# Sample planning

# Sample planning

- **When you plan a study/research/intervention, you should think about the participants that you need**

  **Some basic issues**

- **Representativeness** (inference from sample to population)

- **Generalizability** (inference from sample to population)

- **Robustness** (inference that you can "trust")

- **Feasibility** (something that you can do)

- **Efficiency** (most value for money, "spend" as little as possible)

# Mindsets for sample planning

- **Accuracy**: collect as many participants as needed to have a certain level of accuracy in your parameter estimation

- **Efficiency**: collect as <u>few</u> participants as needed to reach the conclusion that you want to reach

- **Redundancy**: collect as many participants are needed to reach the conclusion that you want to reach with not too much error (trade-off between accuracy and efficiency)

# Some statistical approaches to sample size planning

- **AIPE** (Maxwell, 2008): decide sample size based on a chosen level of Accuracy In Parameter Estimation

- **Sequential designs:**

  **Frequentist** (Lakens, 2014): Start with a planned N and number of interim tests, add N if needed (but adjust alpha)

  **Bayesian** (Schonbrodt et al., 2017, 2018): Start with a minimum N, add N until BF reaches a pre-defined threshold

- **Heuristic:** in different fields there are "magical" rules (N≥20 per cell, N>100, ratio k/N). At best, approximate wise suggestions, at worse misleading

- **Power analysis:** Today and tomorrow

# Some references

UNIVERSITÀ DEGLI STUDI DI MILANO BICOCCA

## Sample Size Planning for Statistical Power and Accuracy in Parameter Estimation

Scott E. Maxwell,[1] Ken Kelley,[2] and Joseph R. Rausch[3]

[1]Department of Psychology, University of Notre Dame, Notre Dame, Indiana 46556;
email: smaxwell@nd.edu
[2]Inquiry Methodology Program, Indiana University, Bloomington, Indiana 47405;
email: kkiii@indiana.edu
[3]Department of Psychology, University of Minnesota, Minneapolis, Minnesota 55455;
email: rausch@umn.edu

**Key Words**

## Special issue article: Methods and statistics in social psychology: Refinements and new developments

### Performing high-powered studies efficiently with sequential analyses

DANIËL LAKENS*

Human Technology Interaction Group, Eindhoven University of Technology, Eindhoven, The Netherlands

**BRIEF REPORT**

## Bayes factor design analysis: Planning for compelling evidence

Felix D. Schönbrodt[1] · Eric-Jan Wagenmakers[2]

## Sequential Hypothesis Testing With Bayes Factors: Efficiently Testing Mean Differences

Felix D. Schönbrodt
Ludwig-Maximilians-Universität München

Eric-Jan Wagenmakers
University of Amsterdam

Michael Zehetleitner
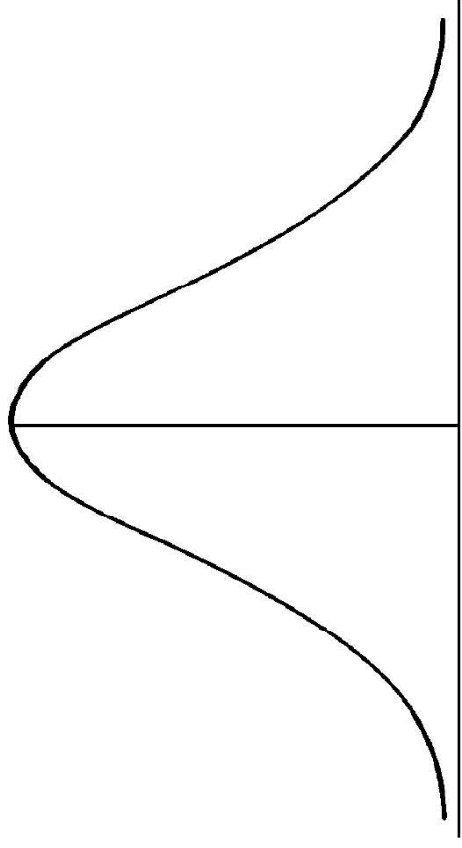Ludwig-Maximilians-Universität München

