

False-Positives, p-Hacking, Statistical Power, and Evidential Value

Leif D. Nelson

**University of California, Berkeley
Haas School of Business**

**Summer Institute
June 2014**



Who am I?

- Experimental psychologist who studies judgment and decision making.
 - And has interests in methodological issues

Who are you?

[not a rhetorical question]

- Grad Student vs. Post-Doc vs. Faculty?
- Psychology vs. Economics vs. Other?
- Have you read any papers that I have written?
 - Really? Which ones?

Things I want you to get out of this

- It is quite easy to get a false-positive finding through p-hacking. (5%)
- Transparent reporting is critical to improving scientific value. (5%)
- It is (very) hard to know how to correctly power studies, but there is no such thing as overpowering. (30%)
- You can learn a lot from a few p-values. (remainder %)

**This will be most helpful to you if you
ask questions.**

**A discussion will be more interesting
than a lecture.**

SLIDES ABOUT P-HACKING

False-Positives are Easy

- It is common practice in all sciences to report less than everything.
 - So people only report the good stuff. We call this *p*-Hacking.
 - Accordingly, what we see is too “good” to be true.
 - We identify six ways in which people do that.

Six Ways to p-Hack

1. Stop collecting data once $p < .05$
2. Analyze many measures, but report only those with $p < .05$.
3. Collect and analyze many conditions, but only report those with $p < .05$.
4. Use covariates to get $p < .05$.
5. Exclude participants to get $p < .05$.
6. Transform the data to get $p < .05$.

OK, but does that matter very much?

- As a field we have agreed on $p < .05$. (i.e., a 5% false positive rate).
- If we allow p-hacking, then that false positive rate is actually 61%.
- Conclusion: p-hacking is a potential catastrophe to scientific inference.

P-Hacking is Solved Through Transparent Reporting

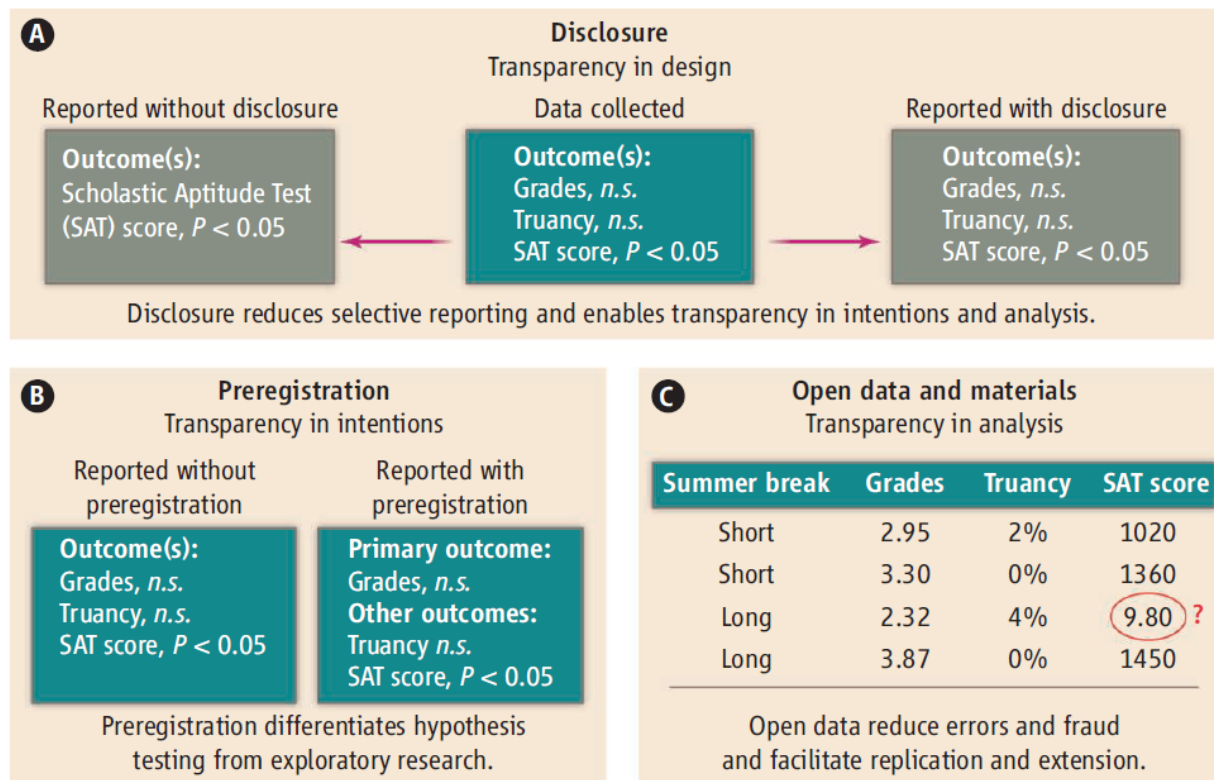
- Instead of reporting only the good stuff, just report all the stuff.

P-Hacking is Solved Through Transparent Reporting

- Solution 1:
 1. Report sample size determination.
 2. $N > 20$ [note: I will tell you later about how this number is insanely low. Sorry. Our mistake.]
 3. List all of your measures.
 4. List all of your conditions.
 5. If excluding, report without exclusion as well.
 6. If covariates, report without.

P-Hacking is Solved Through Transparent Reporting

- Solution 2:



Three mechanisms for increasing transparency in scientific reporting. Demonstrated with a research question: "Do shorter summer breaks improve educational outcomes?" *n.s.* denotes $P > 0.05$.

P-Hacking is Solved Through Transparent Reporting

- Implications:
 - Exploration is necessary; therefore replication is as well.
 - Without p-hacking, fewer significant findings; therefore fewer papers.
 - Without p-hacking, need more power; therefore more participants.

SLIDES ABOUT POWER

Motivation

- With p -hacking,
 - statistical power is irrelevant, most studies work
- Without p -hacking.
 - take power seriously, or most studies fail
- Reminder. Power analysis:
 - **Guess** effect size (d)
 - **Set** sample size (n)
- Our question: Can we make guessing d easier?
- Our answer: **No**
- Power analysis is not a practical way to take power seriously

How to guess d ?

- Pilot

- Prior literature

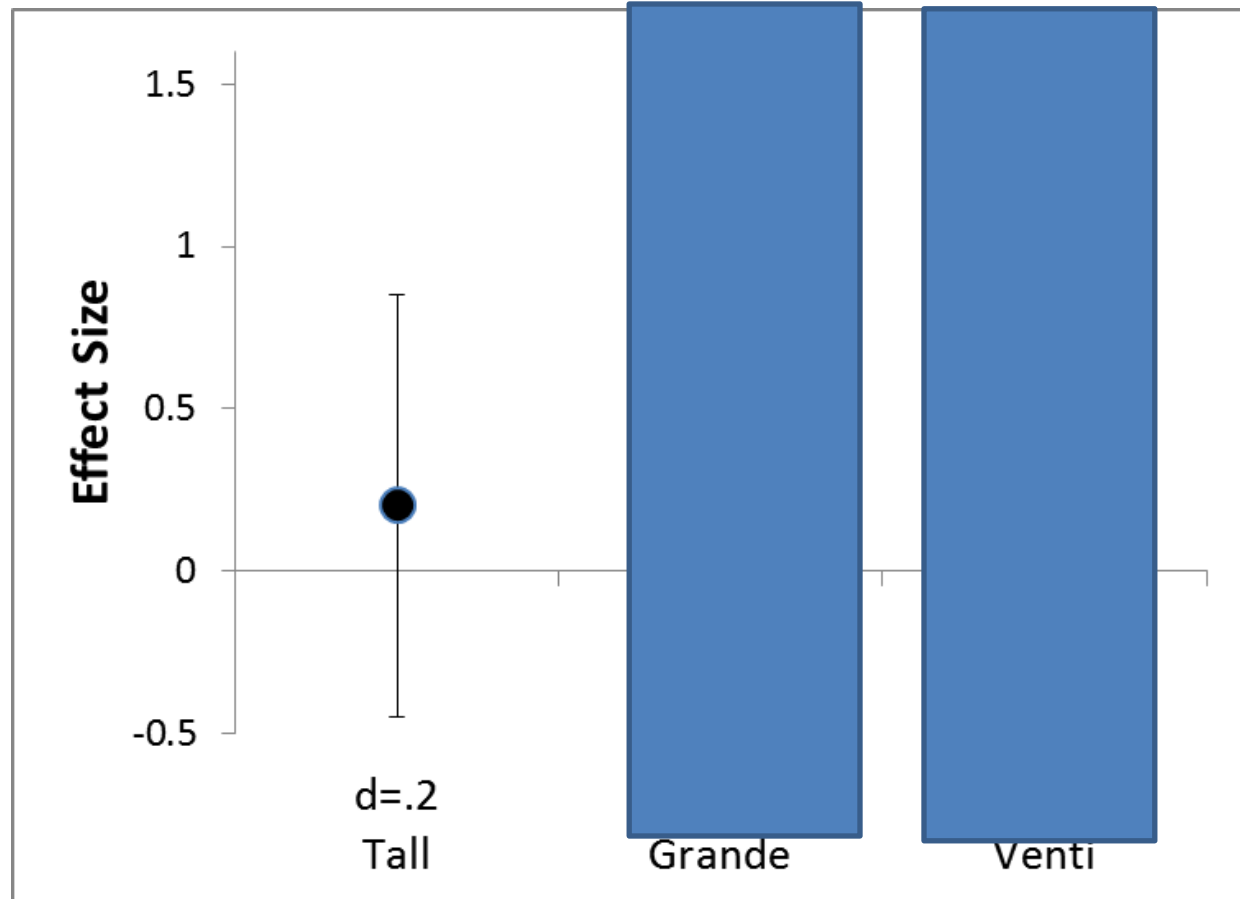
- Theory/gut

Some kind words before the bashing

- Pilots: They are good for:
 - Do participants get it?
 - Ceiling effects?
 - Smooth procedure?
- Kind words end here.

Pilots: useless to set sample size

- Say Pilot: $n=20$
 - $\hat{d} = .2$
 - $\hat{d} = .5$
 - $\hat{d} = .8$



- **In words**
 - Estimates of d have too much sampling error.
- **In more interesting words**
 - Next.

Think of it this way

Say in actuality you need $n=75$

Run Pilot: $n=20$

What will Pilot say you need?

- Pilot 1: "you need $n=832$ "
- Pilot 2: "you need $n=53$ "
- Pilot 3: "you need $n=96$ "
- Pilot 4: "you need $n=48$ "
- Pilot 5: "you need $n=196$ "
- Pilot 6: "you need $n=10$ "
- Pilot 7: "you need $n=311$ "

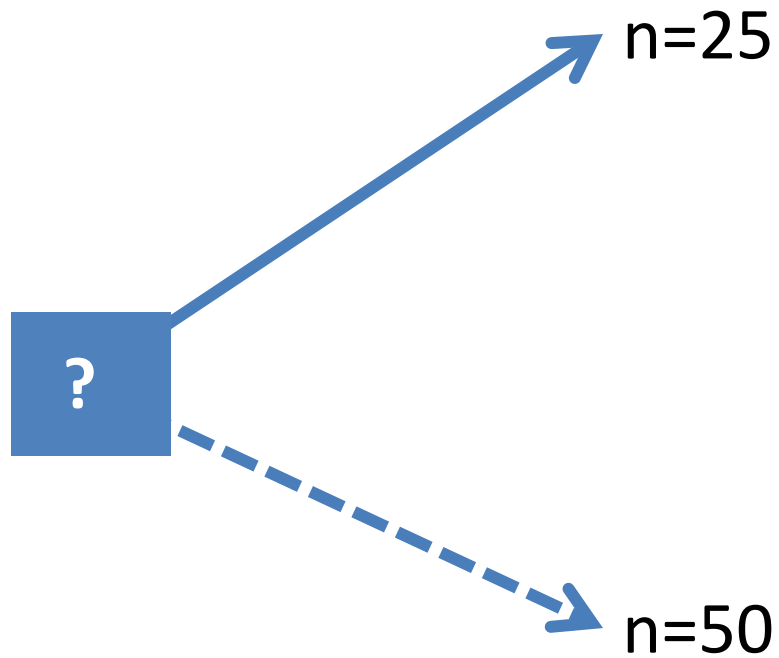
Thanks Pilot!

$n=20$ is not enough.

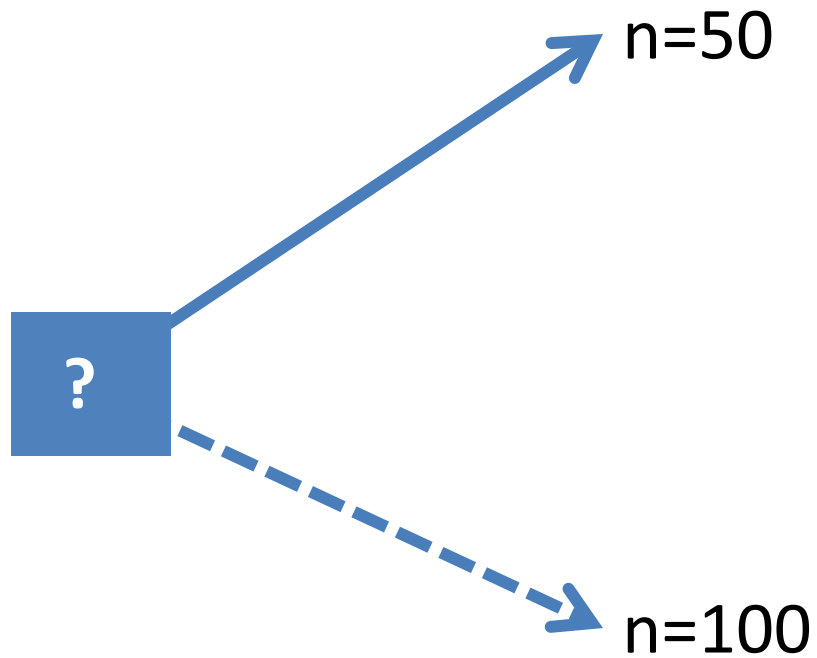
How many subjects do you need

to know

how many subjects you need?

