



*Marco Perugini* Rovereto, 20/05/2016



- Assume that, as scientists, we all want to get it right
- What can we do to increase our chances?

a) Get it right ≠ I am right
b) Get it right ≠ Get it published



- The last few years have probably seen more developments in research methodology than in the previous decades
- Rapid changes in standards for publishing
- Until 2011: Methodology = Boooring!
- From 2012: Methodology = Cool!
- Why?

### (see also Perugini 2014, GIP)



- The year 2011 has been an *annus horribilis* for Psychology
- Three main events:





# Bem (JPSP, 2011)

- Nine experiments showing ESP
- Strong reactions
- Initial article failing to replicate the results was refused by JPSP
- Hard questions on the modal way of analyzing data and on "cherry-picking" results
- Galak et al (2012): 7 failed replication attempts (n=3289)



- Resigned from Dean at Tilburg University (NL)
- Faked data: 53 retracted papers and 10 PhD thesis with invented or dubious data
- Levelt report (2012): proofs beyond doubts of faked data and strong criticisms to the scientific community
- Huge media impact



## Simmons et al (PS, 2011)

#### False-Positive Psychology: Undisclosed Flexibility in Data Collection and Analysis Allows Presenting Anything as Significant

Joseph P. Simmons<sup>1</sup>, Leif D. Nelson<sup>2</sup>, and Uri Simonsohn<sup>1</sup> <sup>1</sup>The Wharton School, University of Pennsylvania, and <sup>2</sup>Haas School of Business, University of California, Berkeley

- They show that commonly used questionable research practices can allow to provide empirical evidence even for null effects (false positive)
- Huge scientific impact (740 citations, 2<sup>nd</sup> most cited paper from 2011 in Psychology)



## And many other events...

- Bargh, Kanheman, Sanna, Smeesters, Simonsohn, Francis, and so on
- Until one buzzword came out

# **Replicability**

• and Psychology was born again ...



- If a result is not replicated, it is not valid
- To be replicated, it needs to be replicable *Replicability*
- A key concept in Science
- Almost forgotten in Psychology
- Now in the forefront
- What does it mean?



European Journal of Personality, Eur. J. Pers. 27: 108–119 (2013) Published online in Wiley Online Library (wileyonlinelibrary.com) DOI: 10.1002/per.1919

#### Recommendations for Increasing Replicability in Psychology<sup>†</sup>

JENS B. ASENDORPF<sup>1</sup>\*, MARK CONNER<sup>2</sup>, FILIP DE FRUYT<sup>3</sup>, JAN DE HOUWER<sup>4</sup>, JAAP J. A. DENISSEN<sup>5</sup>, KLAUS FIEDLER<sup>6</sup>, SUSANN FIEDLER<sup>7</sup>, DAVID C. FUNDER<sup>8</sup>, REINHOLD KLIEGL<sup>9</sup>, BRIAN A. NOSEK<sup>10</sup>, MARCO PERUGINI<sup>11</sup>, BRENT W. ROBERTS<sup>12</sup>, MANFRED SCHMITT<sup>13</sup>, MARCEL A. G. VANAKEN<sup>14</sup>, HANNELORE WEBER<sup>15</sup> and JELTE M. WICHERTS<sup>5</sup>

<sup>1</sup>Department of Psychology, Humboldt University Berlin, Berlin, Germany
 <sup>2</sup>Institute of Psychological Sciences, University of Leeds, Leeds, UK
 <sup>3</sup>Department of Developmental, Personality and Social Psychology, Ghent University, Ghent, Belgium
 <sup>4</sup>Department of Experimental Clinical and Health Psychology, Ghent University, Ghent, Belgium
 <sup>5</sup>School of Social and Behavioral Sciences, Tilburg University, Tilburg, The Netherlands
 <sup>6</sup>Department of Psychology, University of Heidelberg, Heidelberg, Germany
 <sup>7</sup>Max Planck Institute for Research on Collective Goods, Bonn, Germany
 <sup>8</sup>Department of Psychology, University of California at Riverside, Riverside, CA USA
 <sup>9</sup>Department of Psychology, University of Virginia, Charlottesville, VA USA
 <sup>10</sup>Department of Psychology, University of Milano-Bicocca, Milan, Italy
 <sup>12</sup>Department of Psychology, University of Koblenz–Landau, Landau, Germany
 <sup>14</sup>Department of Psychology, University of Koblenz–Landau, Landau, Germany
 <sup>15</sup>Department of Psychology, University of Greifswald, Greifswald, Germany



