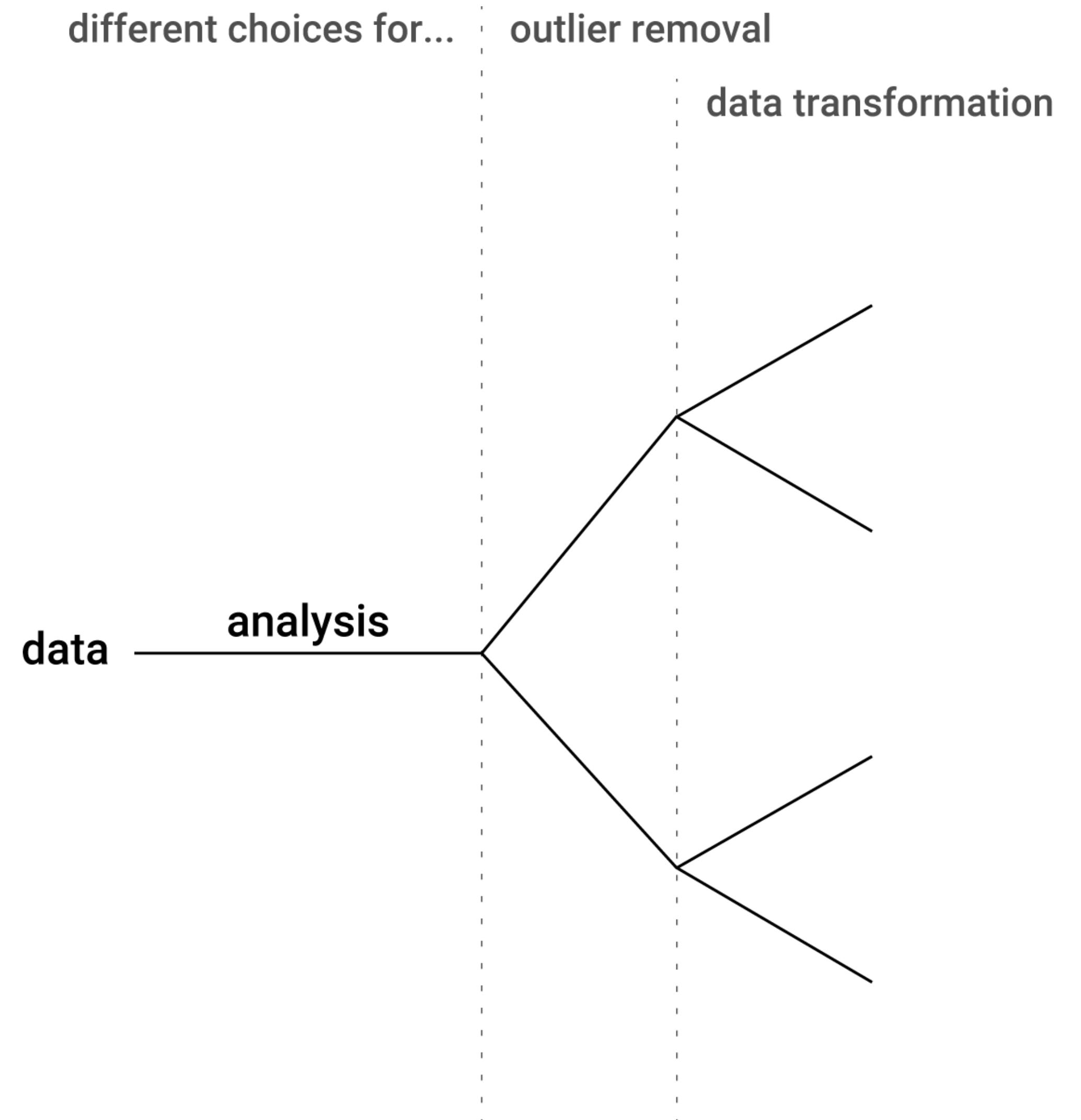
[Gelman and Loken (2016)]