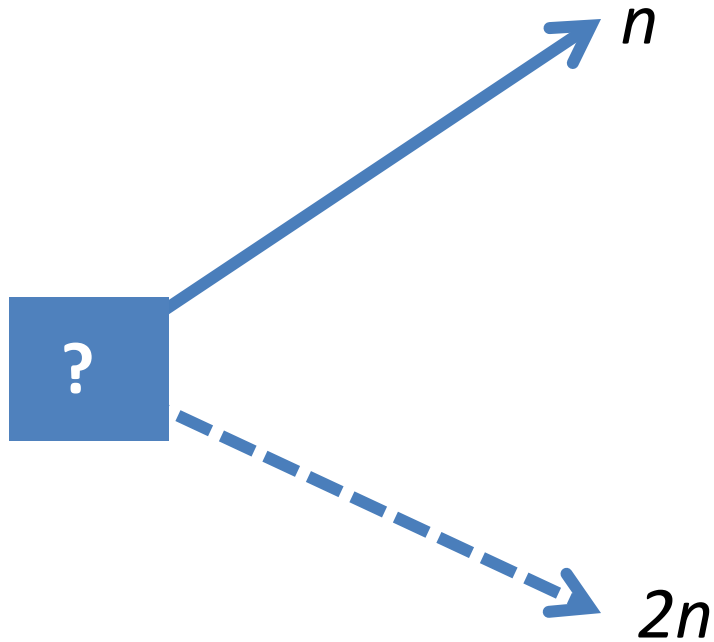


Need a Pilot with...
 $n=133$



Need a Pilot with...
 $n=276$

“Theorem” 1



Need: $5n$

How to guess d ?

- Pilot
- Existing findings
- Theory/gut

Existing findings

- **One hand**
 - Larger samples
- **Other hand**
 - Publication bias
 - More noise
 - \neq sample
 - \neq design
 - \neq measures

Best (im)possible case scenario

- Would guessing d be reasonable based on other studies?

“Many Labs” Replication Project

- Klein et al.,
- 36 labs
- 12 countries
- $N=6344$
- Same 13 experiments

Open Science Framework^{BETA}

Explore ▾

Help ▾

Search

Create an Account or Sign-In

Investigating variation in replicability: The “Many Labs” Replication Project

Public

Watch 5

0

Contributors: [Richard A. Klein](#), [Kate Ratliff](#), [Brian A. Nosek](#), [Michelangelo Vianello](#), [Ronaldo Pilati](#), [Thierry Devos](#), [Elisa Maria Galliani](#), [Mark Brandt](#), [Anna van 't Veer](#), [Abraham M. Rutchick](#), [Kathleen Schmidt](#), [Stepan Bahnik](#), [Marek Vranka](#), [Hans IJzerman](#), [Fred Hasselman](#), [Jennifer Joy-Gaba](#), [Jesse J. Chandler](#), [Leigh Ann Vaughn](#), [Claudia Brumbaugh](#), [Lyn van der Wal](#), [Aaron Wichman](#), [Grant Packard](#), [Beach Brooks](#), [Zeynep Cemalcilar](#), [Justin Storbeck](#), [Konrad Bocian](#), [Carmel Levitan](#), [Michael Jason Bernstein](#), [Lacy Elise Krueger](#), [Matthew Eisner](#), [William E. Davis](#), [Jason A. Nier](#), [Anthony J. Nelson](#), [Troy G. Steiner](#), [Robyn Mallett](#), [Donna Thompson](#), [Jeffrey R. Huntsinger](#), [Wendy Morris](#), [Jeanine Skorinko](#), [Heather Kappes](#)

Date Created: 6/14/2013 3:29 PM | Last Updated: 12/20/2013 2:28 PM

Description: We conducted replications of 13 effects in psychological science with 36 samples and more than 6000 participants. We examined heterogeneity in replicability across sample and setting.

How much TV per day?

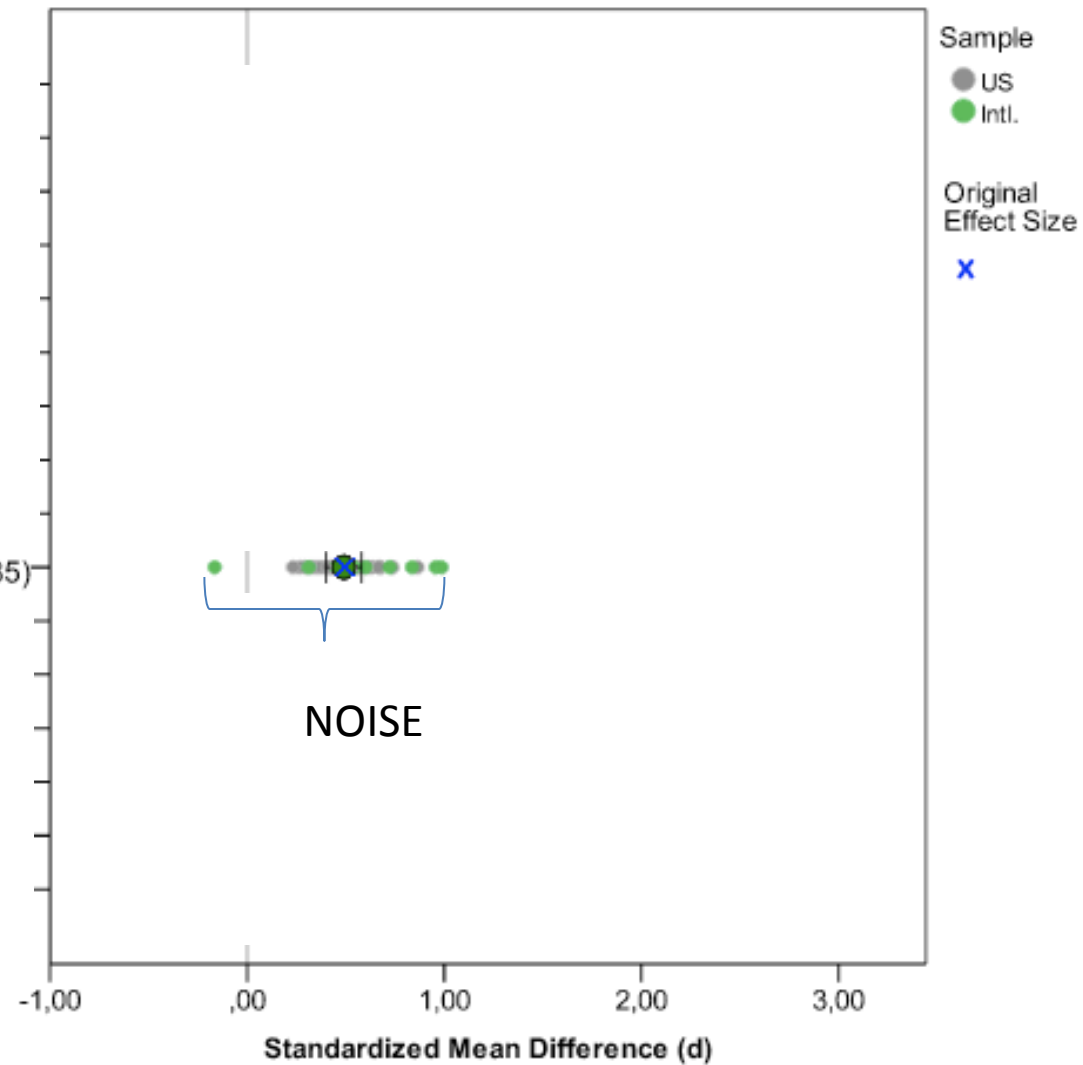
Low-frequency alternatives

Up to ½ hour
 ½ hour to 1 hour
 1 hour to 1½ hours
 1½ hours to 2 hours
 2 hours to 2½ hours
 More than 2½ hours

High-frequency alternatives

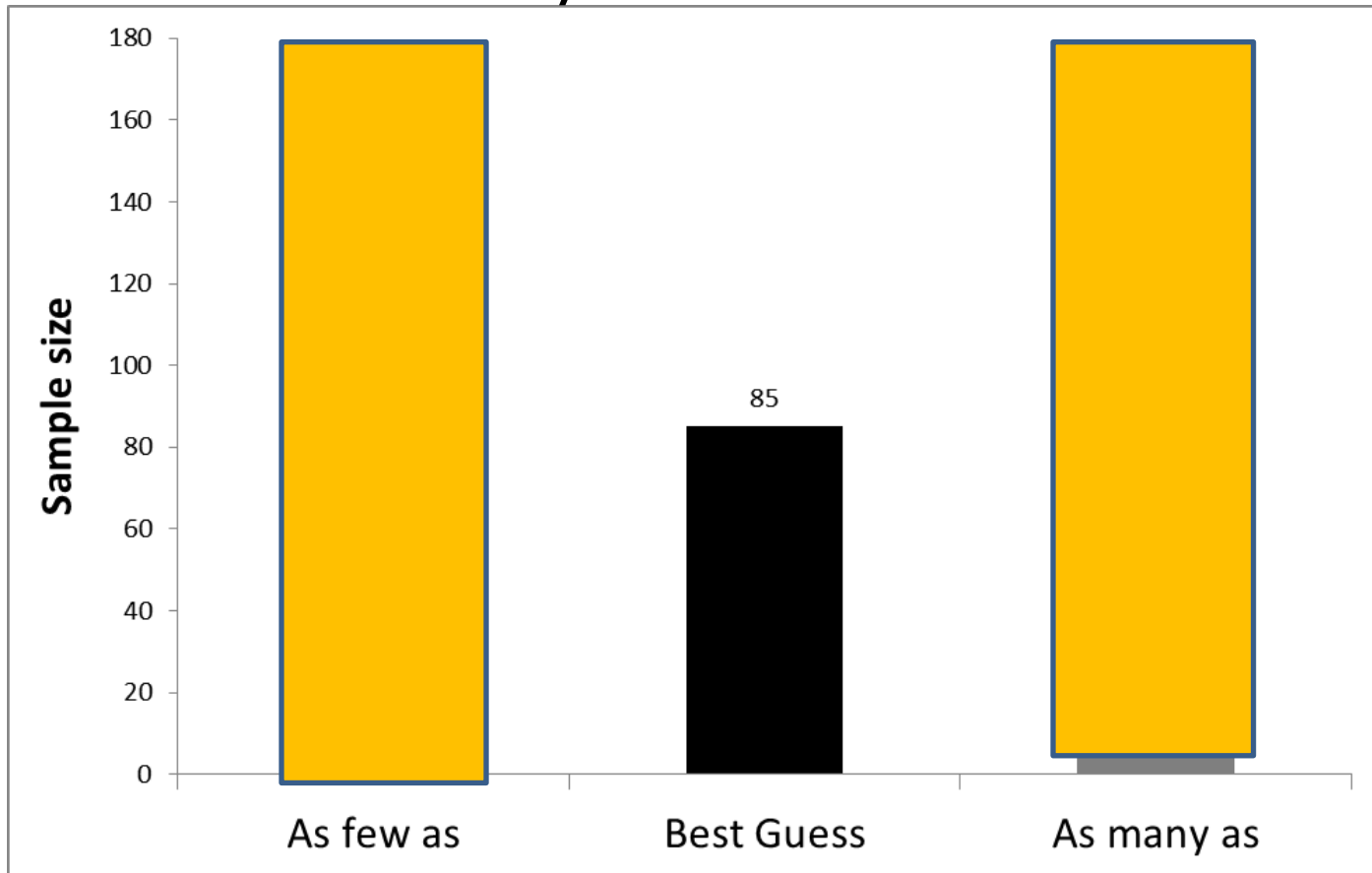
Up to 2½ hours
 2½ hours to 3 hours
 3 hours to 3½ hours
 3½ hours to 4 hours
 4 hours to 4½ hours
 More than 4½ hours

Low-vs.-high category scales (Schwarz et al., 1985)



If 5 identical studies already done

- Best guess: $n=85$
- How sure are you?



Best case scenario gives range 3:1

Reality is massively worse

- Nobody runs 6th identical study.
 - Moderator: Fluency
 - Mediator: Perceived-norms
 - DV: 'Real' behavior
- Publication bias

Where to get d from?

- Pilot
- Existing findings
- Theory/gut

Say you think/feel $d \sim .4$

$d = .44 \sim .4$

$\rightarrow n = 83$

$d = .35, \sim .4$

$\rightarrow n = 130$

Rounding error \rightarrow 100 more participants

Transition (key) slide

- Guessing d is completely impractical
→ Power analysis is also.
- Step back: Problem with underpowering?
- Unclear what failure means.
- Well, when you put it that way:
Let's power so that we know what failure means.

Existing view

1. Goal: **Success**
2. Guess d
3. Set n :
 "80%" success

New View

1. Goal: **Learn from results**
2. **Accept d is unknown**
 If interesting \rightarrow o possible
 If o possible \rightarrow very small possible
3. Set n :
 100% learning
 Works: keep going
 Fails: Go Home

What is “Going *Big*”?

A. Limited resources (most cases)

(e.g., lab studies)

- What n are you **willing to pay** for this effect?
- Run n
 - Fails, too small for me.
 - Works, keep going, adjust n .

B. ‘Unlimited’ resources (fewest cases)

(e.g., Project Implicit, Facebook)

- Smallest effect you *care* about

SLIDES ABOUT P-VALUES

Defining Evidential Value

- **Statistical significance**

Single finding:
unlikely result of chance

Could be caused by selective reporting rather than chance

- **Evidential value**

Set of significant findings:
unlikely result of selective reporting

Motivation: we only publish if $p < .05$

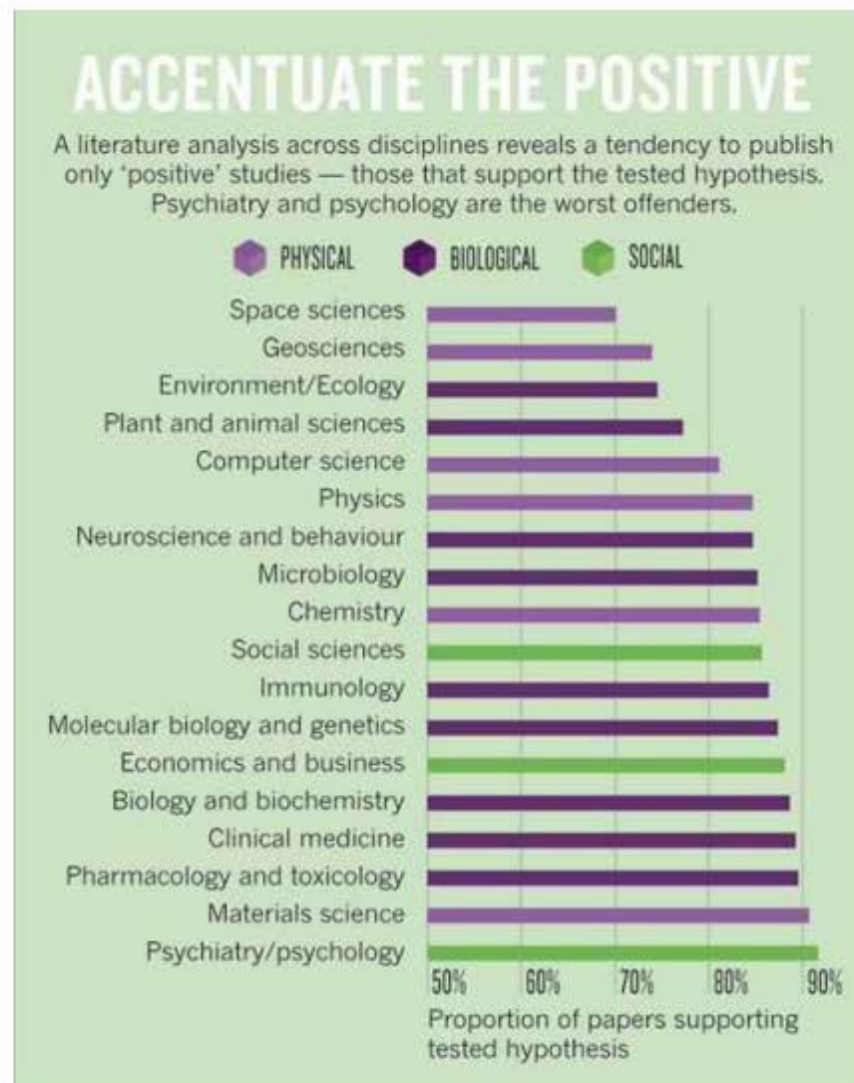


Figure 1: From Fanelli, D. *Scientometrics* 90, 891–904 (2011).