- The study should be described in a way such that everyone qualified can replicate it
- This implies a very detailed method section, including information that often is not disclosed
- Also, the data should be publicly available (at very least, upon request) to replicate the results using appropriate analyses (internal replication)
- Transparency in research



# **Back to basics (and beyond...)**

- As scientists, we all want to get something right
- If we get it right, it is replicable and will be replicated
- But what does it mean "to get it right"?
- What can we do to increase our chances?
- We need first to go back to some basic concepts



- The world is uncertain
- Knowledge is imperfect
- We deal with "samples" rather than "population" (no matter whether you use a Bayesian or a Frequentist approach)
- We try to make inferences from them
- Bertrand Russell's inductivist turkey

*"Essentially, all models are wrong, but some are useful" (Box & Draper, 1987)* 



- Mean, Standard deviation, Standard error
- Confidence intervals
- Effect size
- Errors of statistical inference
- Power analysis
- Replicability of psychological research



- A single value that reflects the central point of a distribution
- If the distribution is normal, it is also the best simple way to summarize it





• Reflects the dispersion (variability) around the mean





• When we measure something, the more the observed data, the less the measurement error

• For example, exit polls are more accurate (less error) the more the sampled voters or polling stations

• The point is that we have a sample but would like to say something about the underlying population (or anyway something that generalizes beyond that sample)



• Standard error does not depend only from how big is a sample size but also from the variability (variance) of the study object

• If everyone answers in the same way, one needs to ask to only one person...

• If people have very different opinions, one need many of them to be able to say something about «what they think»...

•Standard error provides a link between sample and population



• Suppose that most students have a mark between 26 and 27





• Almost all possible samples, even if small, will have values near to

the population mean of 26.5





• Suppose now that there is a lot of variability in the marks





• The samples, especially if small, can be very different from each

other and from the mean population value





• When we estimate a parameter (e.g., mean) in a sample we will make an *estimation error* of the population parameter

• The size of this error is given by the *standard error* 





#### Standard error





The sample estimate does not necessarily correspond to the population value It is possible to estimate a Confidence Interval that provides a range of values that contain the population value with a certain likelihood





The typical CI is 95%, meaning that there is a 95% likelihood that if we were to repeat the study 100 times, 95% of the times we would obtain one of the values within the interval including the sample estimate of the parameter (e.g., mean)

To simplify, CI 95% is roughly equal to the sample mean +/-2 SE.

For example M = 5; DS = 4 N = 100

SE = 
$$\sqrt{\frac{4^2}{100}}$$
 o  $\frac{4}{\sqrt{100}}$  =0.4 Range: 2 x SE = 0.8



A significant effect says little about the size of the effect
The p-value depends on sample size and not only effect size

• There are a number of ways to measure effect size

• Most common: Cohen's *d* 

$$d = \frac{M_1 - M_2}{\text{pooled SD}}, \text{ with pooled SD} = \sqrt{\frac{(n_1 - 1)SD_1^2 + (n_2 - 1)SD_2^2}{n_1 + n_2 - 2}},$$
  
If the two groups have the same sample size,  
$$pooled SD = \sqrt{\frac{SD_1^2 + SD_2^2}{2}}$$



### **Effect size**

#### ➢But also r (correlation coefficient)

$$r = \sqrt{\frac{t^2}{t^2 + df}} \qquad r = \sqrt{\frac{F(1,-)}{F(1,-) + df_R}} \qquad r = \sqrt{\frac{\chi^2(1)}{N}}$$

➢ From d to r and from r to d (e.g.