Marco Perugini
University of Milan–Bicocca

# Basic statistical concepts

# Mean

- **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**

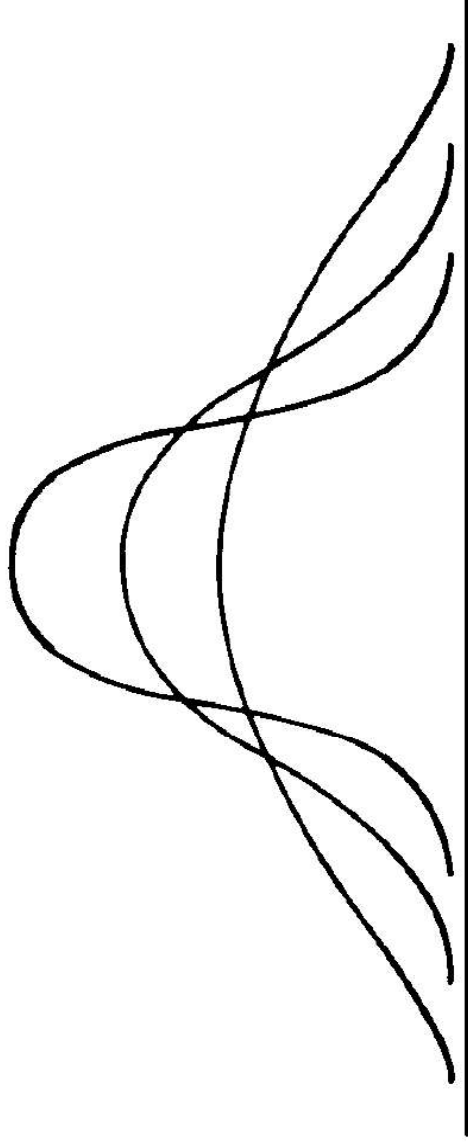$$\overline{X} = \frac{\sum X_i}{N}$$

# Variance and standard deviation

- Reflects the dispersion (variability) around the mean

$$s^2 = \frac{\sum (X_i - \bar{X})^2}{N} = \frac{\sum X^2}{N} - \bar{X}^2$$

$$s = \sqrt{s^2}$$

# Standard error

- When we measure something, more data means less measurement error

- Exit polls are more accurate (less error) the more the sampled voters or polling stations

- We have a sample but would like to say something about the underlying population (or anyway something that generalizes beyond that sample)

# Standard error and variance

- 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»…

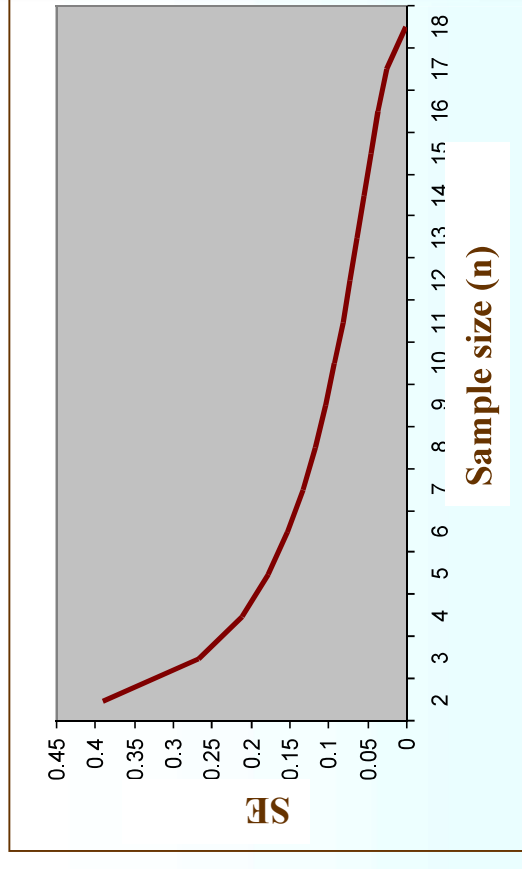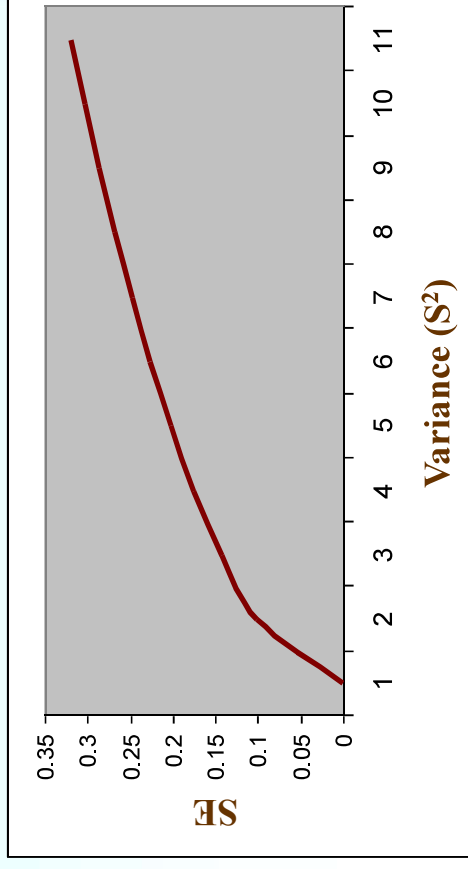- Standard error provides a link between sample and population