Motivation

Nonexisting effects: only see false-positive evidence

Existing effects: only see strongest evidence

**Published scientific evidence is not
representative of reality.**

Outline

- Shape
- Inference
- Demonstration
- How often is p-curve wrong?
- Effect size estimation
- Selecting p -values

p-curve's shape

- **Effect does not exist:** flat
- **Effect exists:** right-skew.
(more lows than highs)
- ***Intensely p-hacked:*** left-skew
(more highs than lows)

Why flat if null is true?

***p*-value:**

prob(result | null is true).

Under the null:

- What percent of findings $p \leq .30$
 - 30%
- What percent of findings $p \leq .05$
 - 5%
- What percent of findings $p \leq .04$
 - 4%
- What percent of findings $p \leq .03$
 - 3%

Got it.

Why more lows than high if true?

(right skew)

- Height: men vs. women
- N = Philadelphia
- What result is more likely?

*In Philadelphia, men taller than women ($p=.047$)
($p=.007$)*

- Not into intuition?

*Differential convexity of the density function
Wallis (Econometrica, 1942)*

Why left skew with p -hacking?

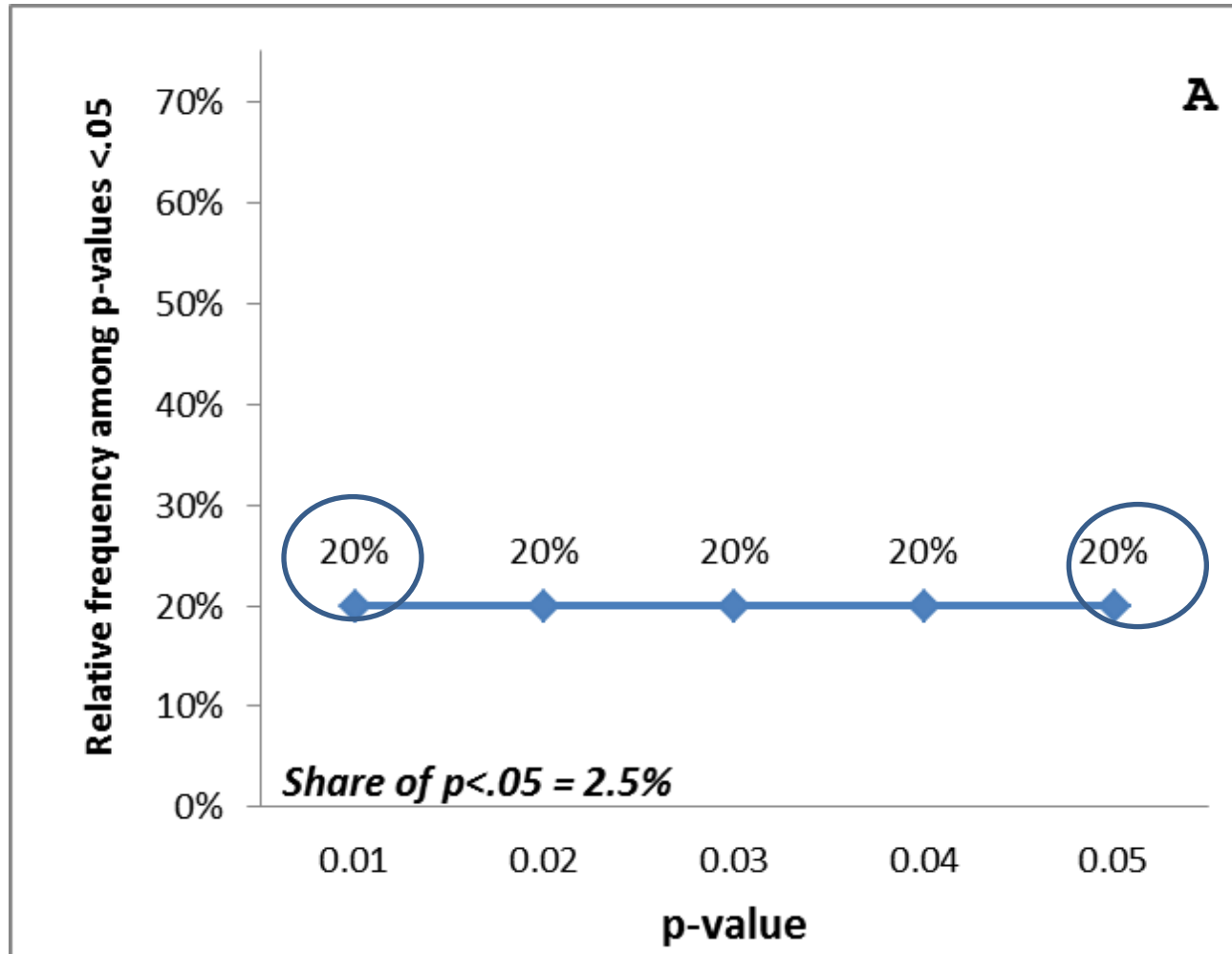
- Because p -hackers have limited ambition
- $p=.21$
 - Drop if >2.5 SD
- $p=.13$
 - Control for gender
- $p=.04$
 - Write Intro
- If we stop p -hacking as soon as $p<.05$,
- Won't get to $p=.02$ very often.

Plotting Expected P -curves

- Two-sample t -tests.
- True effect sizes
 - $d=0, d=.3, d=.6, d=.9$
- p -hacking
 - No: $n=20$
 - Yes: $n=\{20, 25, 30, 35, 40\}$

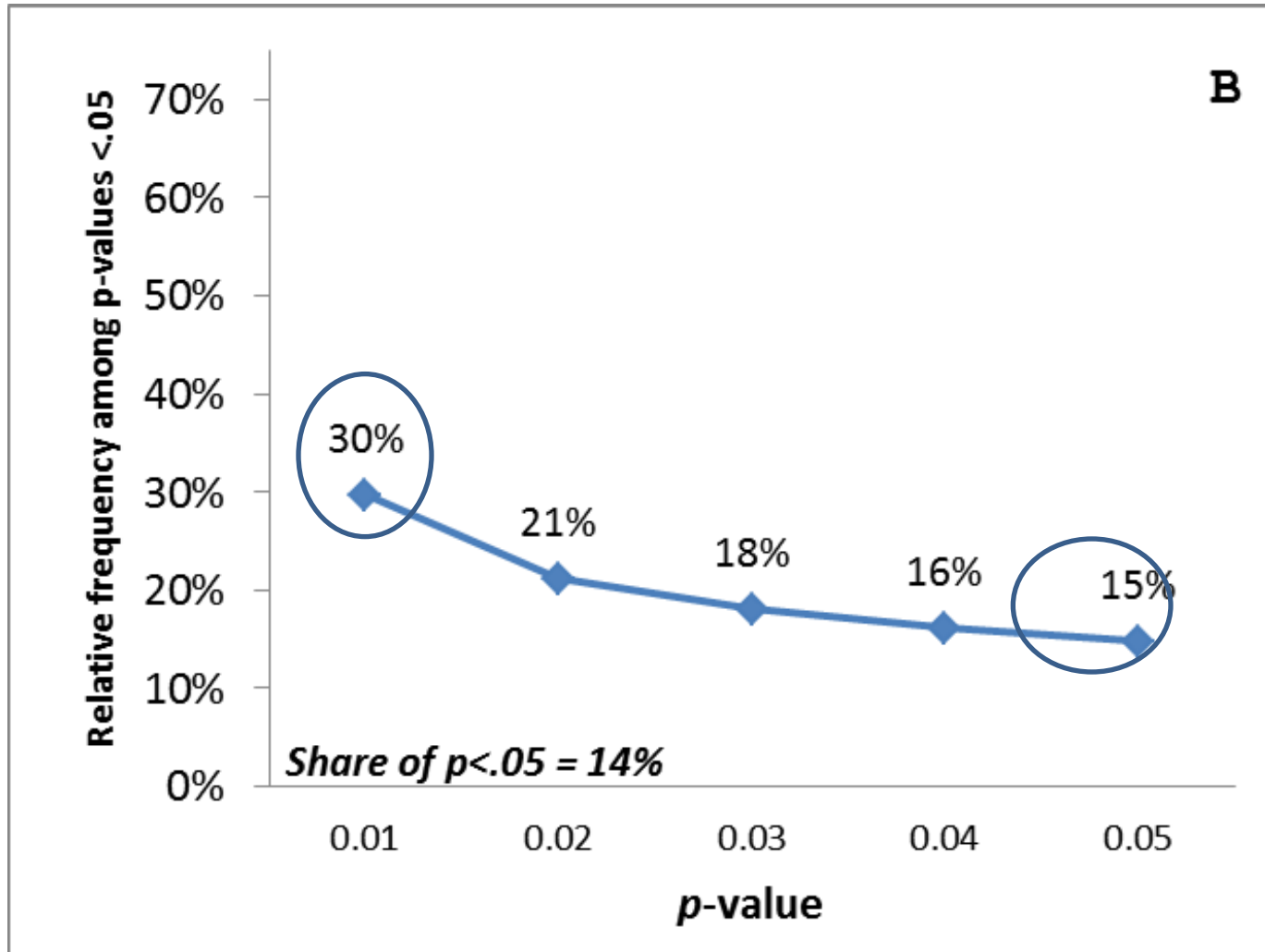
Nonexisting effect ($n=20$, $d=0$)