$$d = \frac{2r}{\sqrt{1-r^2}}$$



### **Example of Cohen's d**

$$pooled SD = \sqrt{\frac{(n_1 - 1)SD_1^2 + (n_2 - 1)SD_2^2}{n_1 + n_2 - 2}} = \sqrt{\frac{(246 - 1)1.107^2 + (219 - 1)1.197^2}{246 + 219 - 2}} = 1.150$$

$$d = \frac{M_1 - M_2}{pooled SD} = \frac{7.79 - 7.57}{1.150} = 0.19$$

Statistiche	di	gruppo
-------------	----	--------

	forma Forma.email	N	Media	Deviazione std.	Errore std. Media
persuasività Giudizio. Persuasività	1 Tu	246	7,57	1,107	,071
	2 Lei	219	7,79	1,197	,081

Conventional values:

0.2 small 0.5 medium 0.8 large





### **Other Effect Size indicators**

Common in factorial designs:

$$\hat{\eta}^2 = \frac{SS_{\text{Effect}}}{SS_{\text{T}}}, \quad \hat{\eta}_P^2 = \frac{SS_{\text{Effect}}}{SS_{\text{Effect}} + SS_{\text{s/Cells}}}, \quad \hat{\omega}_P^2 = \frac{SS_{\text{Effect}} - df_{\text{Effect}}MS_{\text{s/Cells}}}{SS_{\text{Effect}} + (N - df_{\text{Effect}})MS_{\text{s/Cells}}}$$

➤Some bibliografic references:

➢Fritz, C. O., Morris, P. E., & Richler, J. J. (2012). Effect size estimates: Current use, calculations, and interpretation. Journal of Experimental Psychology: General, 141, 2–18.

≻Olejinik & Algina (2003). Generalized eta and omega squared..., Psychological Methods

Cohen (1992). A power primer, Psychological Bulletin

Cohen (1994). The earth is round ( $\rho < .05$ ). American Psychologist

Cohen (1988). Statistical power analysis for the behavioral sciences



- Frequentist approach
- There are three types of errors
- NHST\*: Type I error (False positives) Type II error (False negatives)
- CI: Estimate error (imprecision)

NHST= Null Hypothesis Significance Testing (what you have been most likely taught as a student) H0 vs. H1



### **Errors of inference in NHST**

#### Real World (**POPULATION**)

Null is true (H0 is correct)Null is false (H1 is correct)





- Type I error: *Erroneously rejecting the null hypothesis* (*False positive*).
   The result in the sample is significant (*p* < .05), so the null hypothesis is rejected, but the null hypothesis is actually true in the population.
- **Type II error**: *Erroneously accepting the null hypothesis (<u>False negative</u>). The result in the sample is not significant (p > .05), so the null hypothesis is not rejected, but it is actually false in the population.*



- The Type I error rate (*False positive*) is controlled by the researcher.
- It is called the **alpha rate**, and corresponds to the probability cut-off that one uses in a significance test.
- Conventionally, researchers use an alpha rate (α) of .05. This means that the null hypothesis is rejected when a value such as the one found is likely to occur 5% of the time or less when the null hypothesis is true.
- The test can be two-tailed (more common) or one-tailed (directional)




- The Type II error (*False negative*) can also be controlled by the experimenter.
- The Type II error rate is called **beta**  $(\beta)$  as a complement to alpha.
- How can the beta rate be controlled? The easiest way to control Type II errors is by increase the **statistical power** of a test.
- **Statistical power**= probability of finding an effect, if it exists
- **Power** =  $1 \beta$
- Conventionally a power of at least .80 ( $\beta$ =.20) is considered as acceptable



- Power goes up with larger effect sizes and sample sizes, given a certain decision criterion (e.g.,  $\alpha$ =.05)
- When effect sizes become larger? When the portion of variability (difference) ascribed to the effect of interest grows more than the general (non specific) variability

$$d = \frac{M_1 - M_2}{pooled SD} \qquad \hat{\eta}^2 = \frac{SS_{Effect}}{SS_T}, \qquad r(v, x) = \frac{\operatorname{cov}(v, x)}{sd(v) * sd(x)}$$







• Sample size

• Construct-related (i.e., signal) variance

• Construct-unrelated (i.e., noise) variance





#### How to increase power?

#### Increase sample size

- Administer stronger treatments (e.g., experimental manipulation) BUT be wary of possible reduced ecological validity
- Avoid restrictions of range for dependent variables
- Standardize experimental procedures
- Increase reliability of measures
- Use more homogenous subject samples BUT increased risks to generalizability of results
- Use blocking or repeated measures (within) design BUT sometimes can be inappropriate
- Meta-analytical mindset