different choices for... | outlier removal



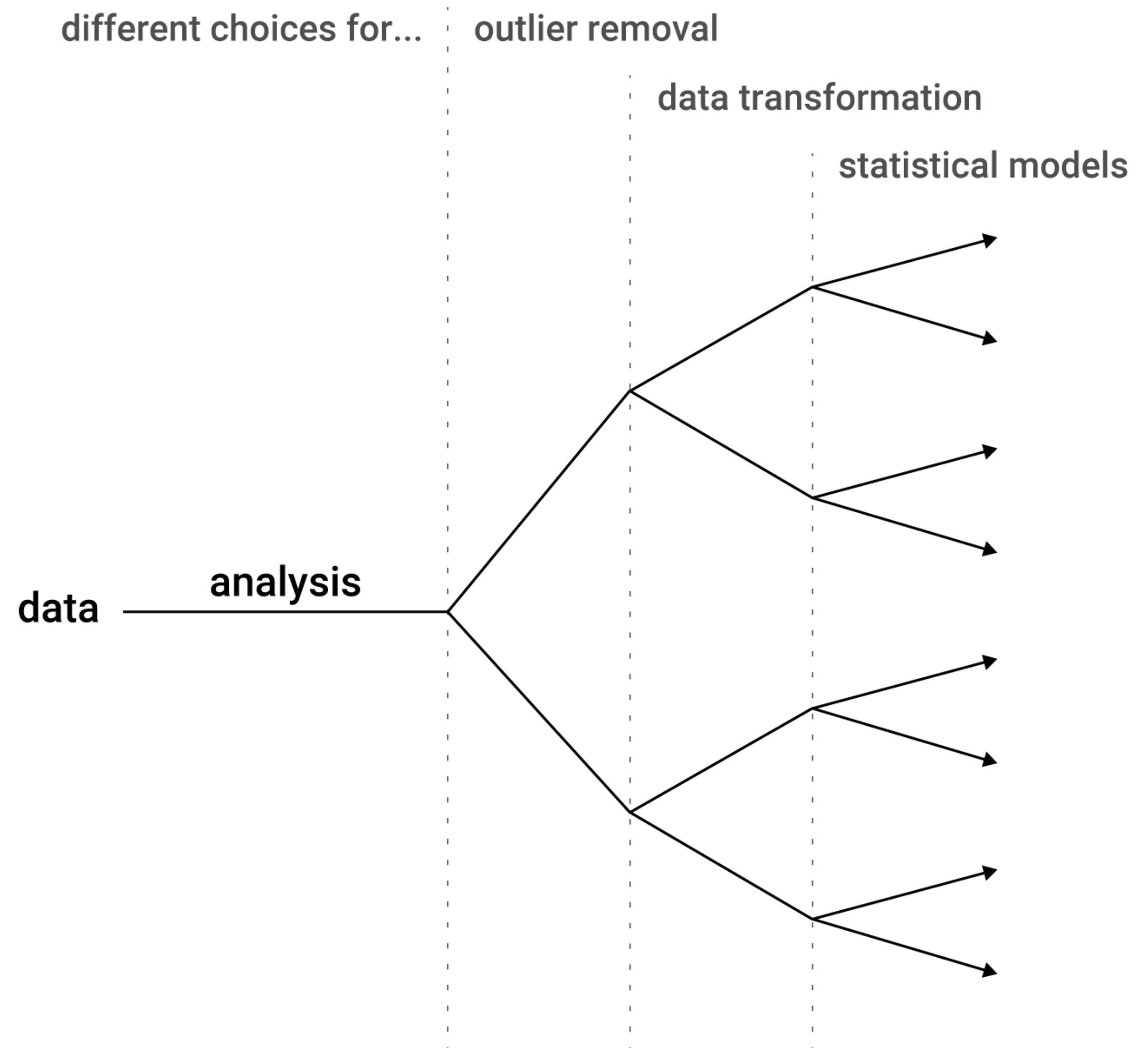
Garden of forking paths

[Gelman and Loken (2016)]



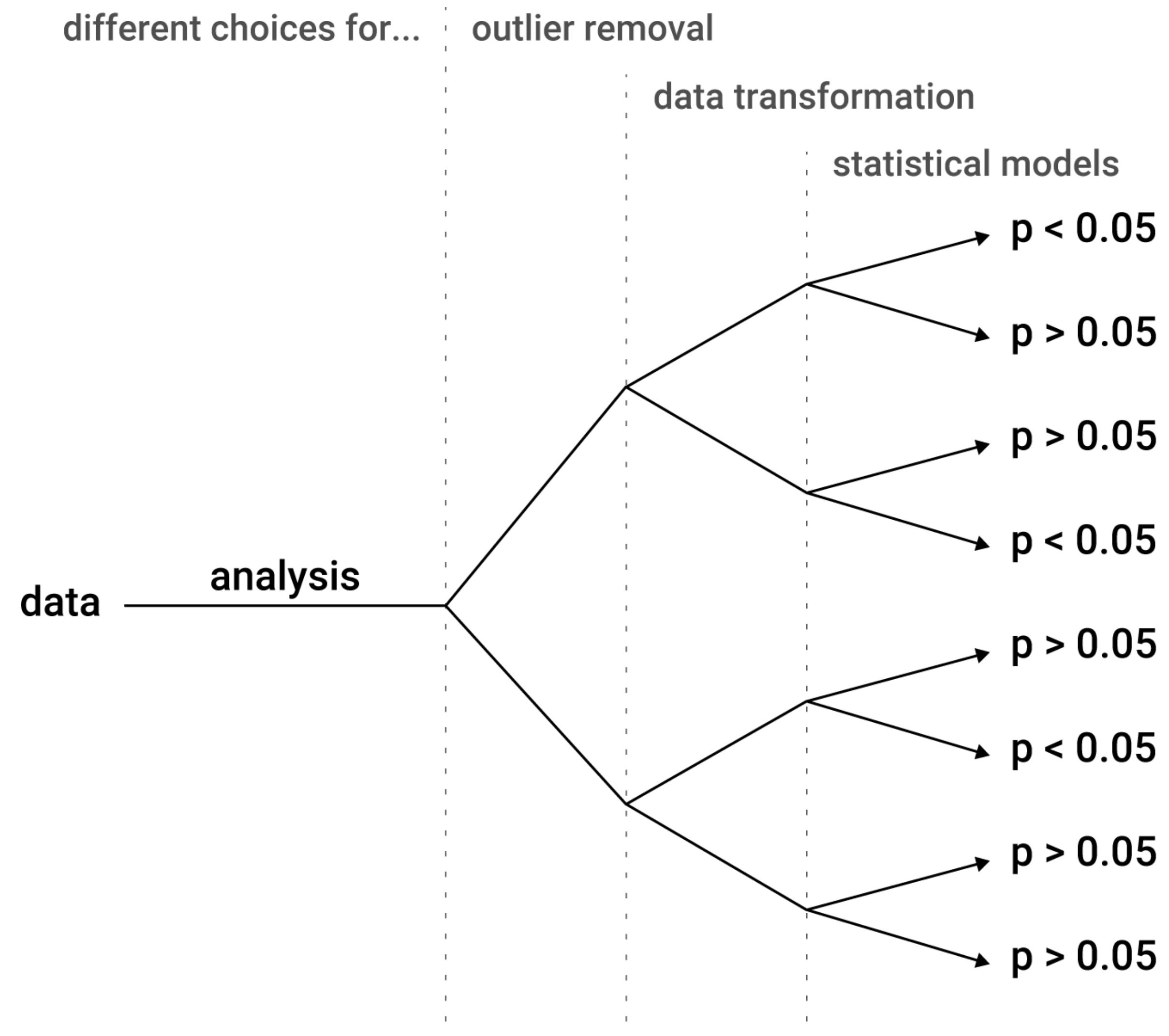
Garden of forking paths

[Gelman and Loken (2016)]