As many $p < .01$ as $p > .04$



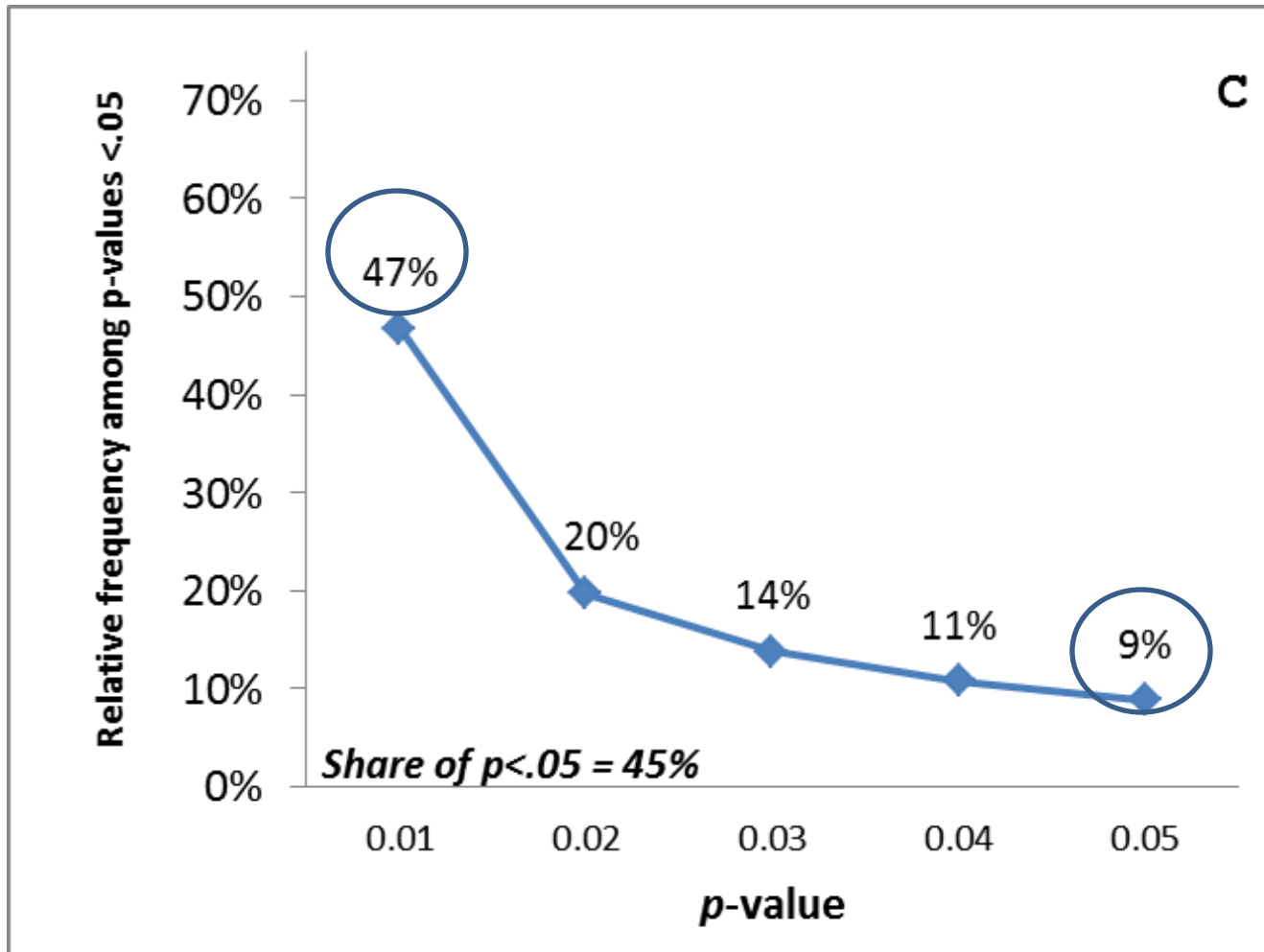
$n=20, d=.3$ / power=14%

Two $p < .01$ for every $p > .04$



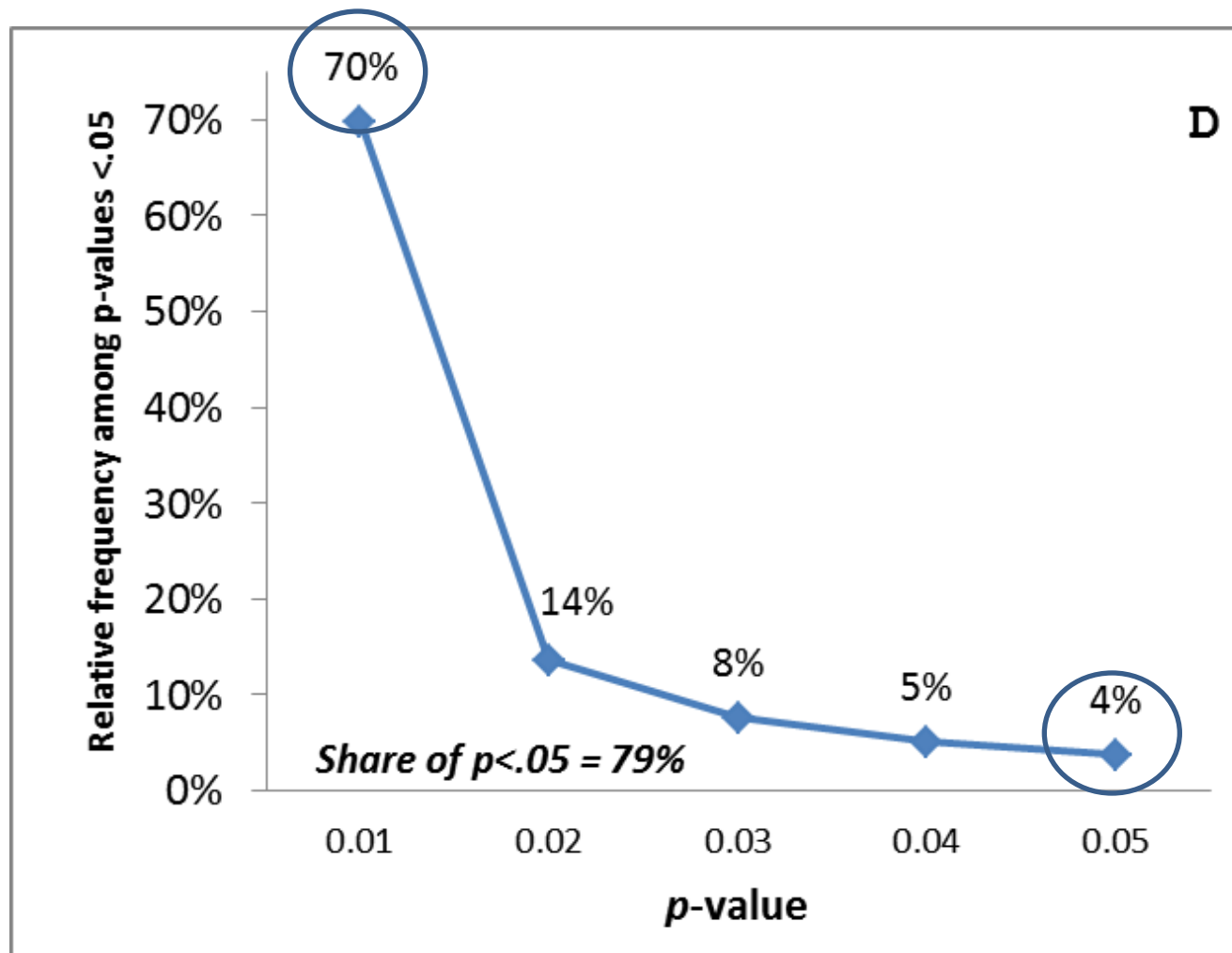
$n=20, d=.6$ / power = 45%

Five $p < .01$ per every one $p > .04$



$n=20, d=.9$ / power=79%

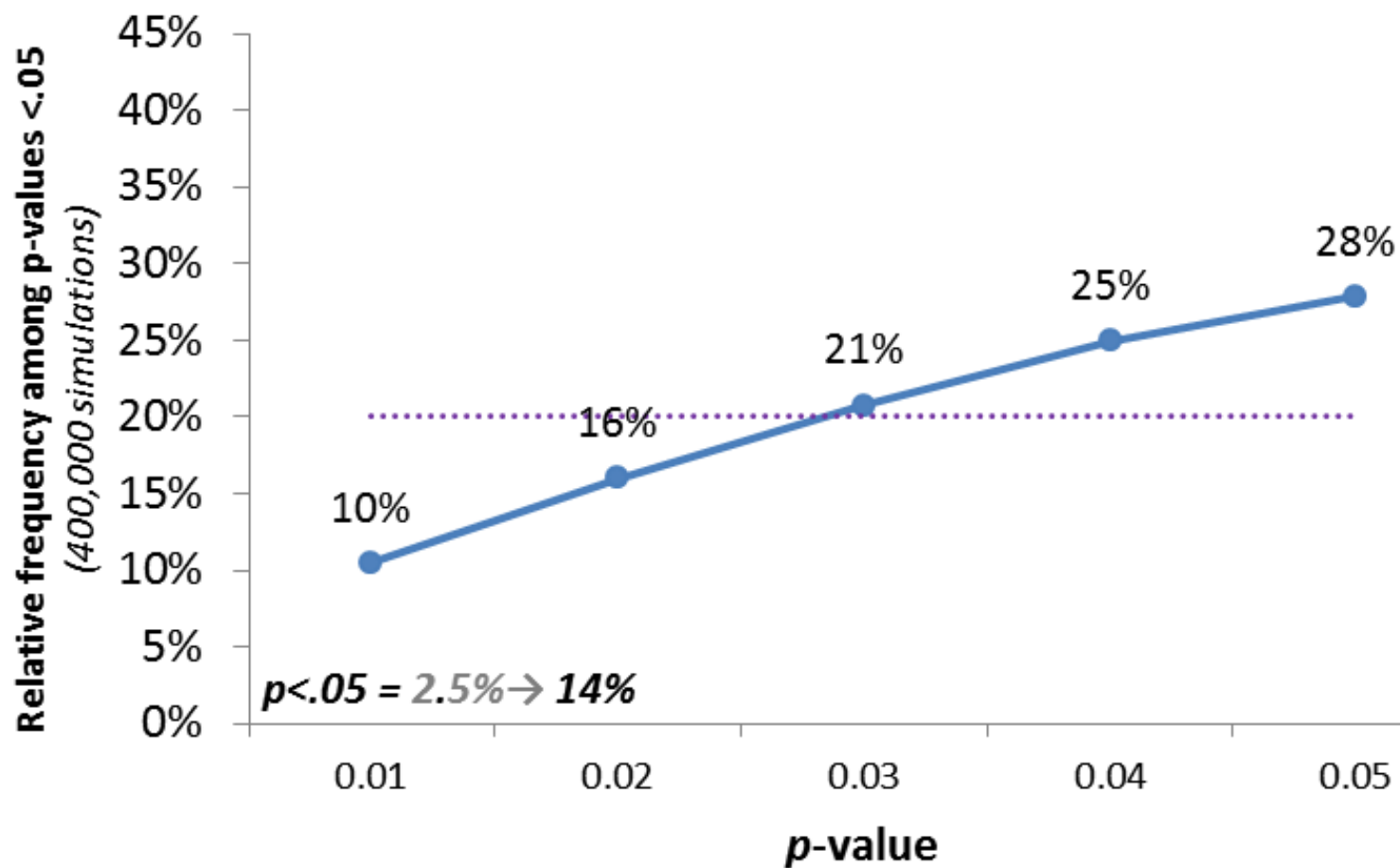
Eigtheen $p < .01$ per every $p > .04$.