#### **Power Between Ss**





#### Gpower, http://www.gpower.hhu.de/en.html



### **Power Within Ss**

• Power for within Ss studies is greater (*ceteris paribus*) but depends on r (e.g., r = .50) between DVs





- Either you expect a large effect size or you need a substantial sample size
- The average ES in Psychology is at best **d=0.5**
- You can have a large sample (e.g., N=274) to reliably detect a small effect size (e.g., d=0.30)
- You can have a small sample (e.g., N=26) to reliably detect a very large effect size (e.g., d=1.0)
- But you <u>cannot</u> have a small sample size to detect a subtle effect (small effect size)
- Within Ss studies can be more powerful



- Frequentist approach
- There are three types of errors
- NHST: Type I error (False positives) Type II error (False negatives)
- CI: Estimate error (imprecision) (e.g., AIPE, Maxwell, ARP 2008)



- In a CI approach, one thing matters a lot: sample size, the bigger, the better (*ceteris paribus*)
- The point is not whether some effect exists (or not) but how precise is our estimate of it
- All effects exist given an infinite sample size (Cohen)
- If you want to get it right, increase sample size



### A practical problem

# • **<u>Big</u>** sample sizes are needed for precise estimates no matter the effect size

AIPE FOR THE STANDARDIZED MEAN DIFFERENCE





# Why NHST can be problematic

- But NHST dominates
- Publications,
   theories,
   discussions,
   fights and careers
   gravitate around
   p<.05</li>

"We are not interested in the logic itself, nor will we argue for replacing the .05 alpha with another level of alpha, but at this point in our discussion we only wish to emphasize that dichotomous significance testing has no ontological basis. That is, we want to underscore that, surely, God loves the .06



nearly as much as the .05. Can there be any doubt that God views the strength of evidence for or against the null as a fairly continuous function of the magnitude of p?"

Rosnow, R.L. & Rosenthal, R. (1989). Statistical procedures and the justification of knowledge in psychological science. American Psychologist, 44, 1276–1284.



- Psychology is now at the forefront of Science in addressing the replicability issue
- Methodological innovations
- Massive effort to estimate replicability
- The Reproducibility Project (OSC)
- 270 volunteers, 64 universities, 11 countries, 100 replicated studies published in three main journals (JPSP, PS, JEP:LMC) in 2008



# Why is it needed?

- Average ES: Cohen's d=0.50
- Average Sample size: n= 40
- Typical power: (1-β)=.35
- This means around 1/3 chance of positive findings
- The literature should be full of non-significant findings
- Really?

#### ACCENTUATE THE POSITIVE

A literature analysis across disciplines reveals a tendency to publish only 'positive' studies — those that support the tested hypothesis. Psychiatry and psychology are the worst offenders.



Proportion of papers supporting tested hypothesis

#### Bakker, van Dijk, & Wicherts, 2012; Asendorpf et al., 2013





- <u>NO !!</u> (Ioannidis, 2005)
- Average power in **Neuroscience**: .21 (Button et al., 2013) This means around 1/5 chance of positive findings
- Cancer Biology: Replication rate of main results from pre-clinical trails (Begley & Ellis, 2012): from 11% (Amgen, 2011) to 25% (Prinz et al. Bayer HC team, 2011)





What one can expect?





Effect itself against null (p-values)
 Effect size comparison
 Meta-analytic precision estimate

Open Science Collaboration (2015). Estimating the reproducibility of psychological science. *Science*, *349*(*6251*). DOI: 10.1126/science.aac4716



### **Effects against null I**

Figure S1: Cumulative P value distributions of original and replication studies.