# Standard error

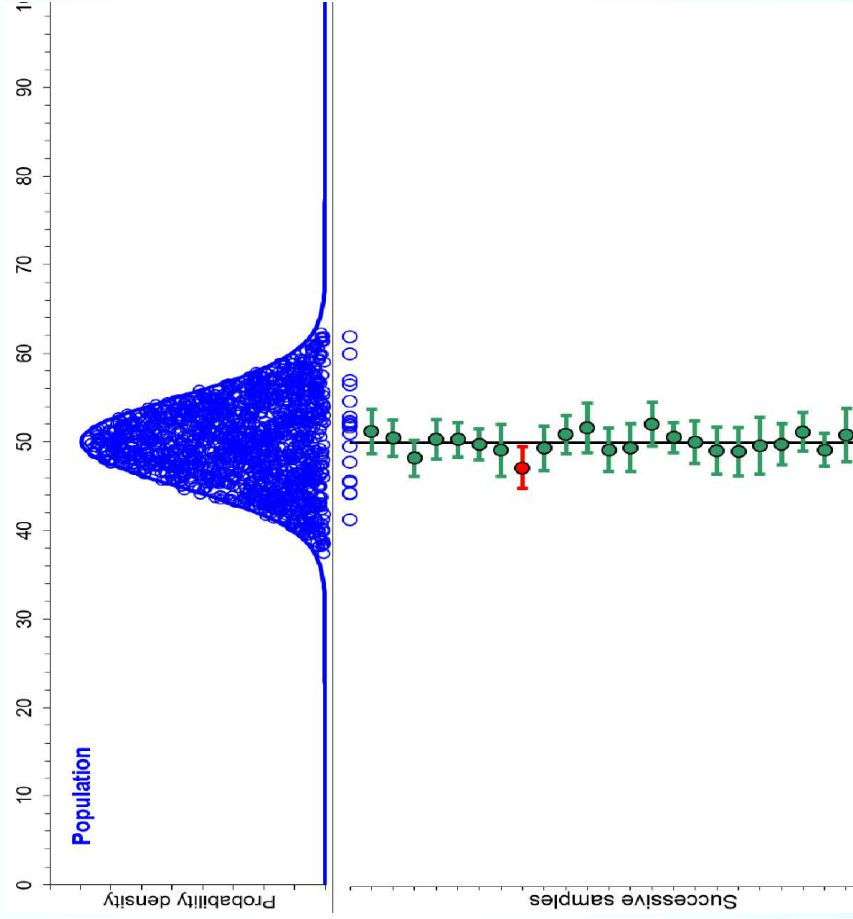➤ Error in estimating a population parameter (e.g., mean) from a sample