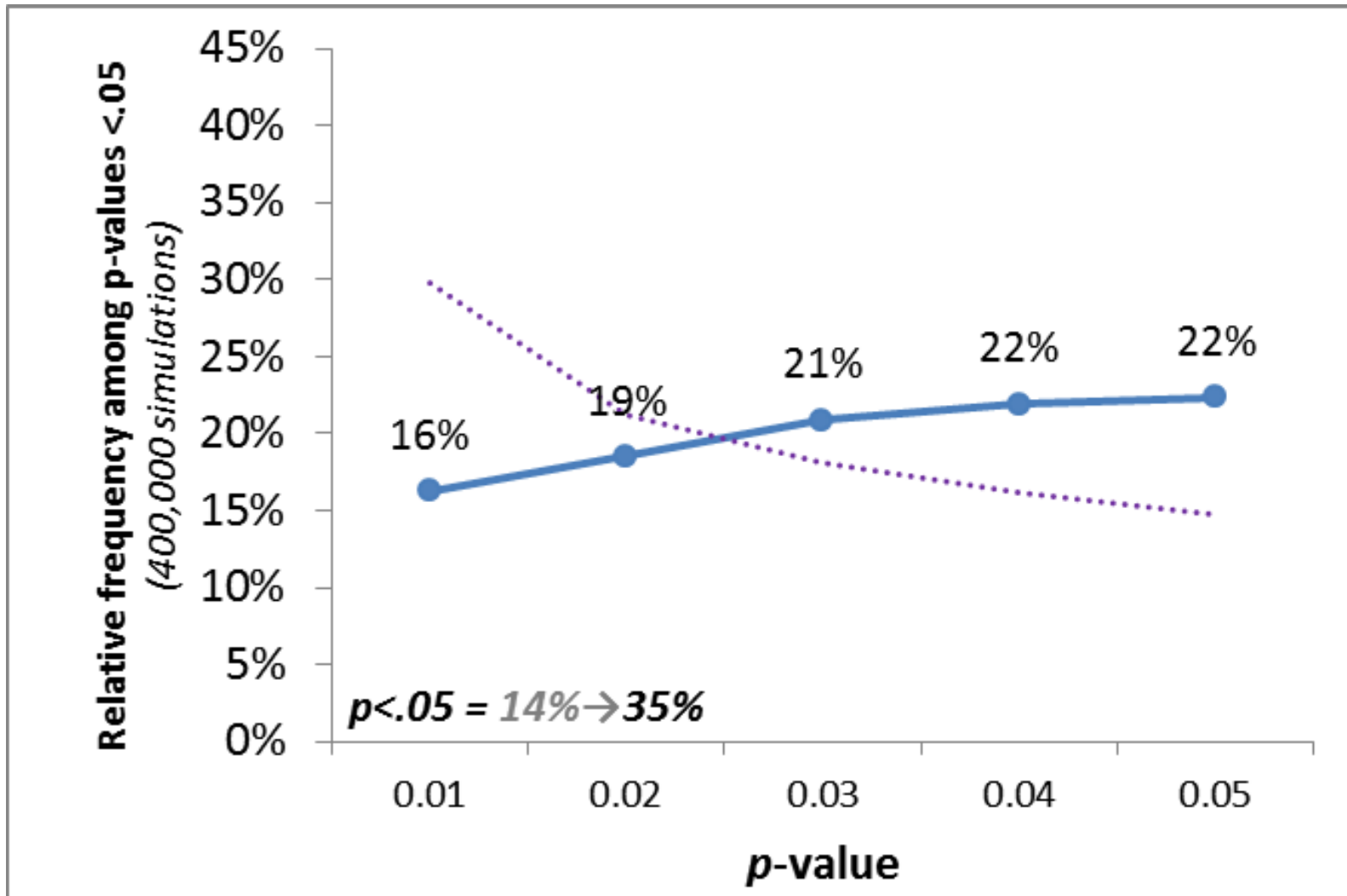
Adding *p*-hacking

$$n=\{20,25,30,35,40\}$$

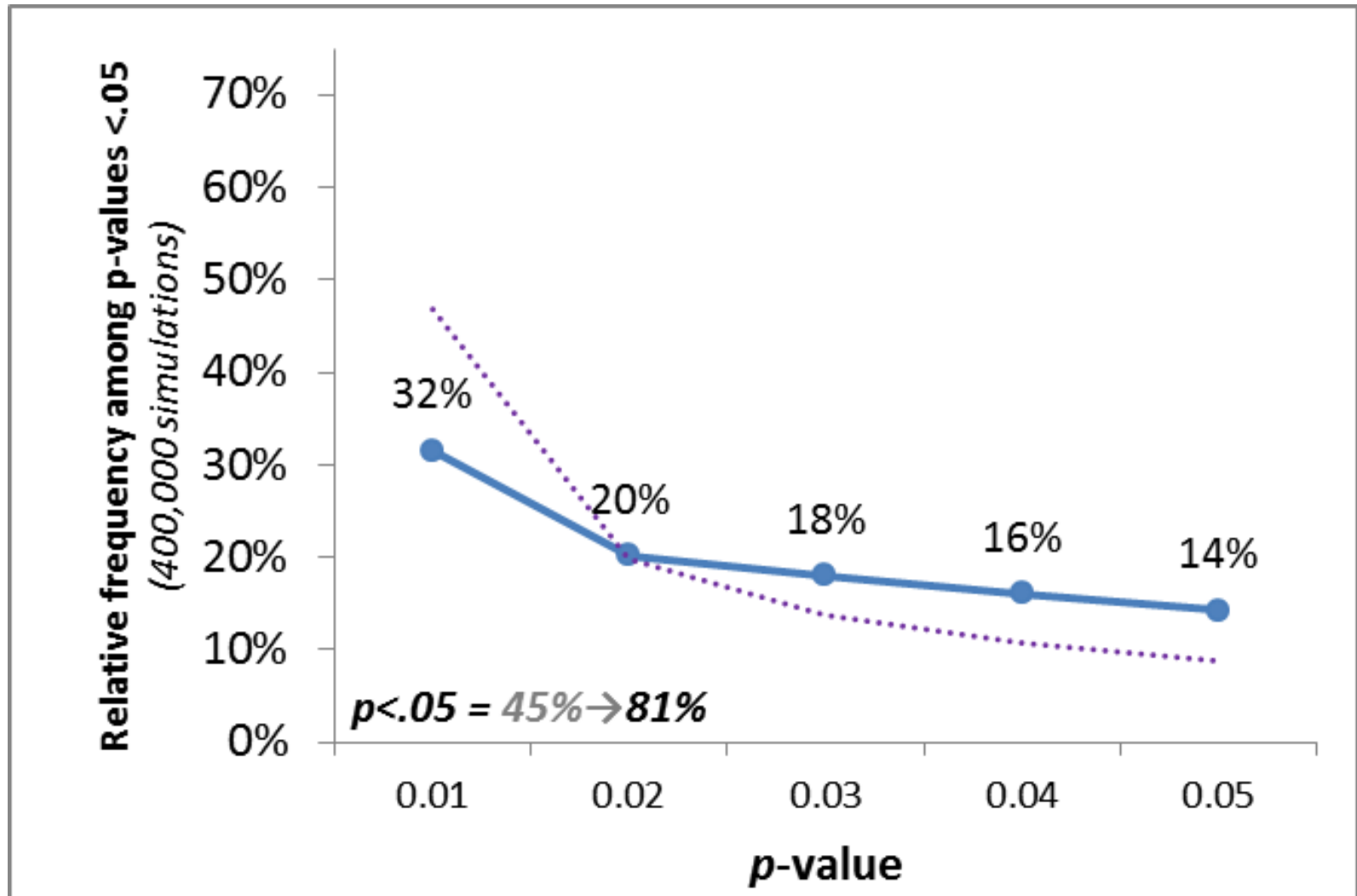
$$d=0$$



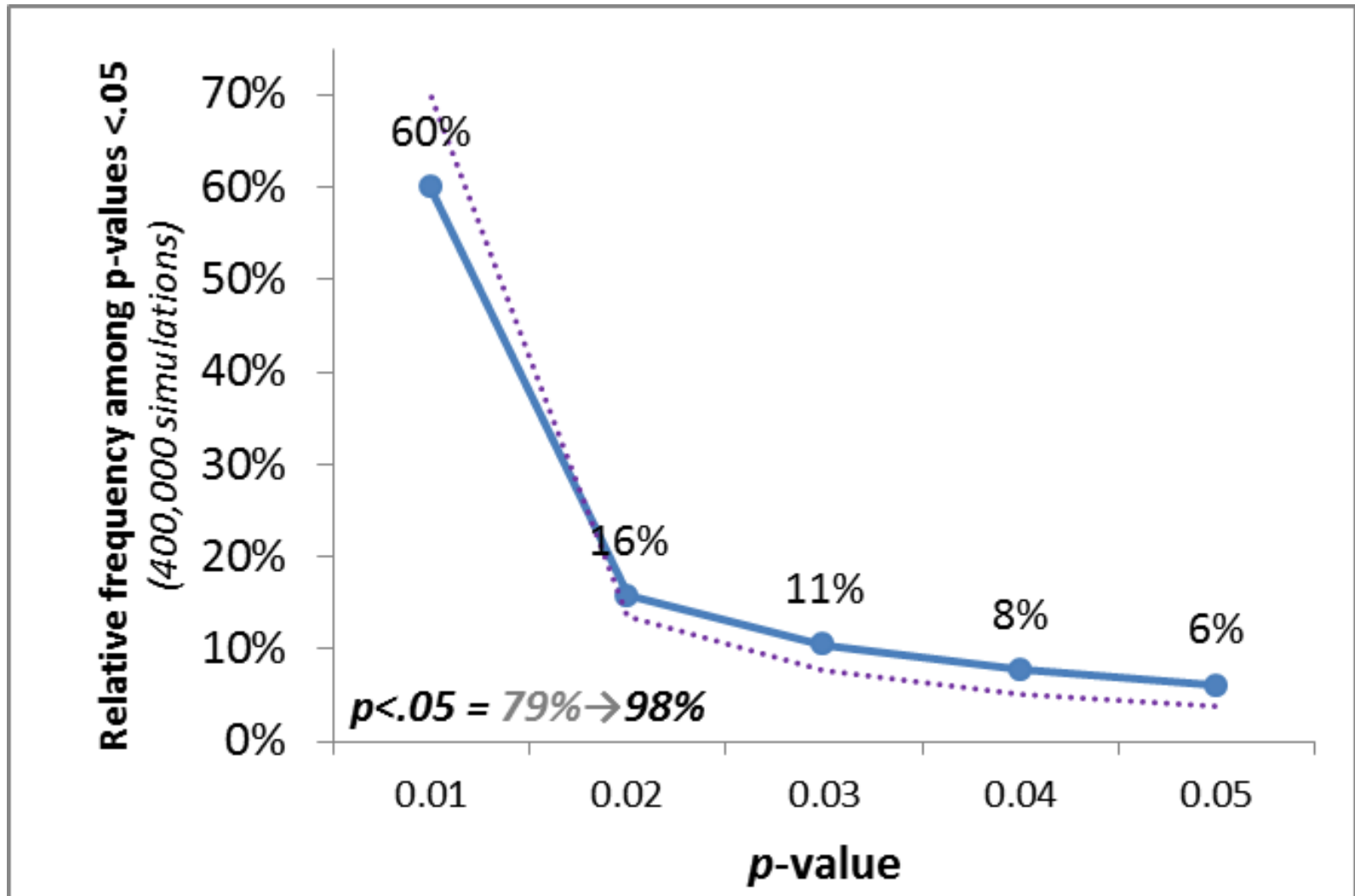
$d=.3$ / original power=14%



$d=.6$ / original-power = 45%



$d=.9$ / original-power=79%



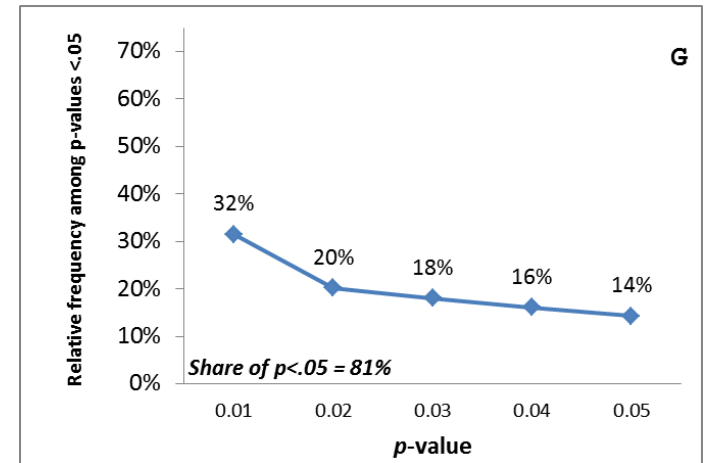
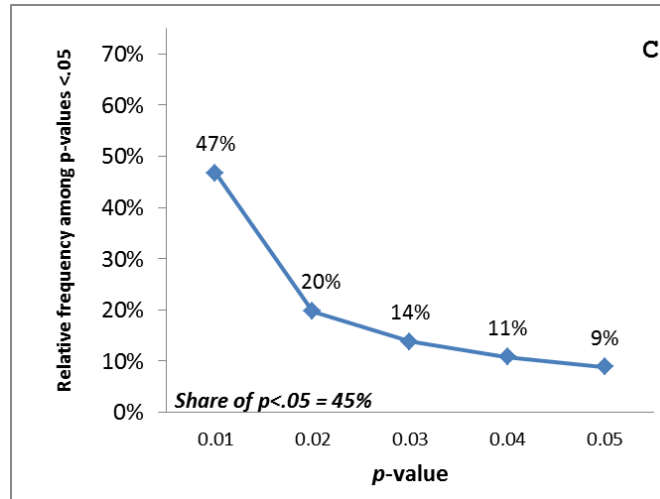
p-hacked findings?

NO

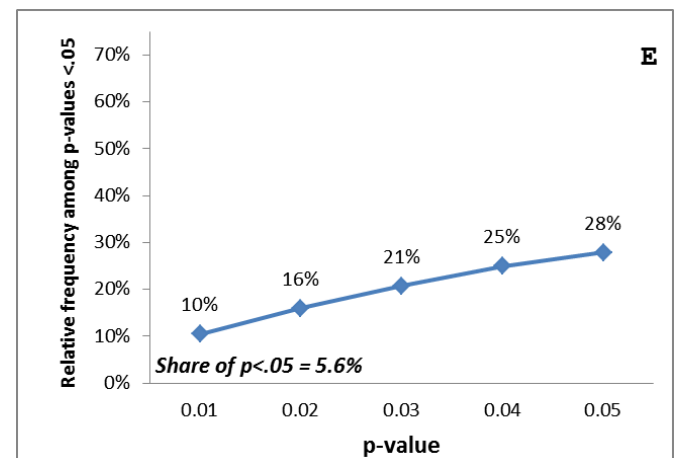
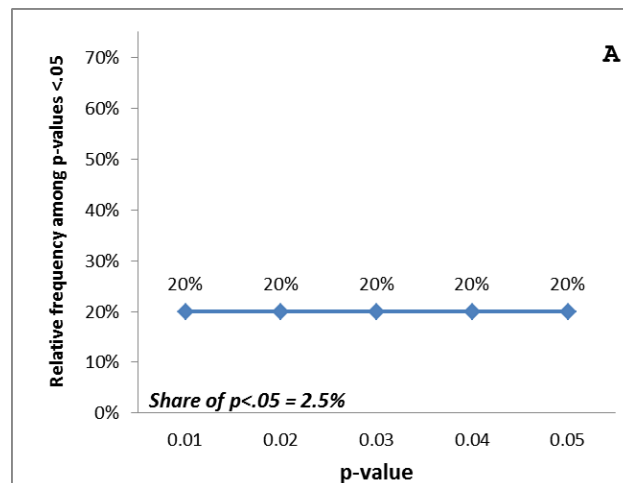
YES

YES

Effect
Exists?



NO



Note:

- p -curve does not test if p -hacking happens.
(it “always” does)

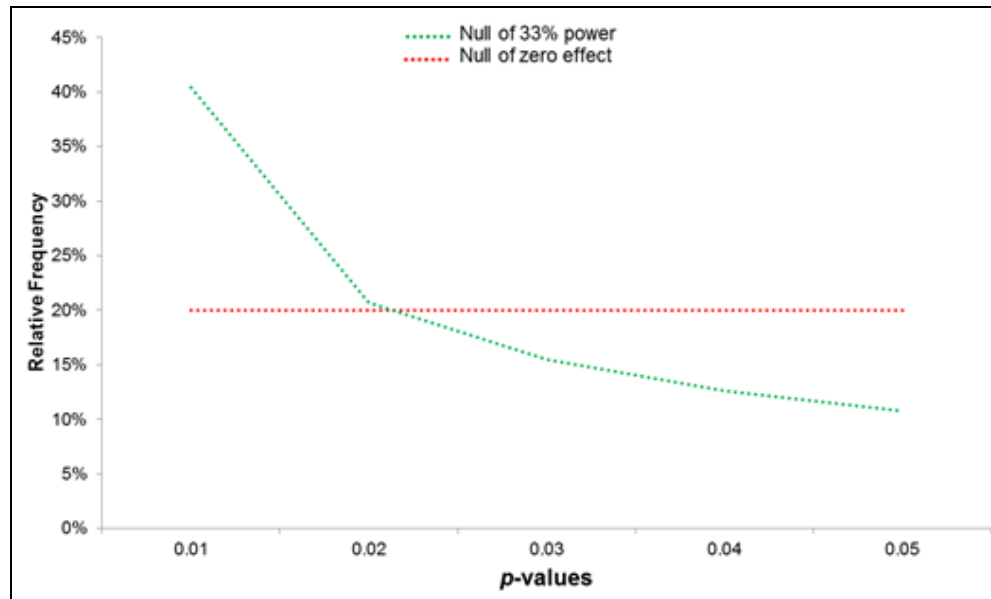
Rather:

- Whether p -hacking was so intense that it eliminated evidential value (if any).

Outline

- Shape
- Inference
- Demonstration
- How often is p-curve wrong?
- Effect-size estimation
- Selecting p -values

Inference with p-curve

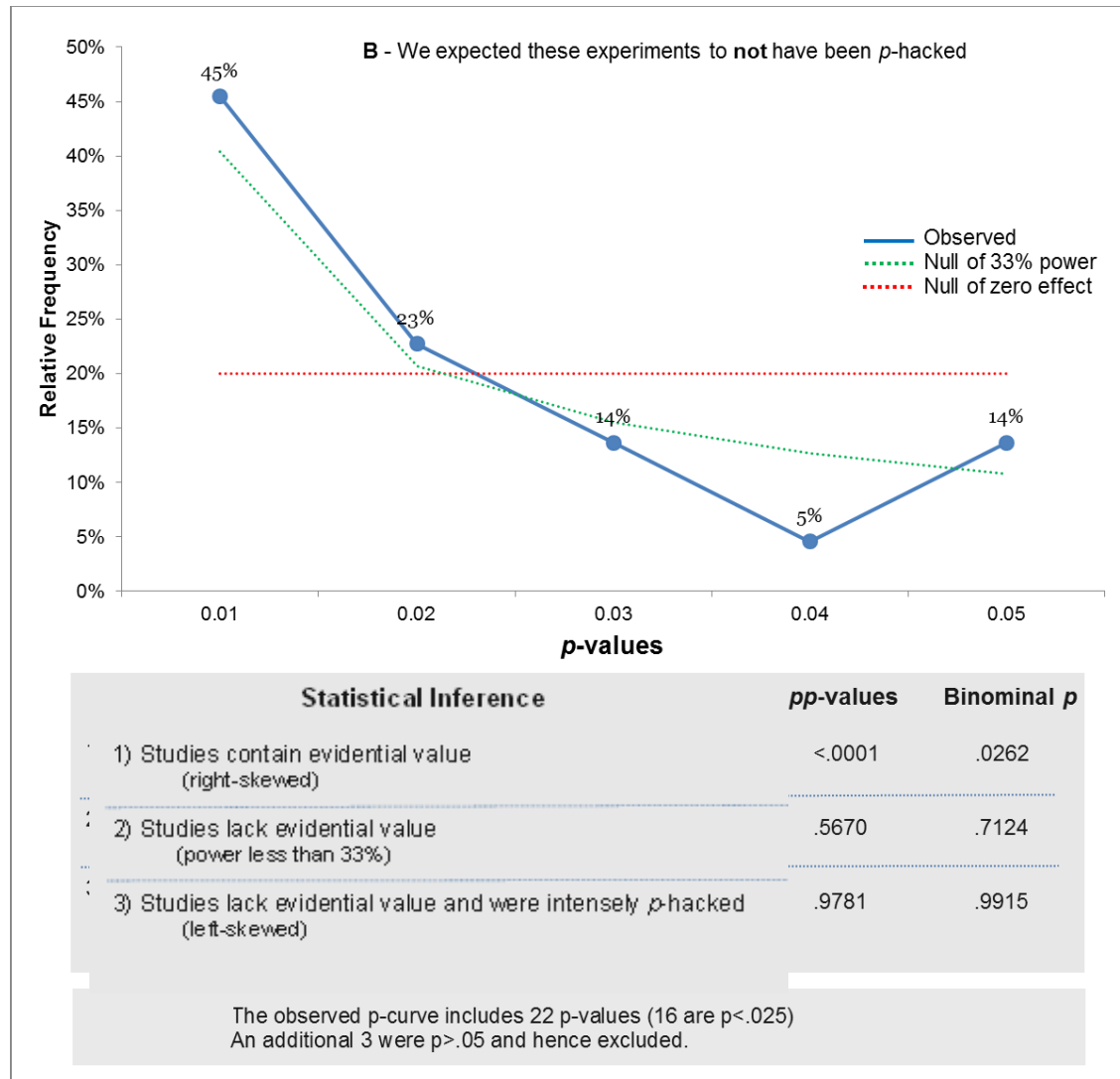


- 1) Right-skewed?
- 2) Flatter than studies powered at 33%?
- 3) Left-skewed?