# **Effects against null II**

Table 1. Summary of reproducibility rates and effect sizes for original and replication studies overall and by journal/discipline. *df*/*N* refers to the information on which the test of the effect was based (for example, *df* of *t* test, denominator *df* of *F* test, sample size –3 of correlation, and sample size for *z* and  $\chi^2$ ). Four original results had *P* values slightly higher than 0.05 but were considered positive results in the original article and are treated that way here. Exclusions (explanation provided in supplementary materials, A3) are "replications *P* < 0.05" (3 original nulls excluded; *n* = 97 studies); "mean original and replication effect sizes" (3 excluded; *n* = 97 studies); "meta-analytic mean estimates" (27 excluded; *n* = 73 studies); "percent meta-analytic (*P* < 0.05)" (25 excluded; *n* = 75 studies); and, "percent original effect size within replication 95% CI" (5 excluded, *n* = 95 studies).

			Effect size comparison					Original and replication combined			
	Replications P < 0.05 in original direction	Percent	Mean (SD) original effect size	Median original df/N	Mean (SD) replication effect size	Median replication df/N	Average replication power	Meta- analytic mean (SD) estimate	Percent meta- analytic (P < 0.05)	Percent original effect size within replication 95% Cl	Percent subjective "yes" to "Did it replicate?"
Overall	35/97	36	0.403 (0.188)	54	0.197 (0.257)	68	0.92	0.309 (0.223)	68	47	39
JPSP, social	7/31	23	0.29 (0.10)	73	0.07 (0.11)	120	0.91	0.138 (0.087)	43	34	25
JEP:LMC, cognitive	13/27	48	0.47 (0.18)	36.5	0.27 (0.24)	43	0.93	0.393 (0.209)	86	62	54
PSCI, social	7/24	29	0.39 (0.20)	76	0.21 (0.30)	122	0.92	0.286 (0.228)	58	40	32
PSCI, cognitive	8/15	53	0.53 (0.2)	23	0.29 (0.35)	21	0.94	0.464 (0.221)	92	60	53



### **Effect size comparison**



**Fig. 3. Original study effect size versus replication effect size (correlation coefficients).** Diagonal line represents replication effect size equal to original effect size. Dotted line represents replication effect size of 0. Points below the dotted line were effects in the opposite direction of the original. Density plots are separated by significant (blue) and nonsignificant (red) effects.



# **Meta-analytic comparison**

Table 1. Summary of reproducibility rates and effect sizes for original and replication studies overall and by journal/discipline. *df/N* refers to the information on which the test of the effect was based (for example, *df* of *t* test, denominator *df* of *F* test, sample size -3 of correlation, and sample size for *z* and  $\chi^2$ ). Four original results had *P* values slightly higher than 0.05 but were considered positive results in the original article and are treated that way here. Exclusions (explanation provided in supplementary materials, A3) are "replications *P* < 0.05" (3 original nulls excluded; *n* = 97 studies); "mean original and replication effect sizes" (3 excluded; *n* = 97 studies); "meta-analytic mean estimates" (27 excluded; *n* = 73 studies); "percent meta-analytic (*P* < 0.05)" (25 excluded; *n* = 75 studies); and, "percent original effect size within replication 95% CI" (5 excluded, *n* = 95 studies).

			Effect size comparison					Original and replication combined				
	Replications P < 0.05 in original direction	Percent	Mean (SD) original effect size	Mediar origina df/N	Mean (SD) replication effect size	Median replication df/N	Average replication power	Meta- analytic mean (SD) estimate	Percent meta- analytic (P < 0.05)	Percent original effect size within replication 95% CI	Percent subjective "yes" to "Did it replicate?"	
Overall	35/97	36	0.403 (0.188)	54	0.197 (0.257)	68	0.92	0.309 (0.223)	68	47	39	
JPSP, social	7/31	23	0.29 (0.10)	73	0.07 (0.11)	120	0.91	0.138 (0.087)	43	34	25	
JEP:LMC, cognitive	13/27	48	0.47 (0.18)	36.5	0.27 (0.24)	43	0.93	0.393 (0.209)	86	62	54	
PSCI, social	7/24	29	0.39 (0.20)	76	0.21 (0.30)	122	0.92	0.286 (0.228)	58	40	32	
PSCI, cognitive	8/15	53	0.53 (0.2)	23	0.29 (0.35)	21	0.94	0.464 (0.221)	92	60	53	



# Summing up RP main results

- 36% replicate at p<.05 (simple answer)
- Effect size are half (publication bias, file drawer effect)
- Less likely to replicate if:
  a) weaker evidence in original study (p<.05 worse than p<.001)</li>
  b) results considered to be "surprising" (*sexy but unreliable findings*)
- Milestone achievement of Psychology
- RP in Cancer Biology has just started



- As scientists, we all want to get something right
- If we get it right, it is replicable and will be replicated
- But what does it mean "to get it right"?
- So, what can we do to increase our chances?