Garden of forking paths

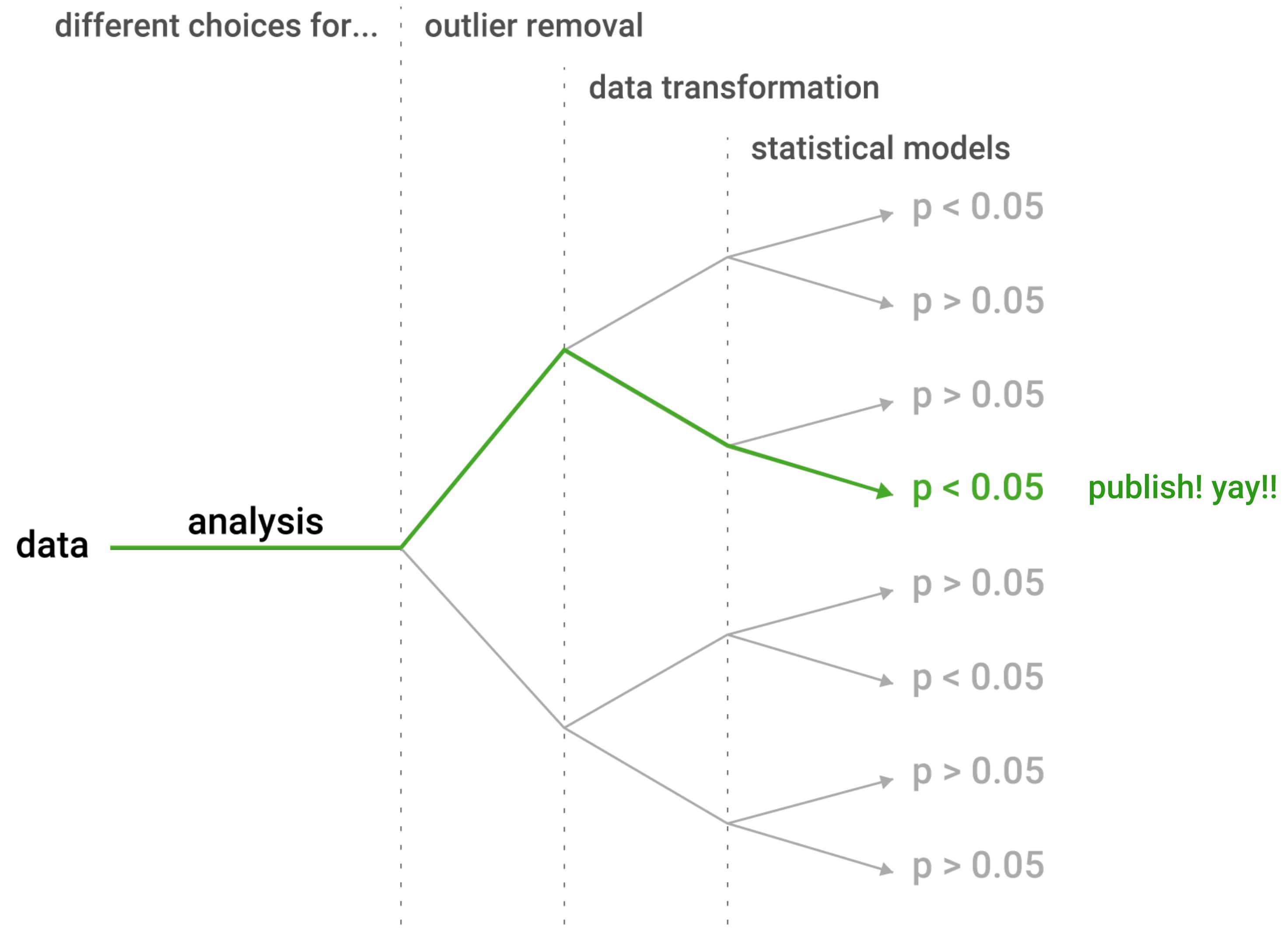
[Gelman and Loken (2016)]



This is
model/specification
uncertainty

Garden of forking paths

[Gelman and Loken (2016)]



[Christie Aschwanden
and Ritchie King (2015):
Science Isn't Broken in
FiveThirtyEight]

Hack Your Way To Scientific Glory



You're a social scientist with a hunch: **The U.S. economy is affected by whether Republicans or Democrats are in office.** Try to show that a connection exists, using real data going back to 1948. For your results to be publishable in an academic journal, you'll need to prove that they are "statistically significant" by achieving a low enough p-value.

1 CHOOSE A POLITICAL PARTY

Republicans

Democrats

2 DEFINE TERMS

Which politicians do you want to include?

- Presidents
- Governors
- Senators
- Representatives

How do you want to measure economic performance?

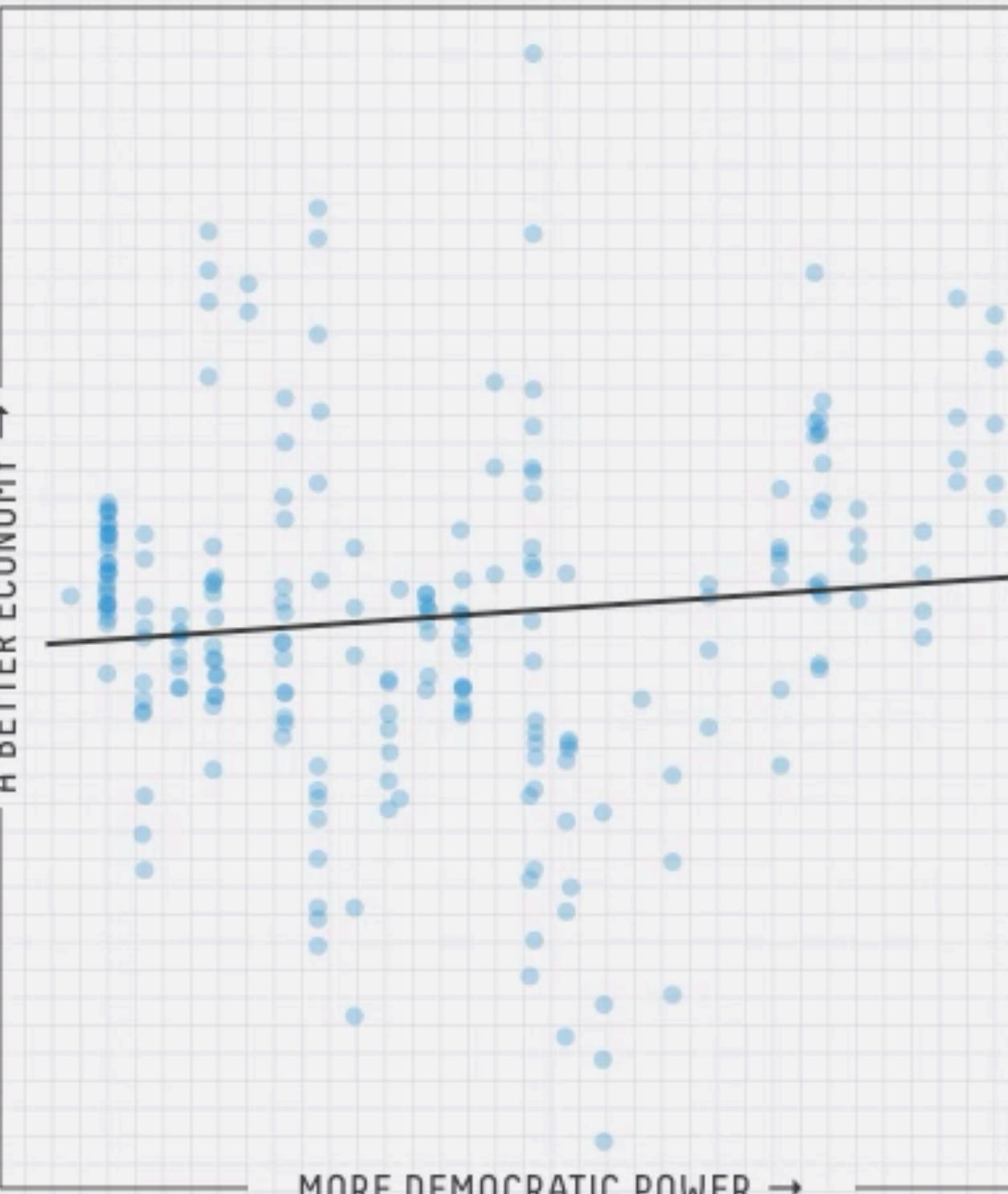
- Employment
- Inflation
- GDP
- Stock prices

Other options

- Factor in power
Weight more powerful positions more heavily
- Exclude recessions
Don't include economic recessions

3 IS THERE A RELATIONSHIP?

Given how you've defined your terms, does the economy do better, worse or about the same when more Democrats are in power? Each dot below represents one month of data.