$$SE = \sqrt{\frac{S^2}{n}}$$

Goes up with increasing variance
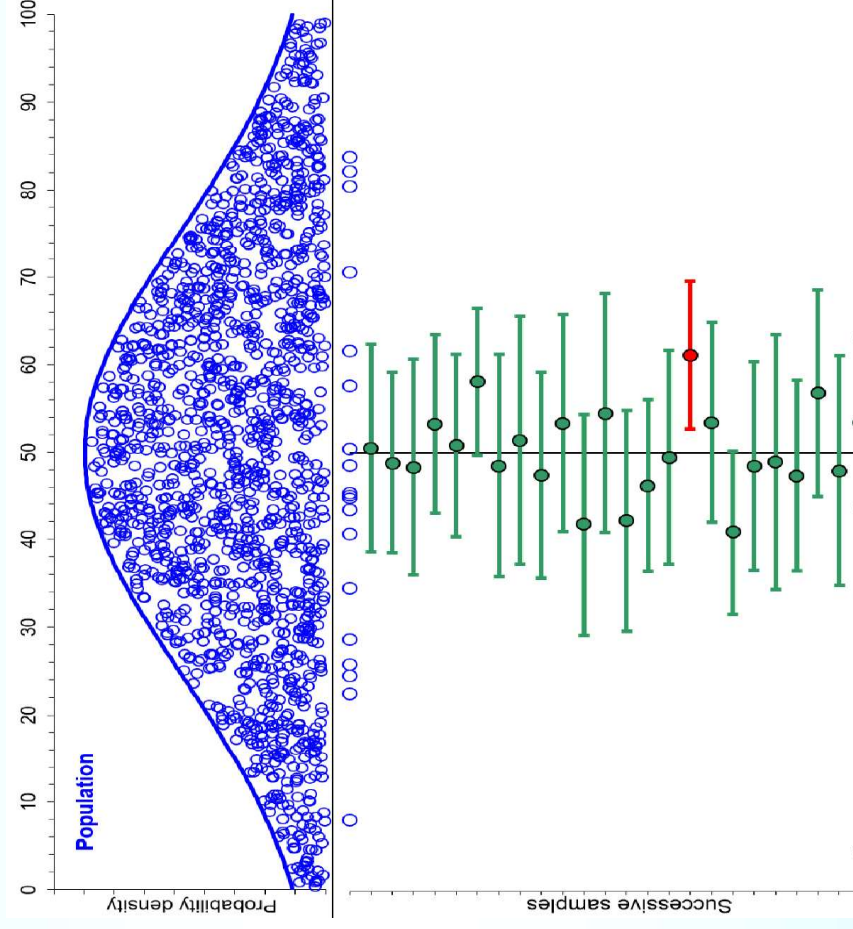
Goes down with increasing sample size