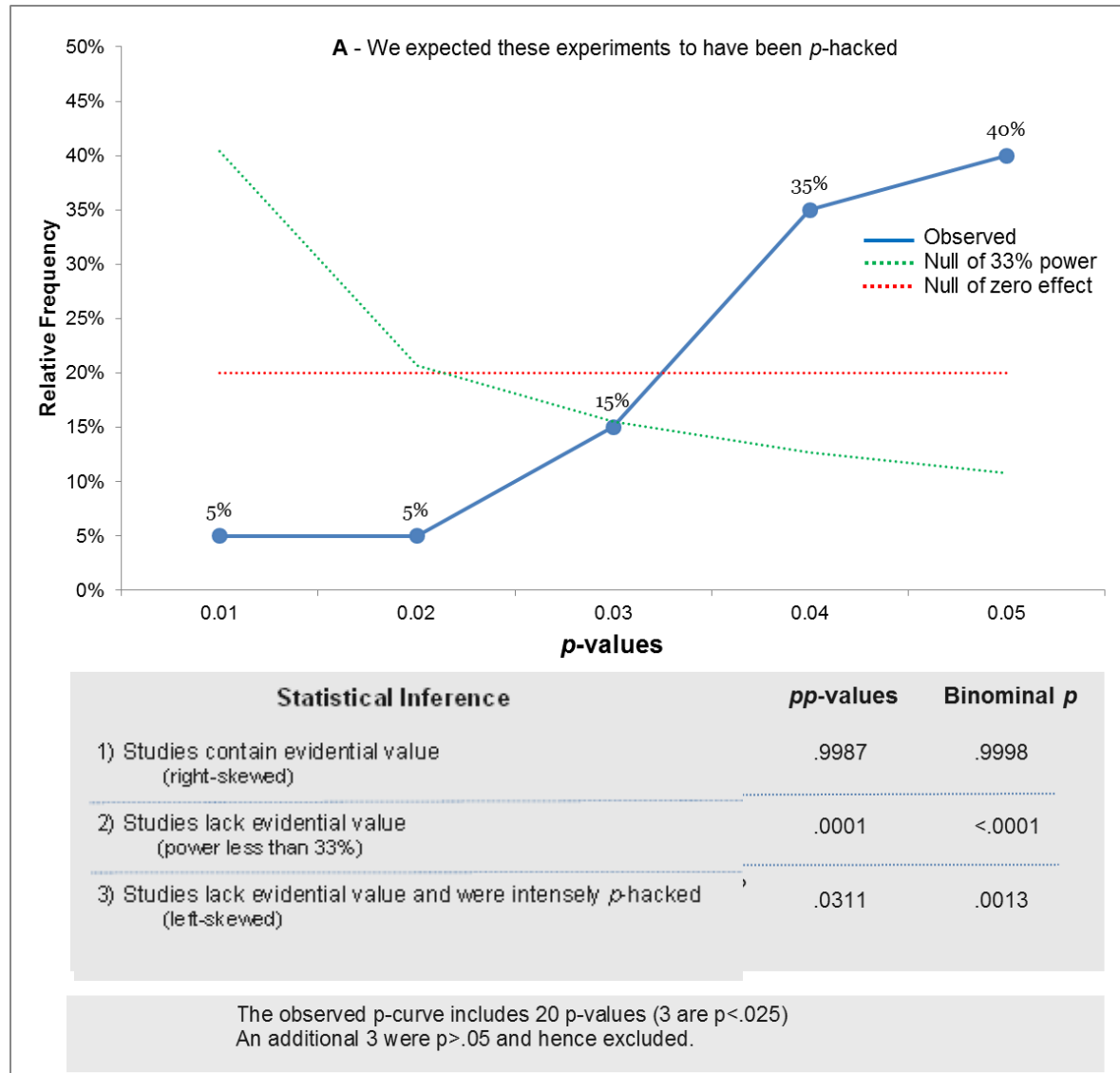
Outline

- Shape
- Inference
- Demonstration
- How often is p-curve wrong?
- Effect-size estimation
- Selecting p -values

Set 1: JPSP with no exclusions nor transformations



Set 2: JPSP result reported only with covariate



- **Next:** New Example

One Swallow Doesn't Make a Summer: New Evidence on Anchoring Effects[†]

By ZACHARIAS MANIADIS, FABIO TUFANO, AND JOHN A. LIST^{*}

Some researchers have argued that anchoring in economic valuations casts doubt on the assumption of consistent and stable preferences. We present new evidence that explores the strength of certain anchoring results. We then present a theoretical framework that provides insights into why we should be cautious of initial empirical findings in general. The model importantly highlights that the rate of false positives depends not only on the observed significance level, but also on statistical power, research priors, and the number of scholars exploring the question. Importantly, a few independent replications dramatically increase the chances that the original finding is true. (JEL D12, C91)



First draft: 2013 10 24

This draft: 2014 04 11



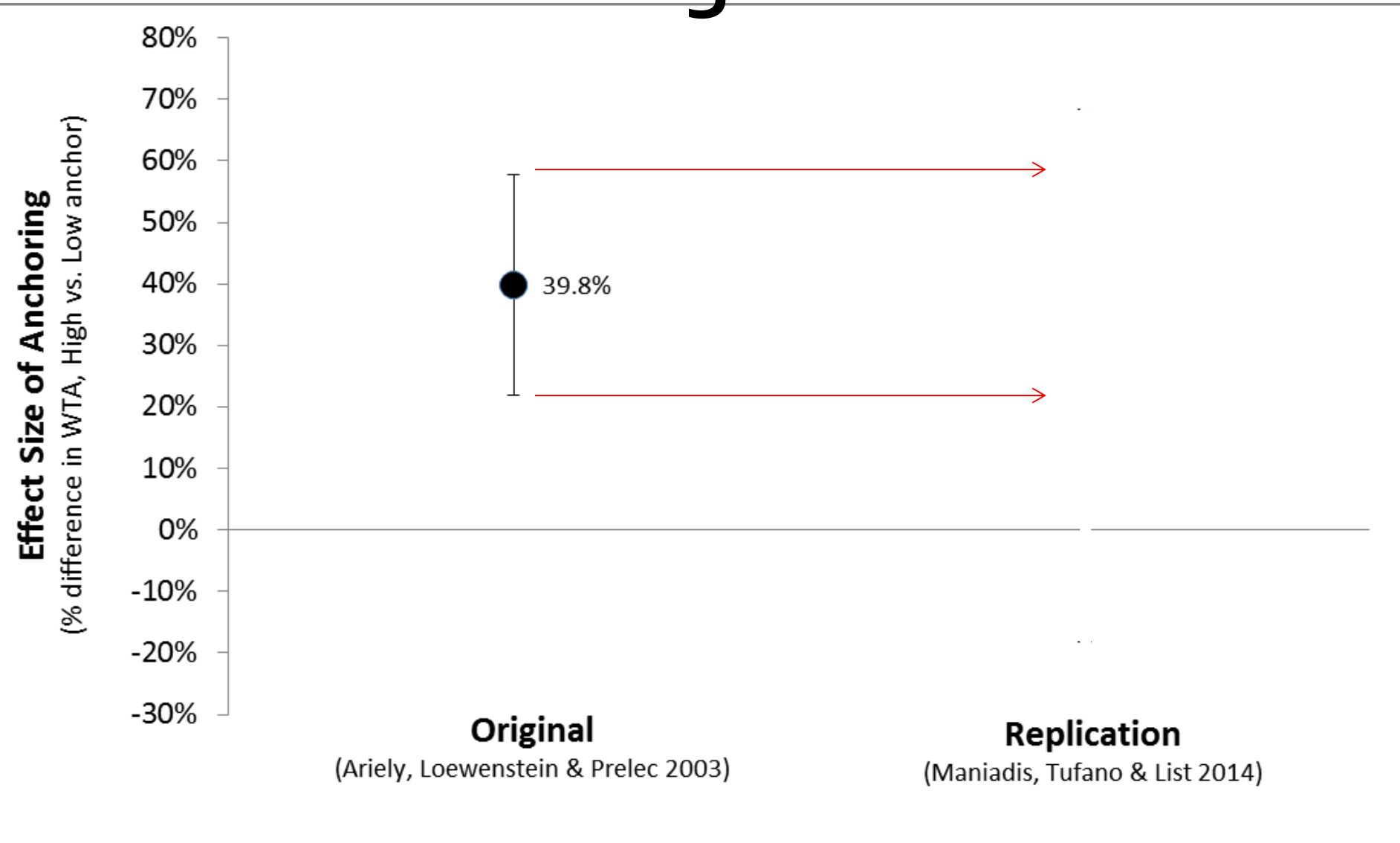
Anchoring is Not a False-Positive: Maniadis, Tufano, and List's (2014) “Failure-to-Replicate” is Actually Entirely Consistent with the Original

Uri Simonsohn
University of Pennsylvania
uws@wharton.upenn.edu

Joseph P. Simmons
University of Pennsylvania
jpsimmo@wharton.upenn.edu

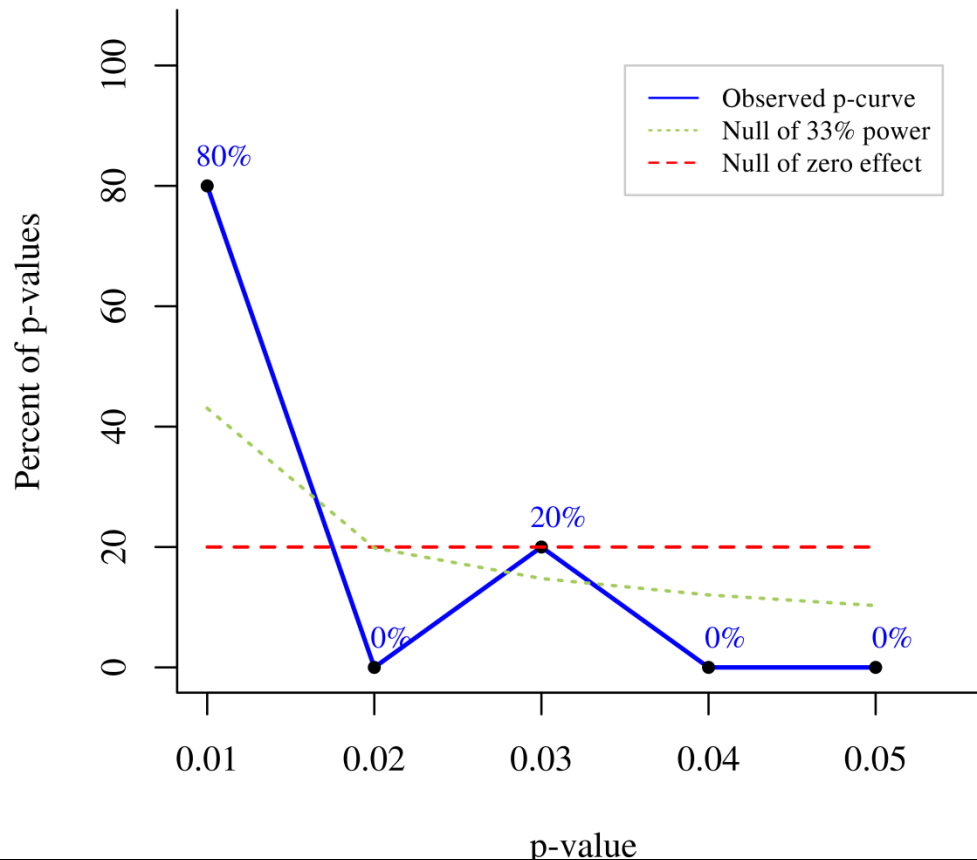
Leif D. Nelson
UC Berkeley
leif_nelson@haas.berkeley.edu

Anchoring and WTA



- Bad replication $\neg \rightarrow$ Good original
- Was original a false-positive?

p-curve results



Statistical Inference

1) Studies contain evidential value
(*right-skewed*)

2) Studies lack evidential value
(*flatter than 33% power*)

3) Studies lack evidential value and were intensely *p*-hacked
(*left-skewed*)