4 IS YOUR RESULT SIGNIFICANT?

If there were no connection between the economy and politics, what is the probability that you'd get results at least as strong as yours? That probability is your p-value, and by convention, you need a **p-value of 0.05 or less** to get published.



Result: Almost

Your **0.06** p-value is close to the 0.05 threshold. Try tweaking your variables to see if you can push it over the line!

If you're interested in reading real (and more rigorous) studies on the connection between politics and the economy, see the work of Larry Bartels and Alan Blinder and Mark Watson.

Data from The @unitedstates Project, National Governors Association, Bureau of Labor Statistics, Federal Reserve Bank of St. Louis and Yahoo Finance.



**Do hurricanes with more feminine
names cause more deaths?**

Female hurricanes are deadlier than male hurricanes

Kiju Jung^{a,1}, Sharon Shavitt^{a,b,1}, Madhu Viswanathan^{a,c}, and Joseph M. Hilbe^d

^aDepartment of Business Administration and ^bDepartment of Psychology, Institute of Communications Research, and Survey Research Laboratory, and ^cWomen and Gender in Global Perspectives, University of Illinois at Urbana-Champaign, Champaign, IL 61820; and ^dDepartment of Statistics, T. Denny Sanford School of Social and Family Dynamics, Arizona State University, Tempe, AZ 85287-3701

Edited* by Susan T. Fiske, Princeton University, Princeton, NJ, and approved May 14, 2014 (received for review February 13, 2014)

Do people judge hurricane risks in the context of gender-based expectations? We use more than six decades of death rates from US hurricanes to show that feminine-named hurricanes cause significantly more deaths than do masculine-named hurricanes. Laboratory experiments indicate that this is because hurricane names lead to gender-based expectations about severity and this, in turn, guides respondents' preparedness to take protective action. This finding indicates an unfortunate and unintended consequence of the gendered naming of hurricanes, with important implications for policymakers, media practitioners, and the general public concerning hurricane communication and preparedness.

gender stereotypes | implicit bias | risk perception | natural hazard communication | bounded rationality

Festimates suggest that hurricanes kill more than 200 people in

violence and destruction (23, 24). We extend these findings hypothesize that the anticipated severity of a hurricane with a masculine name (Victor) will be greater than that of a hurricane with a feminine name (Victoria). This expectation, in turn, will affect the protective actions that people take. As a result, a hurricane with a feminine vs. masculine name will lead to less protective action and more fatalities.

Archival Study

To test this hypothesis, we used archival data on actual fatalities caused by hurricanes in the United States (1950–2012). Nine of the four Atlantic hurricanes made landfall in the United States during this period (25). Nine independent coders who were blind to the hypothesis rated the masculinity vs. femininity of historical hurricane names on two items (1 = very masculine, 11 = very feminine, and 1 = very man-like, 11 = very woman-like), which

“If you torture the data long enough, it will confess.”

- Ronald Chase

Female hurricanes are not deadlier than male hurricanes

Jung et al. (1) assert that hurricanes that made landfall in the United States killed more people when they had female names rather than male names. The article has stirred much controversy. Criticisms range from the inclusion of hurricanes from the era before they were given male names (2) to collective interpretation and the overstatement of their results from the archival test of their hypothesis (3), to the exclusion of their six behavioral experiments on populations in at-risk situations (4). In sum, of this letter,

The criticism of this one: the results of their archi function of the selective incl sors. Using the same data, m variables, I show in Table 1 are not robust to the incl two-way interaction they on analysis. Model 1 reprodu main results. Models 2-4 s that female- and male-n were equally deadly cannot the interaction effect of a metric pressure and its r toll is included. A more a letter should have stated

Models 2-4 show lower barometric pressure tolls and that hurricanes

Table 1. Results from

tolls had smaller death tolls when the hurricanes were strong (lower pressure), but higher death tolls when the hurricanes were weak (higher pressure). The latter result is driven by the pre-1978 sample (model 5). In the post-1978 sample, the interaction effect is insignificant and the damage toll relationship

LETTER

Are female hurricanes really more deadly than male hurricanes?

The reasoning in ref. 1 is fundamentally based on the regression models reported in their table S2, in particular, model 4. However, due to the interaction terms combined with extreme values and weak significance, the analysis is based on a very fragile model; e.g., the model predicts almost 20,000 deaths for hurricane Sandy, which actually caused 158 fatalities. This is far from a plausible

ur claim that the differences between male- or female-named hurricanes for deaths, minimum pressure, category, and damages.

To conclude, the analyses given in ref. 1 are examples of the fact that prediction models using interaction terms have to be handled and interpreted carefully; in particular, using insignificant variables is not expedient and may lead to statistical artifacts.

To summarize, the data do not contain evidence that feminine-named hurricanes cause more deaths than masculine-named hurricanes.

**Björn Christensen^a and
Søren Christensen^{b,1}**



ELSEVIER

Weather and Climate Extremes 12 (2016) 80–84
Contents lists available at ScienceDirect

Contents lists available at ScienceDirect
Extremes 12 (2016) 80

Weather and Climate Extremes

Journal homepage: www.elsevier.com/locate/watex

Hurricane names: A bunch of hot air

Gary Smit

Department of Economics, Pomona College, United States

ARTICLE I



 CrossMark
click for updates

100

Are female hurricanes really deadlier than male hurricanes?

Jung et al. (1) claim to show that “feminine-named hurricanes cause significantly more deaths than do masculine-named hurricanes” (p. 1). This conclusion is mainly obtained by analyzing data on fatalities caused by hurricanes in the United States (1950–2012). By reanalyzing the same data, we show that the conclusion is based on biased presentation and invalid statistics.