- Design your study with adequate power
- It is ok to run initial exploratory/pilot studies (N needs not be too small) to identify an effect and to have a rough estimate of its effect size (don't trust it too much...)
- If you find something, then you need to plan a confirmatory study with adequate power to confirm your effect
- Be careful on overestimations of effect size (*Winner's curse*)



## Winner's curse I

- Inflated effect size from initial study are likely and can affect later confirmatory studies (Ioannides, 2008).
- More likely with asymmetries in publication standards (publish only significant results) and underpowered studies (difficult to establish a priori)
- Under these common conditions, true effect size can be much lower than published effect size
- d=0.80 (n=50), CI [0.22, 1.37]
- N for power 80% can be from 516 to 16, with n=42 for d=0.80

(cf. Safeguard power analysis - PGC, 2014)



## Winner's curse II

This is true also for your own pilot studies, because you will be likely to pursue significant findings and abandon not significant ones (your own "publication bias")



Contents lists available at ScienceDirect

Journal of Experimental Social Psychology



journal homepage: www.elsevier.com/locate/jesp

Exploring Small, Confirming Big: An alternative system to The New Statistics for advancing cumulative and replicable psychological research \*

John Kitchener Sakaluk \*

Department of Psychology, University of Toronto Mississauga, Canada

- It is not good advice
- Better advice "explore reliably, confirm more reliably"



#### Safeguard Power as a Protection Against Imprecise Power Estimates

Marco Perugini, Marcello Gallucci, and Giulio Costantini University of Milan-Bicocca, Italy

- 1. Calculate the effect size in the metric available or preferred for the study (e.g.,  $d_o$ )
- 2. Calculate the two-tail confidence interval (e.g., 60%) of this effect size measure.
- 3. Consider the lower boundary of the confidence interval (e.g., 20th percentile)<sup>5</sup>.
- Use this as the estimate of the true effect size (d<sub>s</sub>).
- 5. Calculate the needed sample size at a chosen power level (e.g., P = .80) for a given analysis at an 80% level of protection;  $N_p$  = sample size computed using standard power analysis;  $N_{s80}$  = sample size considered for safeguard power analysis at an 80% level of protection;  $N_{2.5}$  = sample size considered for safeguard power analysis at an 80% level of protection;  $N_{2.5}$  = sample size considered for safeguard power analysis at an 80% level of protection;  $N_{2.5}$  = sample size considered for safeguard power analysis at an 80% level of protection;  $N_{2.5}$  = sample size considered for safeguard power analysis at an 80% level of protection;  $N_{2.5}$  = sample size considered for safeguard power analysis at an 80% level of protection;  $N_{2.5}$  = sample size considered for safeguard power analysis at an 80% level of protection;  $N_{2.5}$  = sample size considered for safeguard for the 2.5 rule (i.e. and the safe sample size considered for safeguard for the 2.5 rule (i.e. and the safe sample size considered for safeguard for the 2.5 rule (i.e. and the safe sample size considered for safeguard for the 3.5 rule (i.e. and the safe sample size considered for safeguard for the 3.5 rule (i.e. and the safe sample size considered for safeguard for the 3.5 rule (i.e. and the safe sample size considered for safeguard for the safe sample size considered for sample size considered for safe sample size considered for sample sample size considered for sample s

**Table 3.** Application of Safeguard Power Analysis to Some

 Hypothetical Studies

Perspectives on Psychological Science

Study	$d_0$	$N_0$	$d_{\rm s80}$	$N_{\rm p}$	$N_{ m s80}$	N <sub>2.5</sub>	SSR <sub>80</sub>
А	0.70	50	0.45	52	124	125	2.48
В	0.70	200	0.58	52	76	500	0.38
С	1.0	30	0.66	28	58	75	1.93
D	1.5	20	1.05	14	24	50	1.20
Е	0.4	80	0.21	156	570	200	7.12

Note:  $d_0$  = effect size of the original study;  $N_0$  = sample size of the original study;  $d_{s80}$  = effect size considered for safeguard power

analysis at an 80% level of protection;  $N_{\rm p}$  = sample size computed using standard power analysis;  $N_{\rm s80}$  = sample size considered for safeguard power analysis at an 80% level of protection;  $N_{2.5}$  = sample size considered for the 2.5 rule (i.e., sample times 2.5); SSR<sub>80</sub> = safeguard sample ratio for safeguard power analysis at an 80% level of protection.