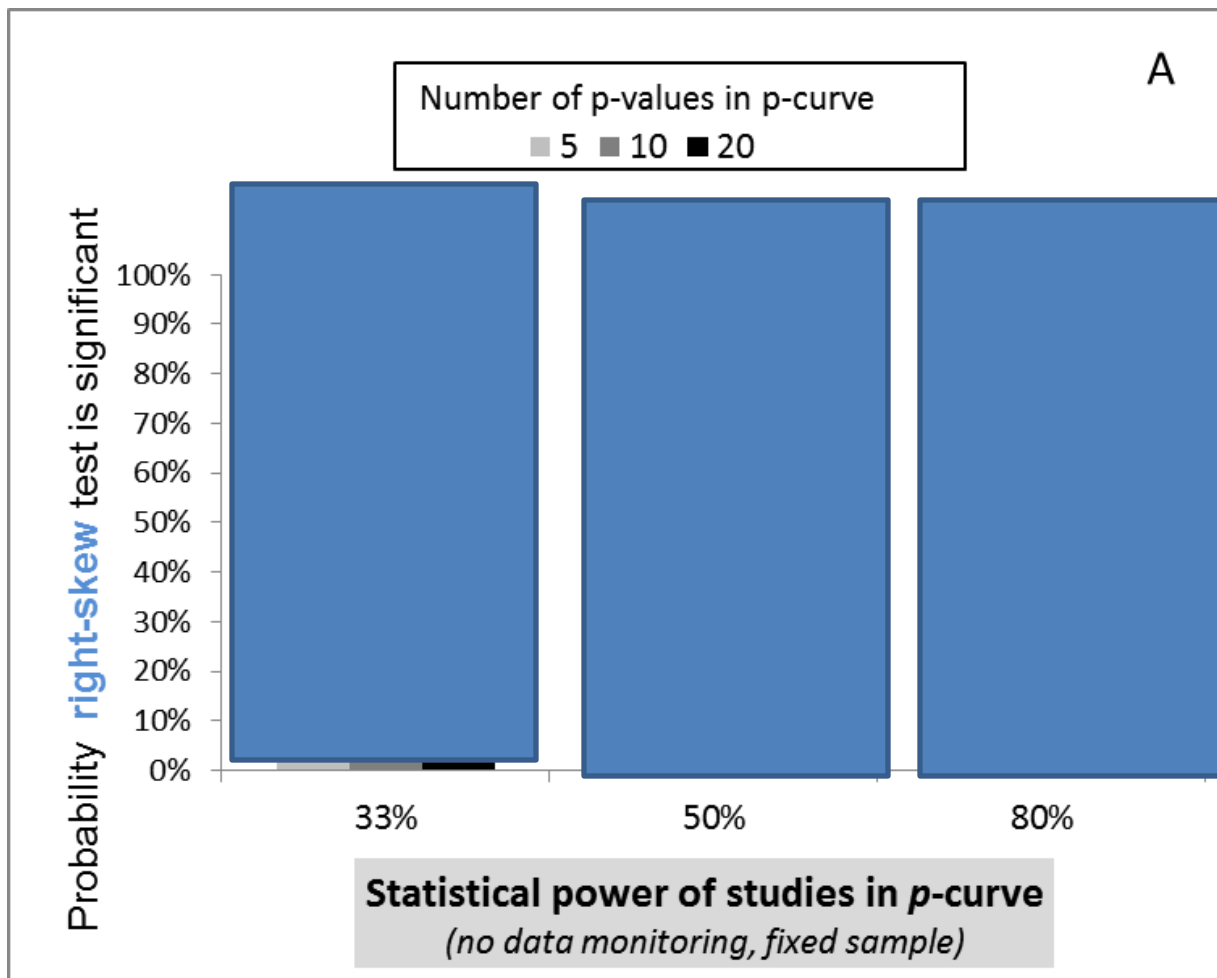
Results

$\chi^2(10)=33.76, p=.0002$

$\chi^2(10)=2.8, p=.9857$

$\chi^2(10)=1.35, p=.9993$

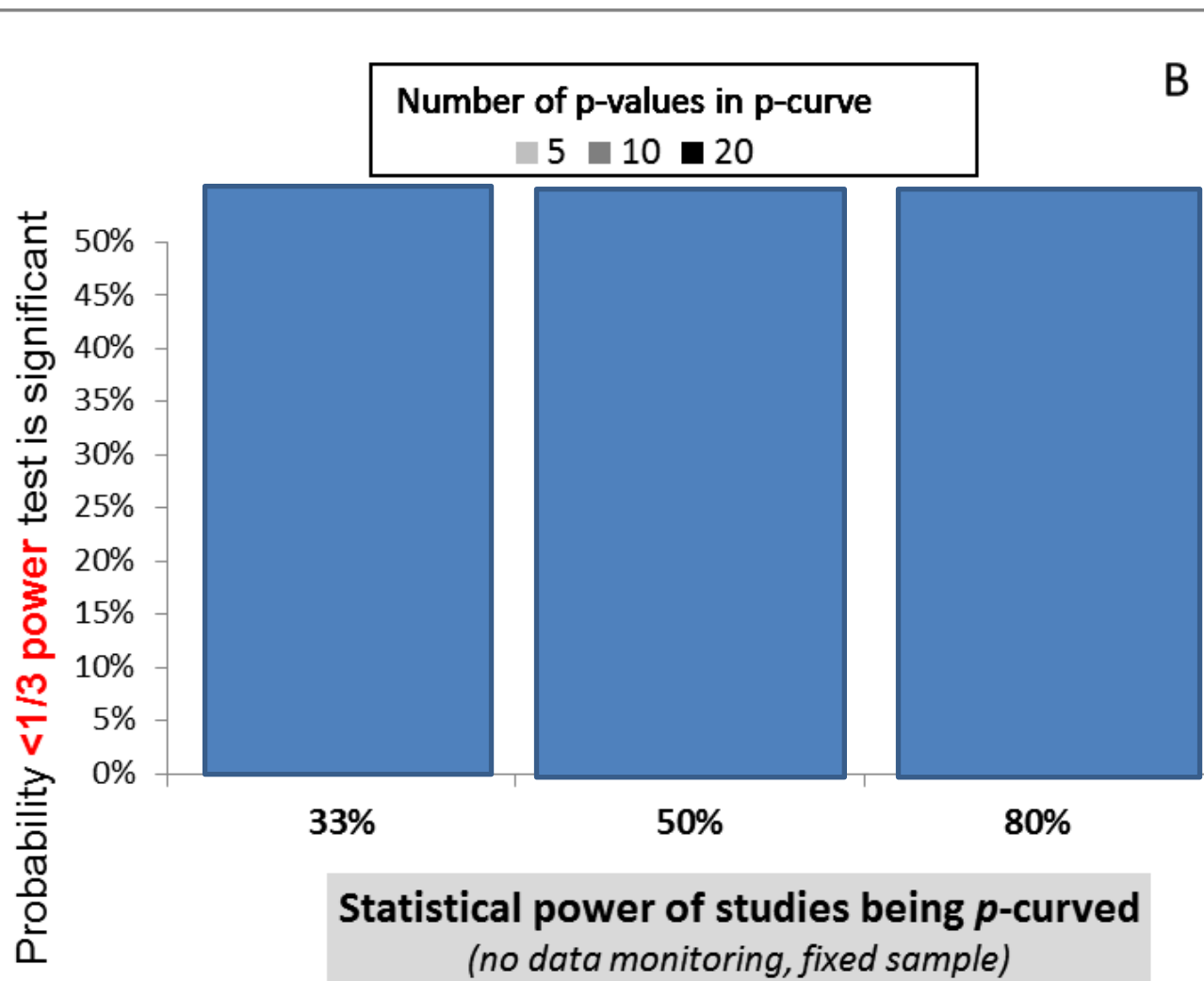
When effect **exists**, how often does p -curve say “evidential value”



Highlights:

More power at 5
Certain with 80%

When effect **exists**, how often does p -curve say “**no evidential value**”



Highlights

- P-curve is ‘never’ wrong on properly powered studies.

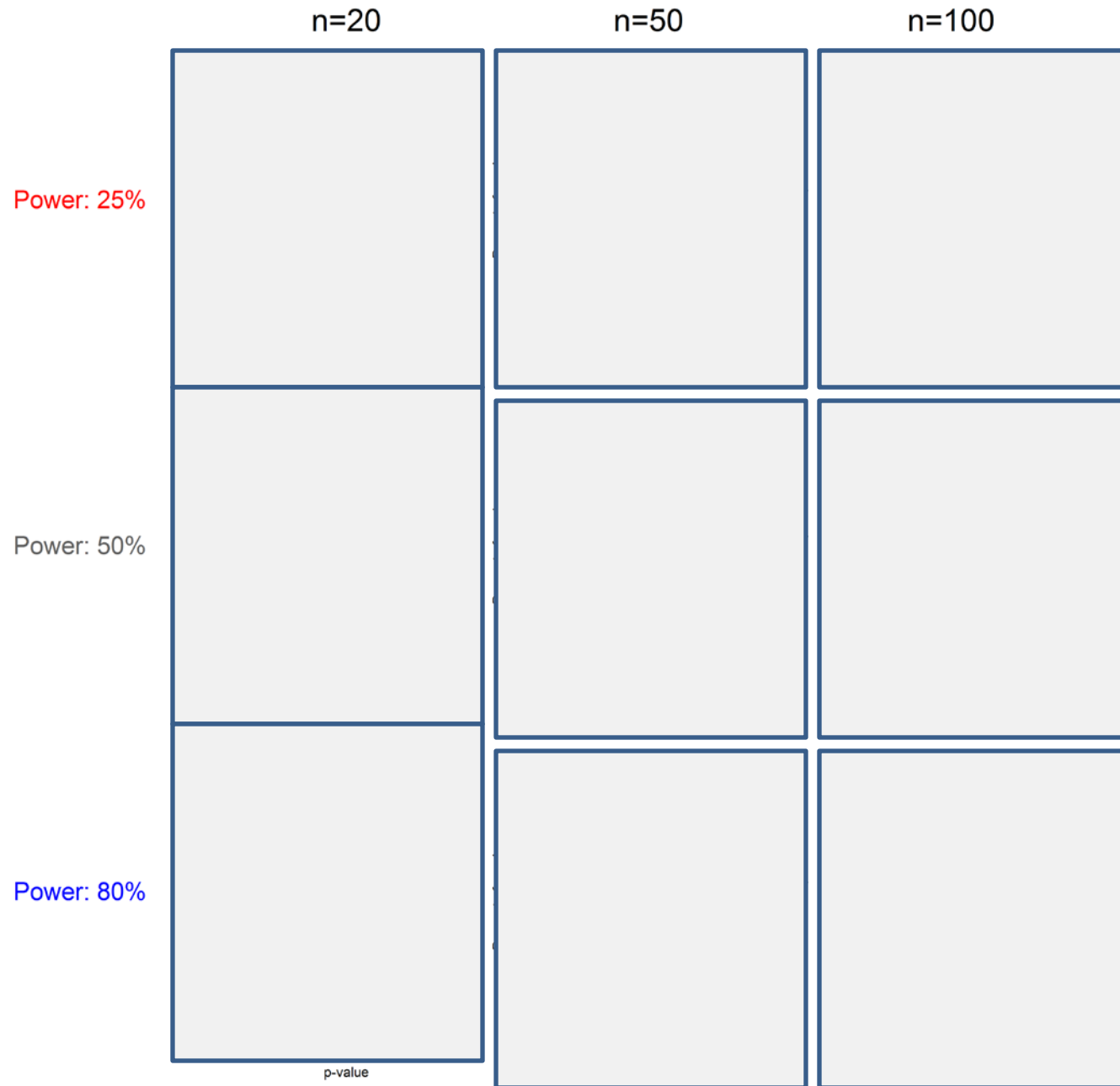
Broad big picture applications

- Possible uses:
 - Meta-analyses of X on Y
 - Meta-analyses of X on anything
 - Meta-analyses of anything on Y
 - Relative truth of opposing findings
 - X is good for Y, vs
 - X is bad for Y
 - Is this journal, on average, true?
 - Universities vs. pharmaceuticals

Everyday applications

(note: 5 p-values can be plenty)

- **Reader:** Should I read this paper?
- **Researcher:** Run expensive follow-up?
- **Researcher:** Explain inconsistent previous finding
- **Reviewer:** Ask for direct replications?



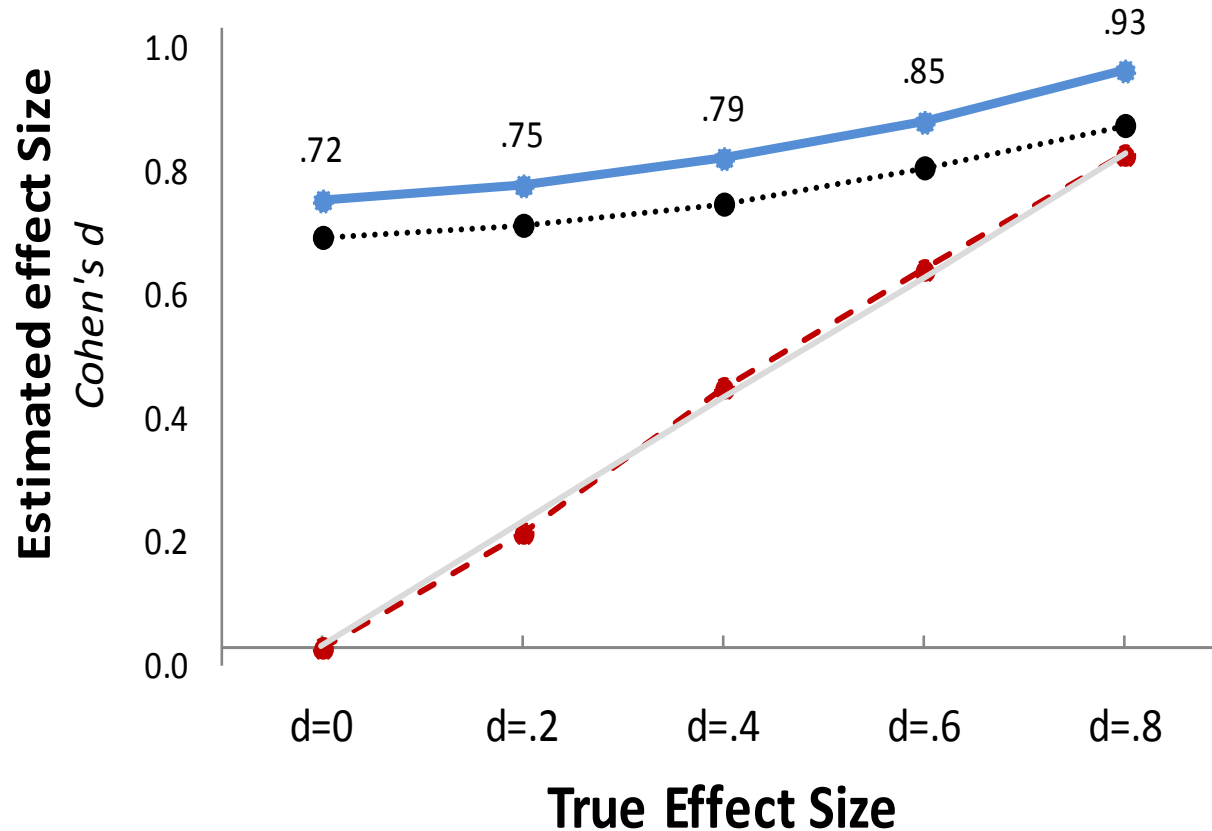
- **Next.**
 - Simulated meta-analysis, file-drawering studies.

—●— Average significant effect ●..... w/ trim-and-fill correction - - - ● - - - p-curve's estimate — True effect size

B

Predetermined sample size: between N=10 & N=70

Fixed effect size: $d_i = d$



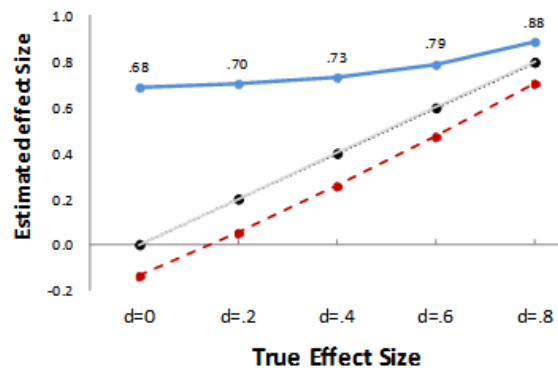
- **Next.**
 - Simulated meta-analysis, p -hacking

—●— Average significant effect —●— Average of all effects - -●- p-curve's estimate — True effect size

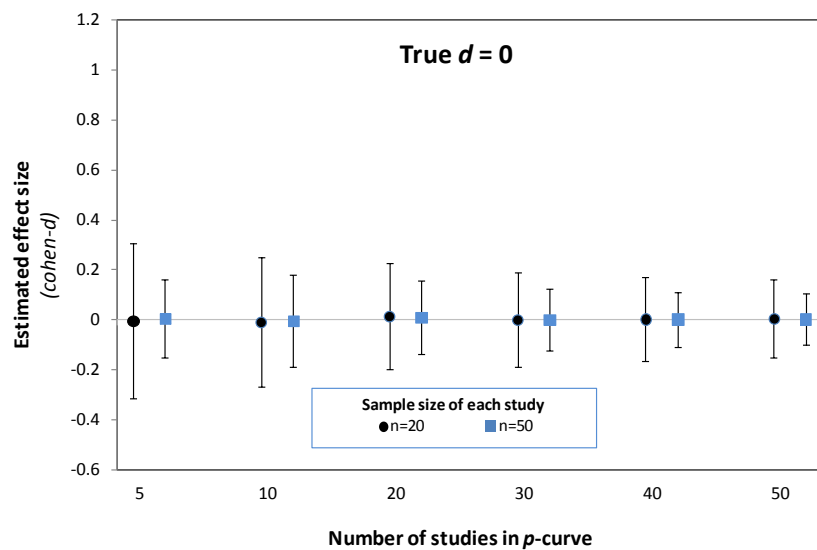
A

Data-peeking

(if $p > .05$ with $n=20$, add 10 observations)