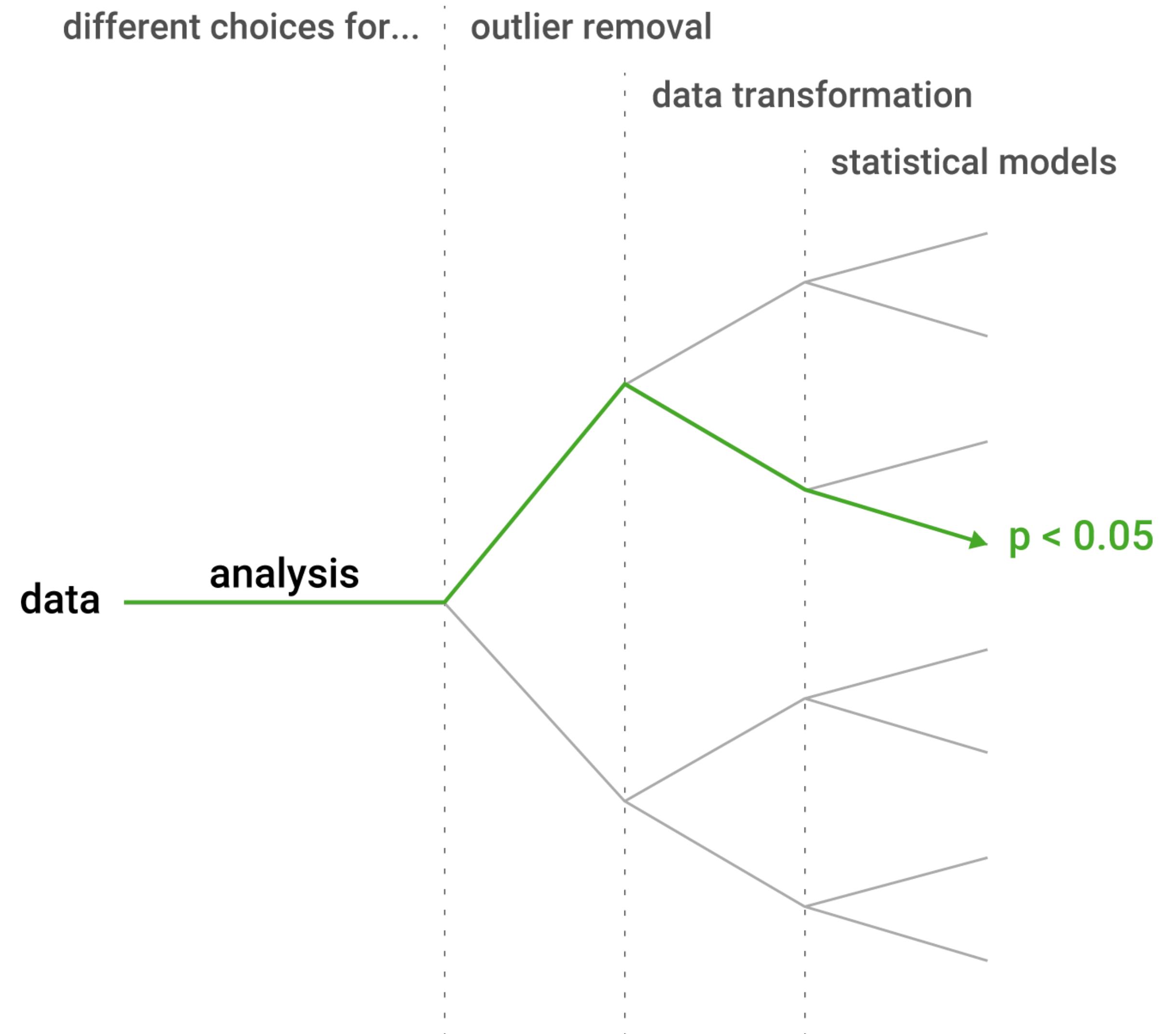
The reasoning in ref. 1 is fundamentally based on the regression models reported in their table S2, in particular, model 4. However, due to the interaction terms combined with extreme values and weak significance, the analysis is based on a very fragile model; e.g., the model predicts almost 20,000 deaths for hurricane Sandy, which actually caused 150 fatalities. This is far from a plausible result.

Now, we explain our claim that the results are presented in a biased way. By holding the minimum pressure at its mean in prediction of counts of deaths, the authors only report the influence of MFI and normalized damage (figure 1 in ref. 1). This ignores the influence of the second interaction term MFI minimum pressure, which shows an opposite influence (see the estimated parameters on p. 11 first paragraph). By considering the counts of deaths under constant normalized damage, the results are contrary: male-named hurricanes with a low minimum pressure (strong hurricanes) are associated with more deaths than female ones (Fig. 1).

In the light of an alternating male-female

-named hurricanes are deadlier because people do not take them seriously based on a questionable statistical analysis of a narrowly defined dataset not robust in that it is not confirmed by a straightforward analysis of data. B.V. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

Pre-registration



A referee in a black shirt and dark shorts stands on a soccer field, holding a red card high in his right hand. He is looking towards the camera. In the background, several players are visible, including one in a red and white jersey with 'MAN UTD' on it, and others in light blue jerseys with 'DE JONG 34' and 'KOMPANY 4' on them. The scene is set at night with stadium lights illuminating the field.

Do soccer referees give more red cards to dark-skinned players than light-skinned ones?