© The Author(s) 2014 Reprints and permissions: sagepub.com/journalsPermissions.nav DOI: 10.1177/1745691614528519 pps.sagepub.com

2014, Vol. 9(3) 319-332





- Distinguish between exploratory and confirmatory studies
- If you find an "unexpected effect", confirm it with another well powered study before building on it
- Results can be significant simply out of random sampling



# **Third pointer**

- Results stabilize with bigger sample sizes
- Try to have a decent sample size
- Sometimes results can be significant in opposite directions given small sample sizes
- For example, stability of correlation coefficients (cf. Schonbrodt & Perugini, 2013)



- Examples from my own research
- S1: one effect that exists in the full sample (r between H/H<sub>quest</sub> and H/H<sub>adjec</sub>=.48) and one that does not (r between H/H<sub>quest</sub> and  $Ext_{quest}$ =-.05)
- Correlations calculated adding Ss at each step starting from N=10 to full sample (evolution of r)
- Real Ss order
- Boostrapped (s=1000) CI 95%



## H vs. H

Correlation evolution for hon\_Adject & hon\_Hexaco





### H vs. H

Correlation evolution for hon\_Adject & hon\_Hexaco





# The results: H vs. E

Correlation evolution for hon\_Hexaco & ext\_Hexaco





#### H vs. E

Correlation evolution for hon\_Hexaco & ext\_Hexaco



![](_page_71_Picture_0.jpeg)

![](_page_71_Picture_1.jpeg)

Contents lists available at SciVerse ScienceDirect

#### Journal of Research in Personality

journal homepage: www.elsevier.com/locate/jrp

Brief Report

#### At what sample size do correlations stabilize?

#### Felix D. Schönbrodt <sup>a,\*</sup>, Marco Perugini <sup>b</sup>

<sup>a</sup> Department of Psychology, Ludwig-Maximilians-Universität, Leopoldstr. 13, 80802 München, Germany <sup>b</sup> Department of Psychology, University of Milan, Bicocca, Piazza dell'Ateneo Nuovo 1 (U6), 20126 Milan, Italy

![](_page_71_Figure_9.jpeg)


- Sequential effects can be devastating for small samples (N≤60)
- Estimates start to stabilize for N≥100 (but it strictly depends on the expected correlation; e.g., N≈180 for r = .4)
- Small samples (N≤60) can give many false positives/negatives, especially for small effects (but see Sequential testing with Bayes Factors, *Schönbrodt*, *Wagenmakers, Zehetleitner, & Perugini, in press, PM*)



### **Practical implications**

Lakens, D., & Evers, E. R. K. (2014). Sailing from the seas of chaos into the corridor of stability: Practical recommendations to increase the informational value of studies. *Perspectives in Psychological Science*, 9, 278–292. doi:10.1177/ 1745691614528520

 $d = \frac{2r}{\sqrt{1 - r^2}}$ 

In columns 3 and 4 of Table 1 on p. 280, the recommended sample sizes are for the total sample, not for the sample size per condition. The corrected table appears below.