- **Next.** Precision from few studies



Number of studies in p -curve

- **Next.** Demonstration 1: Many Labs Replication project
 - Real study, participants, data
 - But, see all attempts

Open Science Framework BETA

Explore ▾

Help ▾

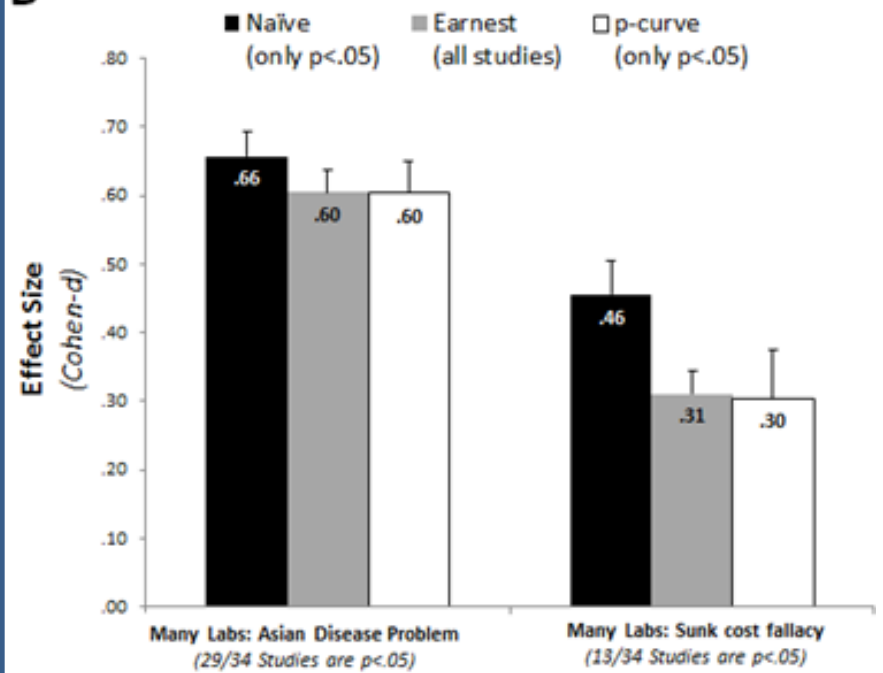
Search

Investigating Variation in Replicability: A “Many Labs” Replication Project

Contributors: [Richard A. Klein](#) | [Kate Ratliff](#) | [Michelangelo Vianello](#) | [Reginald B. Adams, Jr.](#) | [Stepan Bahnik](#) | [Michael Mark Brandt](#) | [Beach Brooks](#) | [Claudia Brumbaugh](#) | [Zeynep Cemalcilar](#) | [Jesse J. Chandler](#) | [Winnee Cheong](#) | [William Matthew Eisner](#) | [Natalia Frankowska](#) | [David Furrow](#) | [Elisa Maria Galliani](#) | [Fred Hasselman](#) | [Joshua A. Hicks](#) | [James Jeffrey R. Huntsinger](#) | [Hans IJzerman](#) | [Melissa-Sue John](#) | [Jennifer Joy-Gaba](#) | [Heather Kappes](#) | [Lacy Elise Krueger](#) | [Robyn Mallett](#) | [Wendy Morris](#) | [Anthony J. Nelson](#) | [Jason A. Nier](#) | [Grant Packard](#) | [Ronaldo Pilati](#) | [Abraham M. Rutch](#) | [Skorinko](#) | [Robert W. Smith](#) | [Troy G. Steiner](#) | [Justin Storbeck](#) | [Lyn van swol](#) | [Donna Thompson](#) | [Anna van 't Veer](#) | [Aaron Wichman](#) | [Julie A. Woodzicka](#) | [Brian A. Nosek](#)

- 36 labs
- 13 “effects”
 - Example 1. Sunk Cost (Significant: 50% labs)
 - Example 2. Asian Disease (86%)

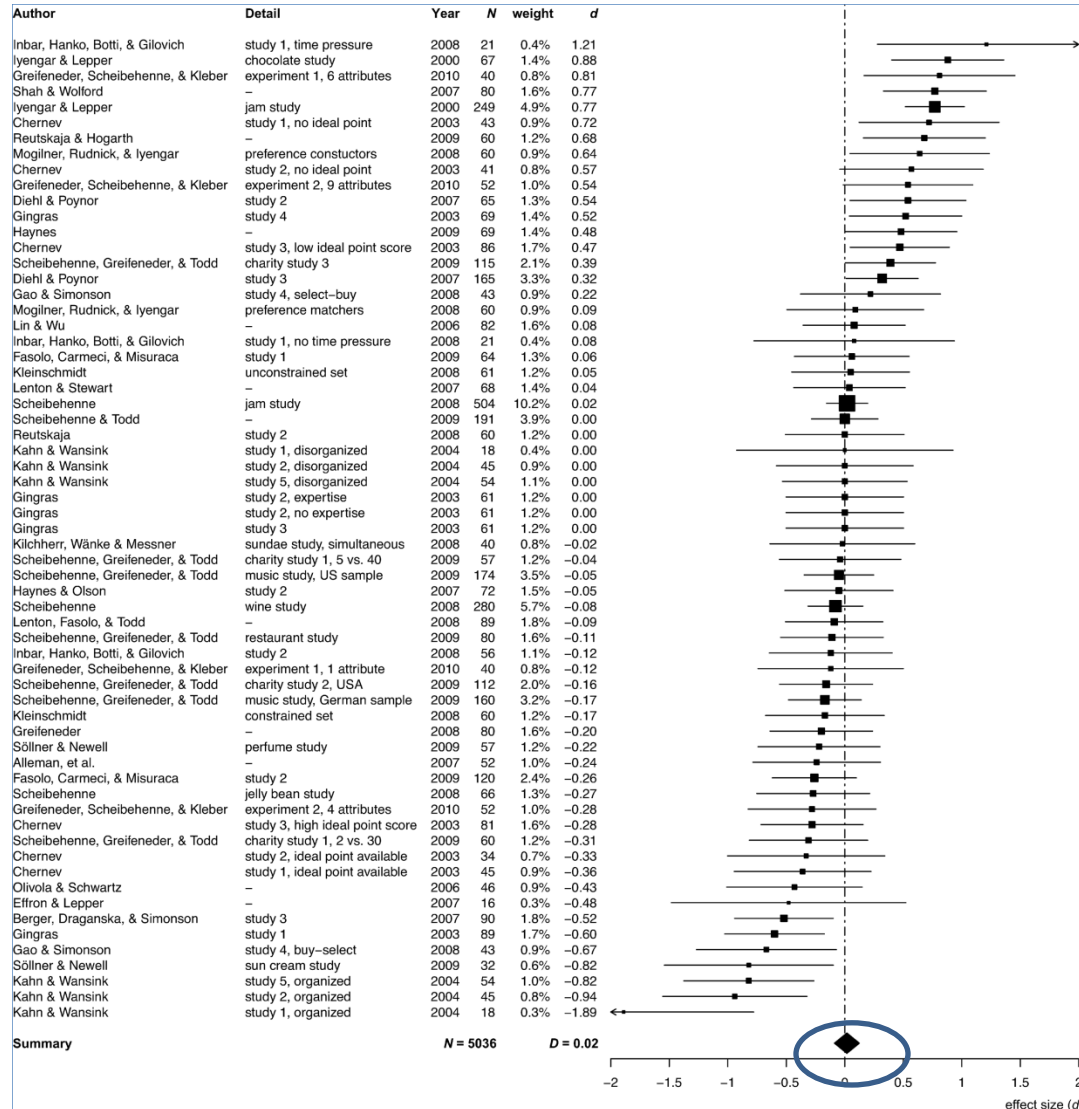
B



- **Next.** Demonstration 2: Choice Overload

A demonstration

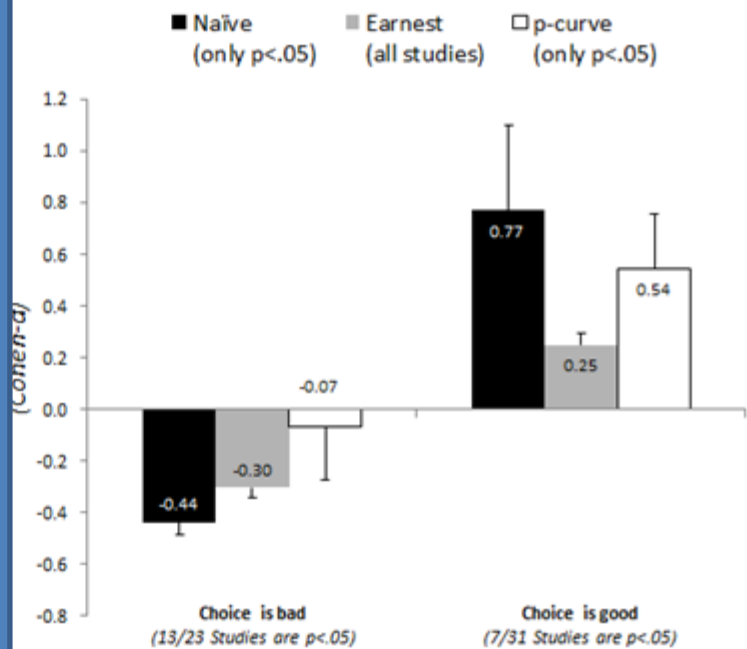
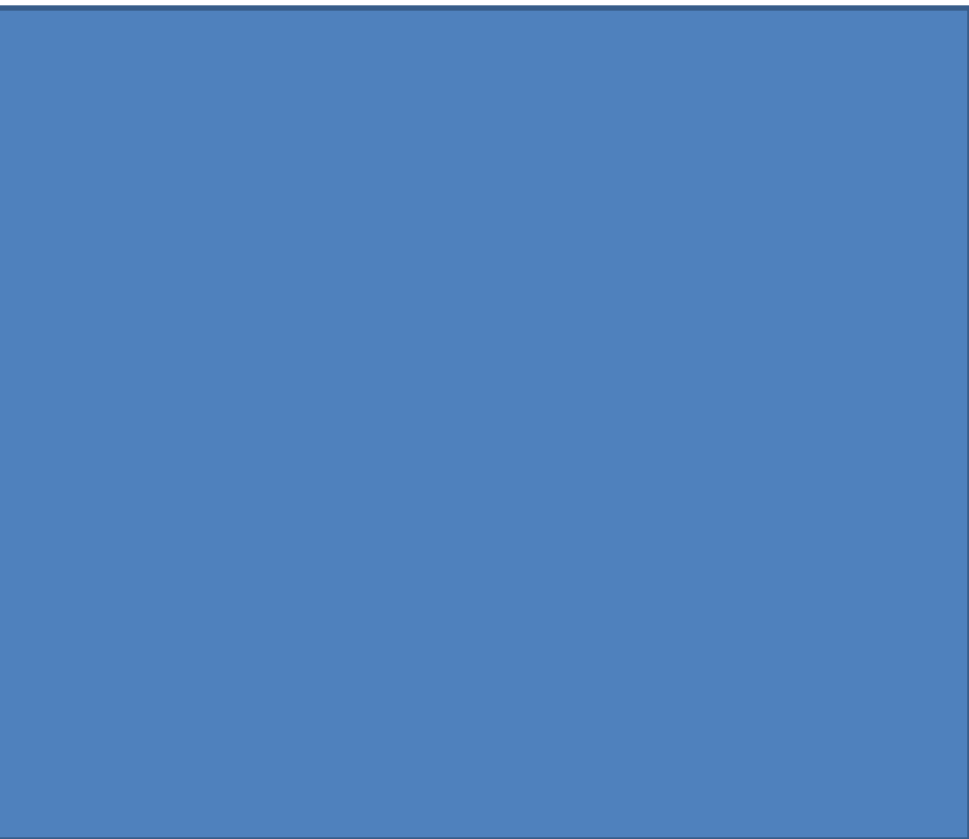
Choice Overload meta-analysis



**
Choice is bad

Choice is good

**



How to think about p-values

- When a study has lots of statistical power (big effect + big sample), expect to see very small p-values.
- When you see a really big p-value ($p = .048$), you should be concerned.
- Unexpected thought: When the p-values are really small in the absence of statistical power, you can have different (more unsettling) concerns.

I don't have any more slides, but I have many more thoughts and opinions.

Ask.



datacolada.org

P-curve

p-curve.com

Paper 1 - Evidential Value	Paper 2 - Effect size	The online app	The User Guide	Supp Materials
<p>This draft: 2017 04 29</p> <p>Figure 1</p> <p>P-curve: A Key To The File Drawer</p> <p>Leo D. Nelson, Joseph P. Simmons, David A. Lakens, University of California, University of Pennsylvania, University of Pennsylvania, Berkeley</p> <p>Word count: 6220</p> <p>* Corresponding author: Leo D. Nelson, leo@berkeley.edu, phone: 510 895 0200, 100 University Hall, 210 University Avenue, Berkeley, CA 94720</p>	<p>P-Curve Fixes Publication Bias: Obtaining Unbiased Effect Size Estimates from Published Studies Alone</p> <p>Leo D. Nelson, David A. Lakens, Joseph P. Simmons, University of California, University of Pennsylvania, University of Pennsylvania, Berkeley</p>	<p>The p-curve results</p> <p>Statistical Inferences</p> <p>1. Observed number of significant studies: 10 (right-skewed)</p> <p>2. Observed mean p-value: 0.50 (right-skewed)</p> <p>3. Observed number of significant studies: 10 (right-skewed)</p> <p>4. Observed mean p-value: 0.50 (right-skewed)</p> <p>5. Observed number of significant studies: 10 (right-skewed)</p> <p>6. Observed mean p-value: 0.50 (right-skewed)</p> <p>7. Observed number of significant studies: 10 (right-skewed)</p> <p>8. Observed mean p-value: 0.50 (right-skewed)</p> <p>9. Observed number of significant studies: 10 (right-skewed)</p> <p>10. Observed mean p-value: 0.50 (right-skewed)</p>	<p>Official User-Guide to the P-curve</p> <p>Four steps to a valid p-curve:</p> <ol style="list-style-type: none"> 1. Create and export a study selection table 2. Create a P-curve (Statistical Table PDF) to select statistical results 3. Read statistical results to p-curve app 4. Copy/paste app's output into your paper 	<pre> function peeking(nsimon, n0, n1, every, d, seed, nk): r = random() k = 1 while k <= nk: k = k + 1 if (k % every == 0): data = peeking(nsimon, n0, n1, every, d, seed, nk) output = output + data else: data = peeking(nsimon, n0, n1, every, d, seed, nk) output = output + data return output </pre>