**Different researchers
may create very different
pre-registration
documents**

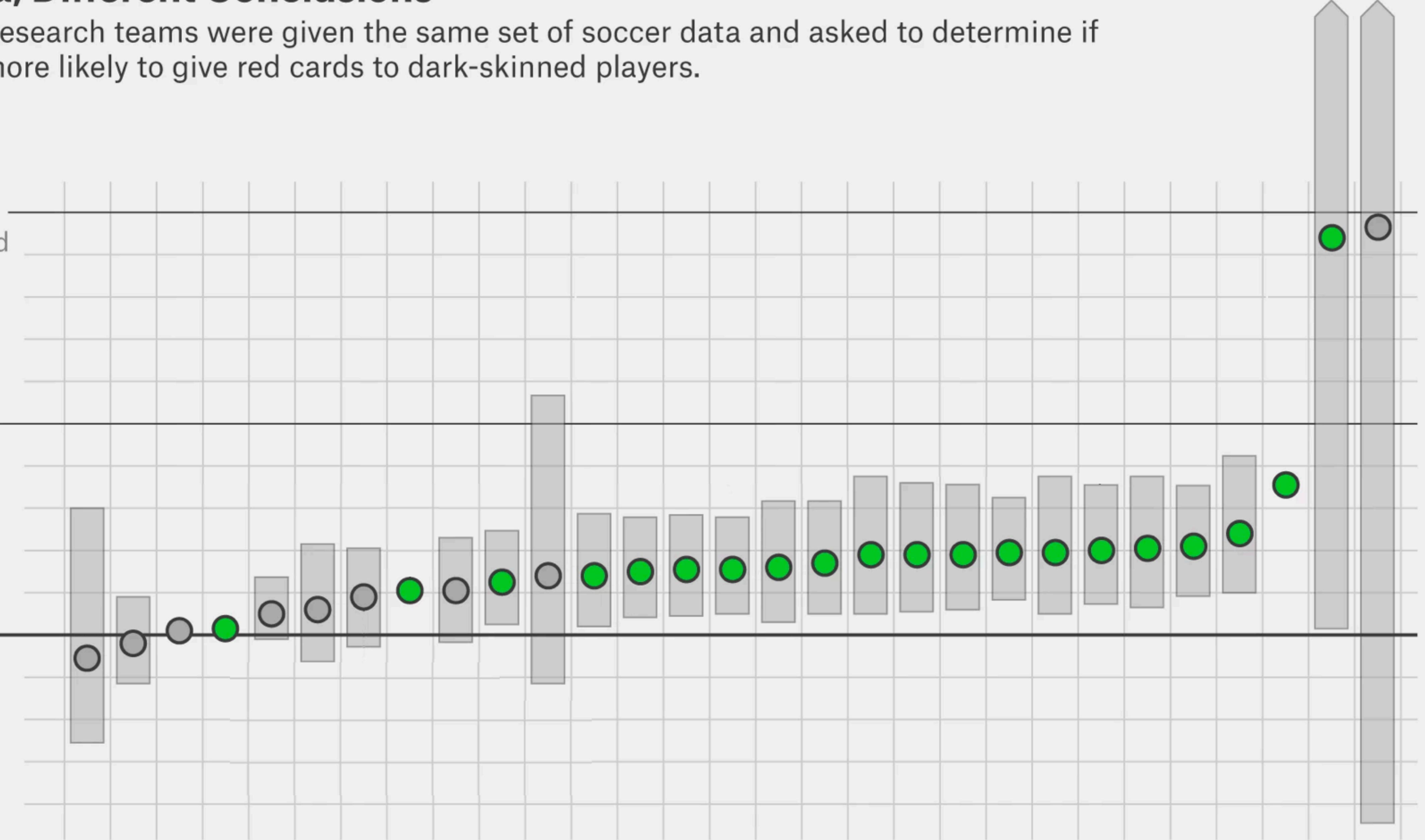
Same Data, Different Conclusions

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.