**Table 1.** Recommended Sample Size Per Condition When Comparing Two Independent Groups Based on Different Effect Sizes (r and Cohen's  $d_{pqp}$ ) to Achieve the Point of Stability (POS) With 80% Confidence and Corridor Widths of .2 and .1 (See Part 1), to Achieve 80% or 90% Power to Observe the Effect With an Alpha of .05 and to Achieve a v Statistic Higher Than .5 (See Part 3)

r	dpop	POS, 80% w = .2	POS, 80% w = .1	80% power	90% power	v > .5
.1	0.20	31	126	394	527	404
.2	0.41	29	119	95	126	99
.3	0.63	26	106	41	54	43
.4	0.87	22	91	22	29	23
.5	1.15	17	72	13	17	14
.6	1.50	13	52	9	11	9
.7	1.96	10	33	6	7	6



- Results stabilize with smaller standard errors
- Standard errors depend on N and SD
- Smaller SD means smaller SE
- The SD can be reduced (*ceteris paribus*) with more reliable measures, more precise experimental design, less within Ss variability
- Plan your design as simple and as clean as possible



## **Fifth pointer**

- Be very careful with Questionable Research Practices (QRP)
- Do not cherry-pick DVs among many that you have
- Do not exclude cases as is
- Do not make multiple interim analyses to decide whether to collect additional Ss
- Read Simmons et al. (2011): not every recommendation is perfect, but they do give many good ones



## **Sixth pointer**

- Some of this stuff is already implemented in top level journals
- For example, Psychological Science

#### DISCLOSURE QUESTIONS:

For all studies in your recently published article titled [publication title], please endorse the following statements: (please type an X to indicate your answer)

We reported the total number of observations which were excluded (if any) and the criterion for doing so. (If no observations excluded, please indicate Yes) Yes: \_\_\_\_\_No: \_\_\_\_

If no, please report this information here (e.g., data from 3 participants in Study 2 excluded due to computer malfunction; 4 participants in Study 1 excluded for not following instructions):

We reported all tested experimental conditions, including failed manipulations. Yes: \_\_\_\_ No: \_\_\_\_

If no, please provide brief explanation for not reporting this information (e.g., critical software implementation error; editorial request):

We reported all administered measures/items. Yes: \_\_\_\_ No: \_\_\_\_

If no, please provide brief explanation for not reporting this information (e.g., measures not related to research question; scores from unreported measure insufficiently reliable):

We reported (a) how we determined our sample size and (b) our data collection stopping rule. Yes: \_\_\_\_ No: \_\_\_\_

If no, please describe (a) the basis for the sample sizes used and (b) how you decided to stop collecting data (e.g., decided ahead of time to collect data until minimum sample size achieved and this was followed; sample size determined by power analysis but did not achieve it by the end of term):



- a) Average effect size in published studies are small to medium (d  $\approx 0.50$ )
  - b) Average sample sizes are small (n  $\approx 40$ )
  - c) Typical power is low  $(1-\beta \approx .35)$
  - d) In neuroscience even lower (1- $\beta \approx .21$ )
- This means that, assuming that all effects are true, there should be approx. 35% of positive findings but they actually are **over 90%** and above all other sciences (Fanelli, 2012)
- How this is possible?







### How To Win Lotto 8.7 Times Out Of Every 10 Times?



For The Most Respected Lottery System On The Internet 100% Guaranteed

# MIRACLES HAPPEN



- **Increase sample size if you want to get it right**
- Decrease "noise" in the study
- Dichotomous thinking does not help: everything happens under some circumstances
- Ask what, when, how much, how something happens
- Get it right  $\neq$  I am right
- To get it right means to reduce False positives (Type I error), False negatives (Type II error) and to have reasonably precise estimates



### Tversky & Kanheman (1971) (Belief in the law of small numbers, PB, 76, 105-110

We submit that people view a sample randomly drawn from a population as highly representative, that is, similar to the population in all essential characteristics. Consequently, they expect any two samples drawn from a particular population to be more similar to one another and to the population than sampling theory predicts, at least for small samples.

In review, we have seen that the believer in the law of small numbers practices science as follows:

 He gambles his research hypotheses on small samples without realizing that the odds against him are unreasonably high. He overestimates

power.

2. He has undue confidence in early trends (e.g., the data of the first few subjects) and in the stability of observed patterns (e.g., the number and identity of significant results). He overestimates significance.

 In evaluating replications, his or others', he has unreasonably high expectations about the replicability of significant results. He underestimates the breadth of confidence intervals.

4. He rarely attributes a deviation of results from expectations to sampling variability, because he finds a causal "explanation" for any discrepancy. Thus, he has little opportunity to recognize sampling variation in action. His belief in the law of small numbers, therefore, will forever remain intact.