Referees are
three times as
likely to give red
cards to
dark-skinned
players

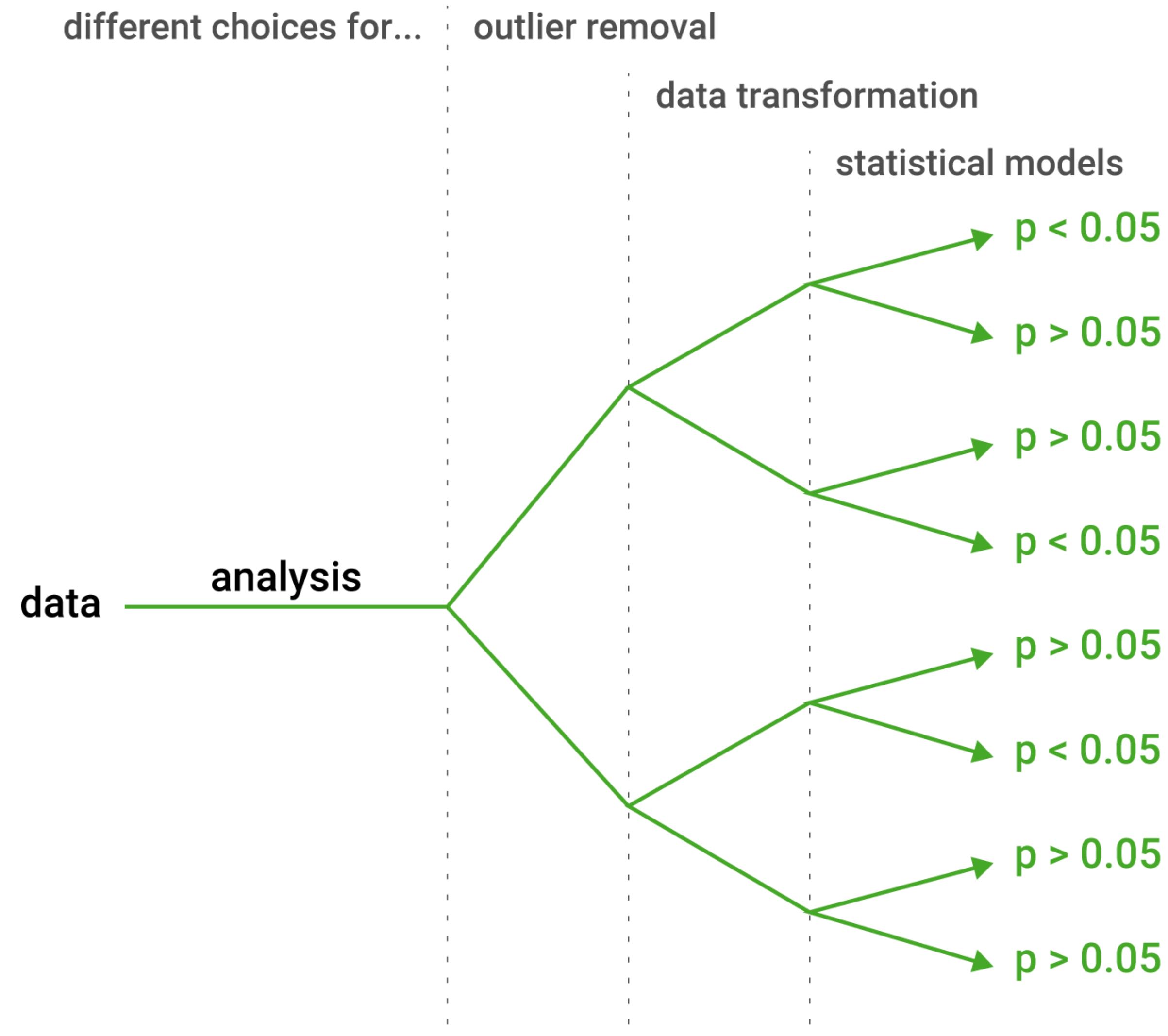
Twice as likely

Equally likely



Multiverse analysis

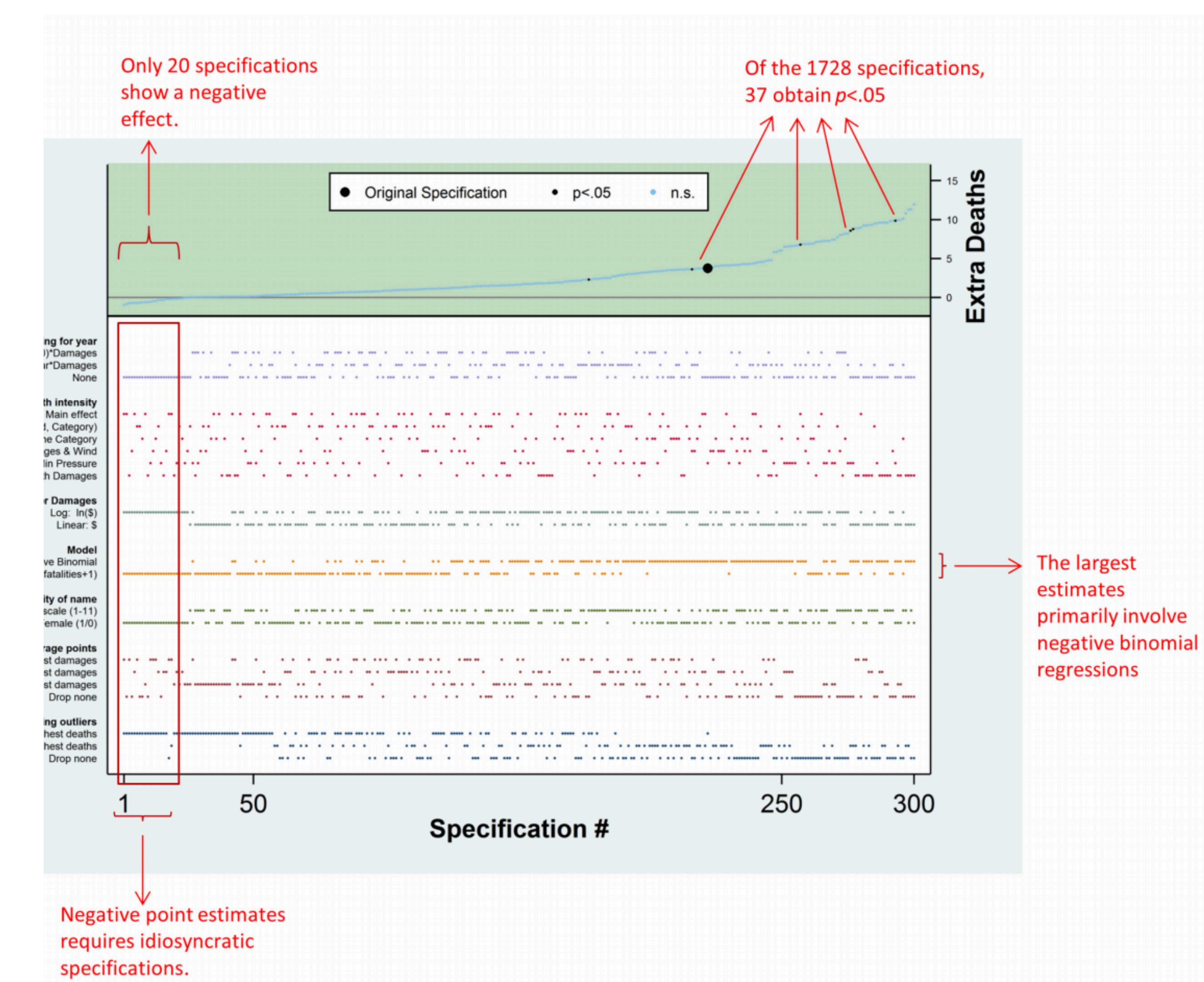
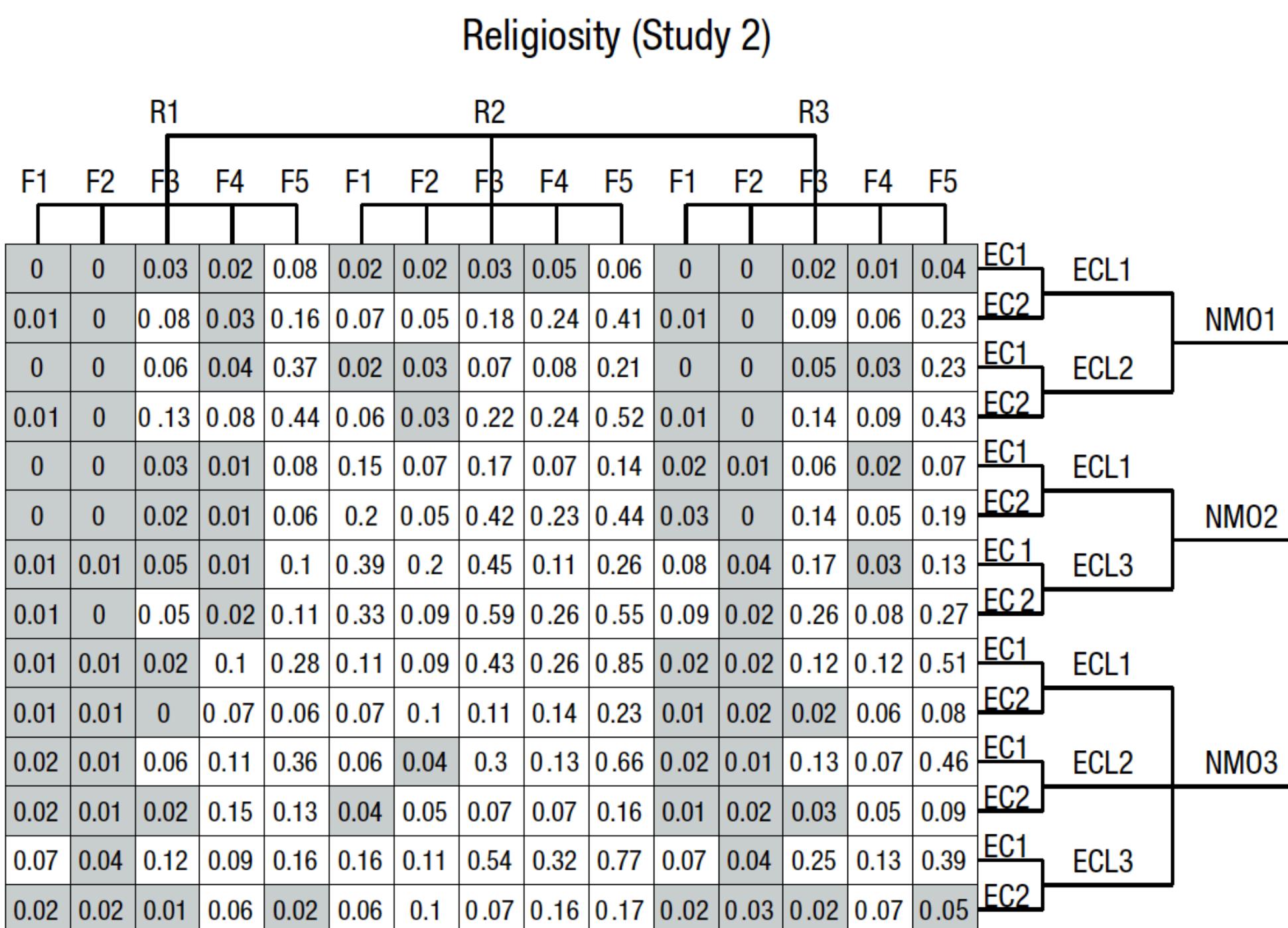
[Steegen, Tuerlinckx, Gelman, and Vanpaemel (2016):
Increasing Transparency Through a Multiverse Analysis]



**Performing and
reporting all
reasonable analysis
scenarios.**

How to report all of
these analyses?

Visual summaries?



[Steegen, Tuerlinckx, Gelman, and Vanpaemel and Loken (2016): Increasing Transparency Through a Multiverse Analysis]

[Simonsohn, Simmons, and Nelson (2015)
Specification curve: Descriptive and inferential statistics on all reasonable specifications]

Can we do better?

Explorable Multiverse Analysis Reports (EMARs)

[Dragicevic, Jansen, **Sarma**, Kay, and Chevalier (2019): Increasing the Transparency of Research Papers with Explorable Multiverse Analyses]

Re-Evaluating the Efficiency of Physical Visualizations: A Simple Multiverse Analysis

Pierre Dragicevic

Inria

pierre.dragicevic@inria.fr

Yvonne Jansen

CNRS & Sorbonne Université

jansen@isir.upmc.fr

ABSTRACT

A previous study has shown that moving 3D data visualizations to the physical world can improve users' efficiency at information retrieval tasks. Here, we reanalyze a subset of the experimental data using a multiverse analysis approach. Results from this multiverse analysis are presented as explorable explanations, and can be interactively explored in this paper. The study's findings appear to be robust to choices in statistical analysis.

AUTHOR KEYWORDS

Physical visualization; multiverse analysis.

ACM CLASSIFICATION KEYWORDS

H5.2 User Interfaces: Evaluation/Methodology

GENERAL TERMS

Human Factors; Design; Experimentation; Measurement.

INTRODUCTION

Whereas traditional visualizations map data to pixels or ink, physical visualizations (or "data physicalizations")

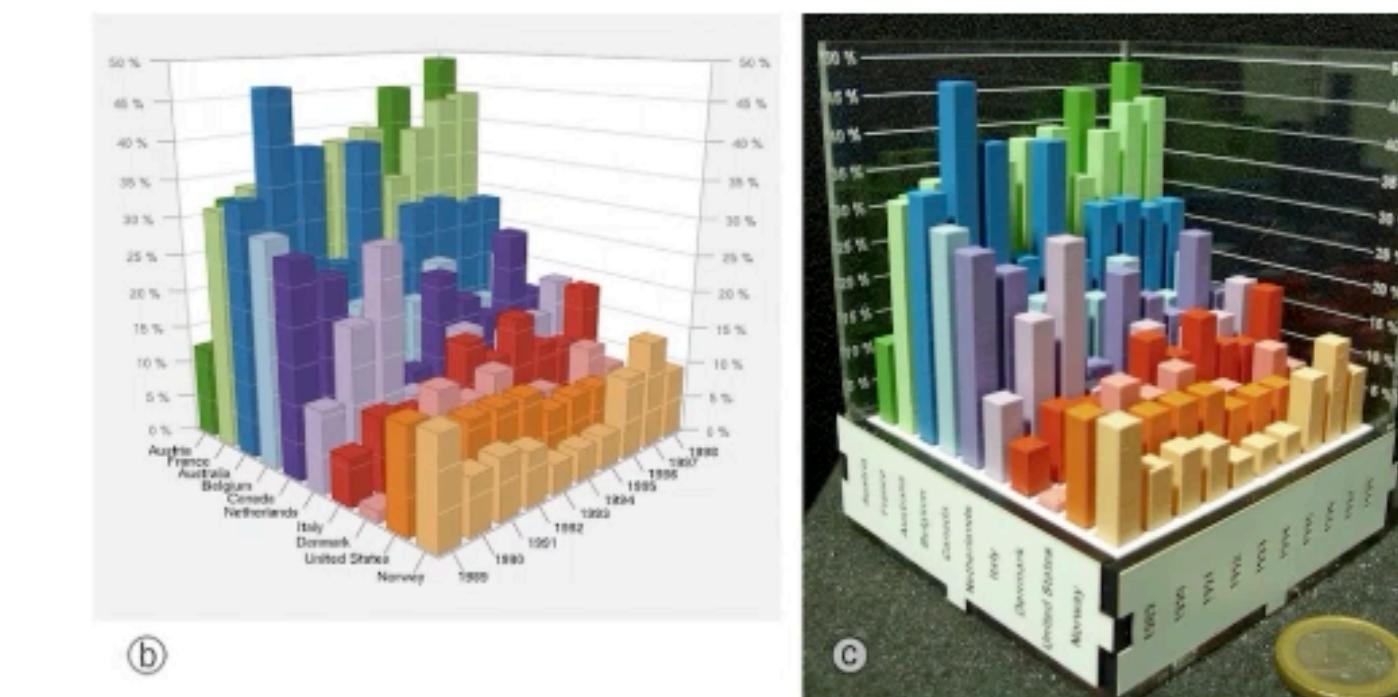


Figure 1. 3D bar chart, on-screen and physical.

STUDY

The study consisted of two experiments. In the first experiment, participants were presented with 3D bar charts showing country indicator data, and were asked simple questions about the data. The 3D bar charts were presented both on a screen and in physical form (see Figure 1). The on-screen bar chart could be rotated in all directions with the mouse. Both a regular and a stereoscopic display were tested. An interactive 2D bar chart was also used as a control condition. Accuracy was high across all conditions,

An Explorable Multiverse Analysis of Durante et al. (2013)

Pierre Dragicevic

Inria

pierre.dragicevic@inria.fr

ABSTRACT

In this paper, we reproduce a small part of Steegen et al.'s multiverse analysis of Durante et al.'s study using explorable explanations. The data processing options can be selected interactively, which allows us to show the interaction plot reported in Durante et al. in addition to the p-value.

AUTHOR KEYWORDS

Multiverse analysis.

ACM CLASSIFICATION KEYWORDS

H5.2 User Interfaces: Evaluation/Methodology

GENERAL TERMS

Human Factors; Design; Experimentation; Measurement.

INTRODUCTION

Steggen and colleagues [1] introduced the concept of *multiverse analysis*, which they illustrated by re-analyzing data from a 2013 paper by Durante and colleagues [2] entitled “The fluctuating female vote: Politics, religion, and the ovulatory cycle”. Here, we report the initial part of the

- days 9–17 for high fertility and 18–25 for low fertility [5],
- days 8–14 for high fertility and 1–7 and 15–28 for low fertility [6], and
- days 9–17 for high fertility and 1–8 and 18–28 for low fertility [7].

Second, there are different reasonable ways of estimating a woman's next menstrual onset, which is an intermediate step in determining cycle day. ☒ A woman's cycle day can be based on the number of days before next menstrual onset, which in turn is based on cycle length, which is computed as the difference between the start date of the woman's last menstrual period and the start date of the woman's previous menstrual period [2]. ☐ Another way to estimate next menstrual onset is based on the women's reported estimate of their typical cycle length [8].

Relationship status

There are at least three options for the dichotomization of women's relationship status into single or committed.

- ☒ Women who selected response Option 1 or 2 on the relationship status item can be assigned to the group of single women, whereas women who selected response

All these and more
examples can be found at:

explorablemultiverse.github.io/

**We need to promote and
support transparent
statistical reporting**

Thanks!

And thanks to Matt Kay, Pierre Dragicevic,
Yvonne Jansen, Fanny Chevalier



abhsarma.github.io



@abhsarma

Abhraneel Sarma
University of Michigan
School of Information



mucollective.co