# Parameter estimation:
## Error and variability



N=20

N=20

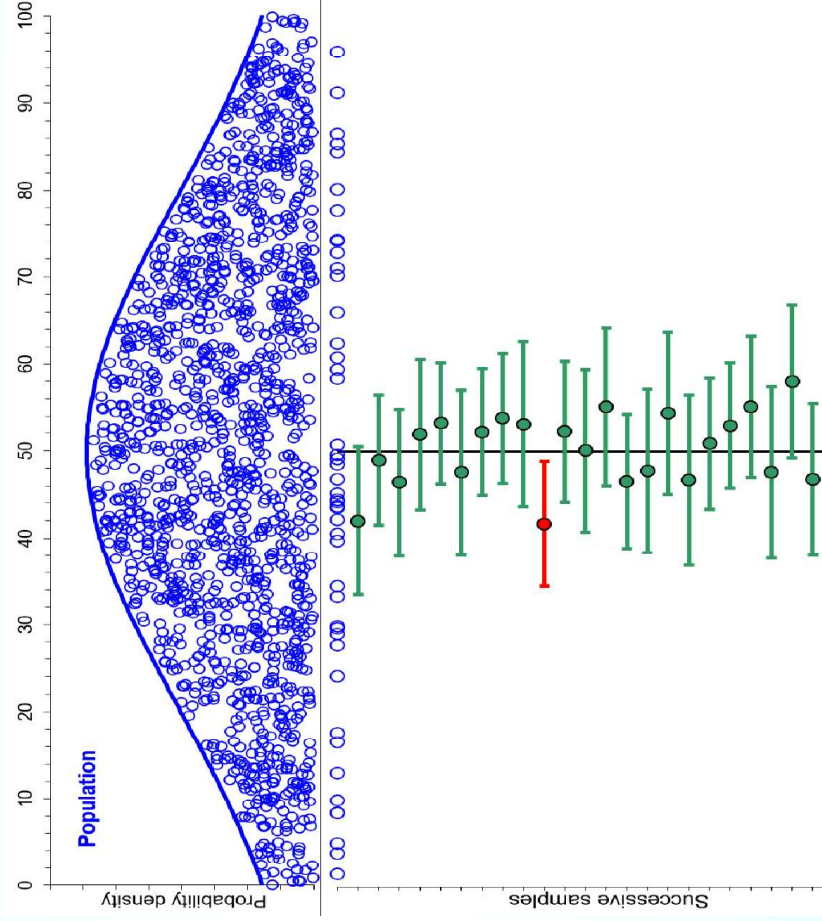Small variability = small SE

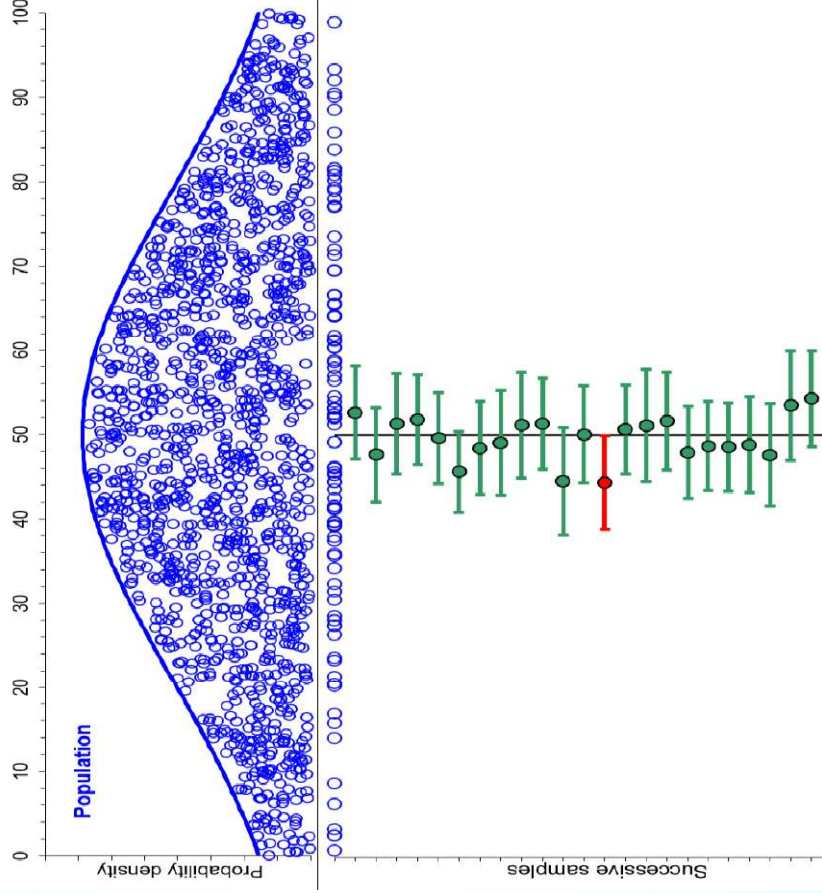Large variability = large SE

# Error and variability



N=50

N=100

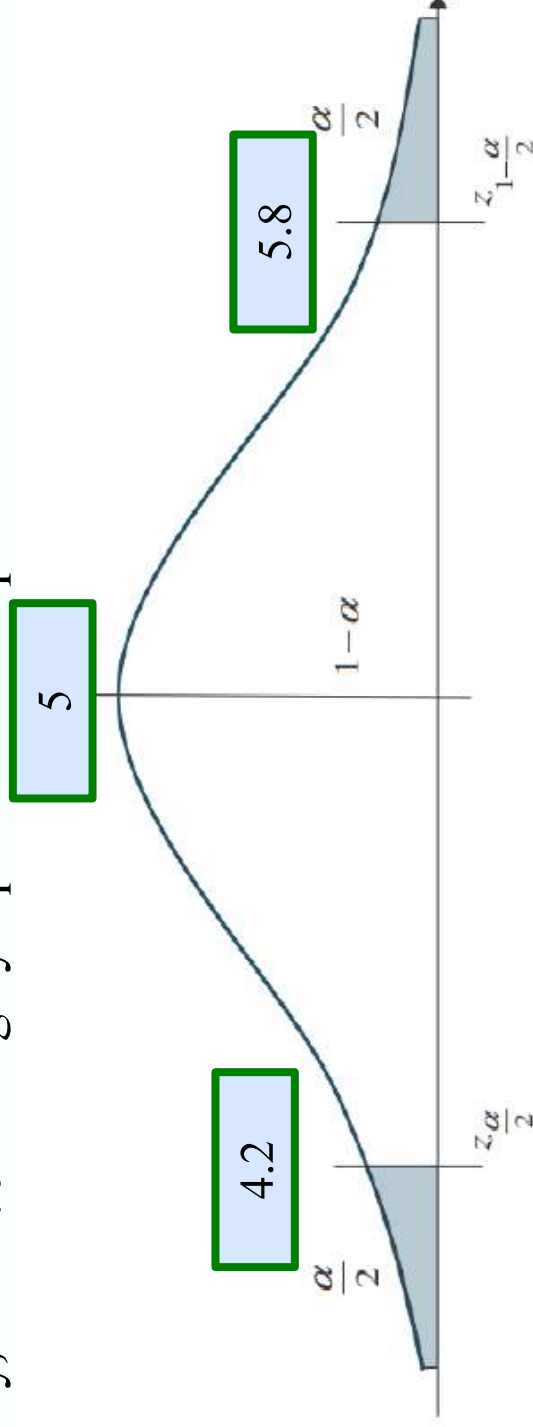Large variability = large SE

Small variability = small SE

# Error and variability

- **Remember**: basically you have almost always results from samples and not from populations

- There is an error in inferring results from samples as if they apply to a population

- Greater variability means more errors

# From SE to Confidence Interval (CI)

The sample estimate does not correspond to the population value.
Confidence Interval provides a range of values that contain the population value with a certain likelihood (e.g., 95%), should the study be repeated many times
To simplify, CI 95% is roughly equal to the sample mean +/- 2 SE

4.2

5

5.8

$$\frac{\alpha}{2}$$

$$1-\alpha$$

$$\frac{\alpha}{2}$$

$$z_{\frac{\alpha}{2}}$$

$$z_{1-\frac{\alpha}{2}}$$

Standard
Error

$$\mu \in \left( \left[ \overline{X}_n \pm t^{(n-1)}_{1-\frac{\alpha}{2}} \sqrt{\frac{\overline{s}_n^2}{n}} \right. \right.$$

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

$$SE = \sqrt{\frac{4^2}{100}} \text{ o } \frac{4}{\sqrt{100}} = 0.4$$

Range: 2 x SE = 0.8
95% CI = [4.2, 5.8]

# The Confidence Interval (CI)

The CI reflects the concept of **accuracy** in estimating a parameter

Imagine this research scenario. We want to understand the efficacy of 2 ads for a product (e.g., snack). N=100

We computed the mean evaluation of the two ads

A)  M = +3.10; DS = 15, p<.05
B)  M = +2.50; DS = 10, p<.05

Which is the best ad? It is not obvious that it is A

A) 95% CI= [0.16, 6.04]
B) 95% CI= [0.54, 4.46]

A can be 6.04, but it can also be 0.16.

B is more accurate, so its possible values are less spread: it is very unlikely that its mean is lower than 0.54

# Also correlations have confidence intervals

- Confidence intervals can be calculated for many statistical parameters
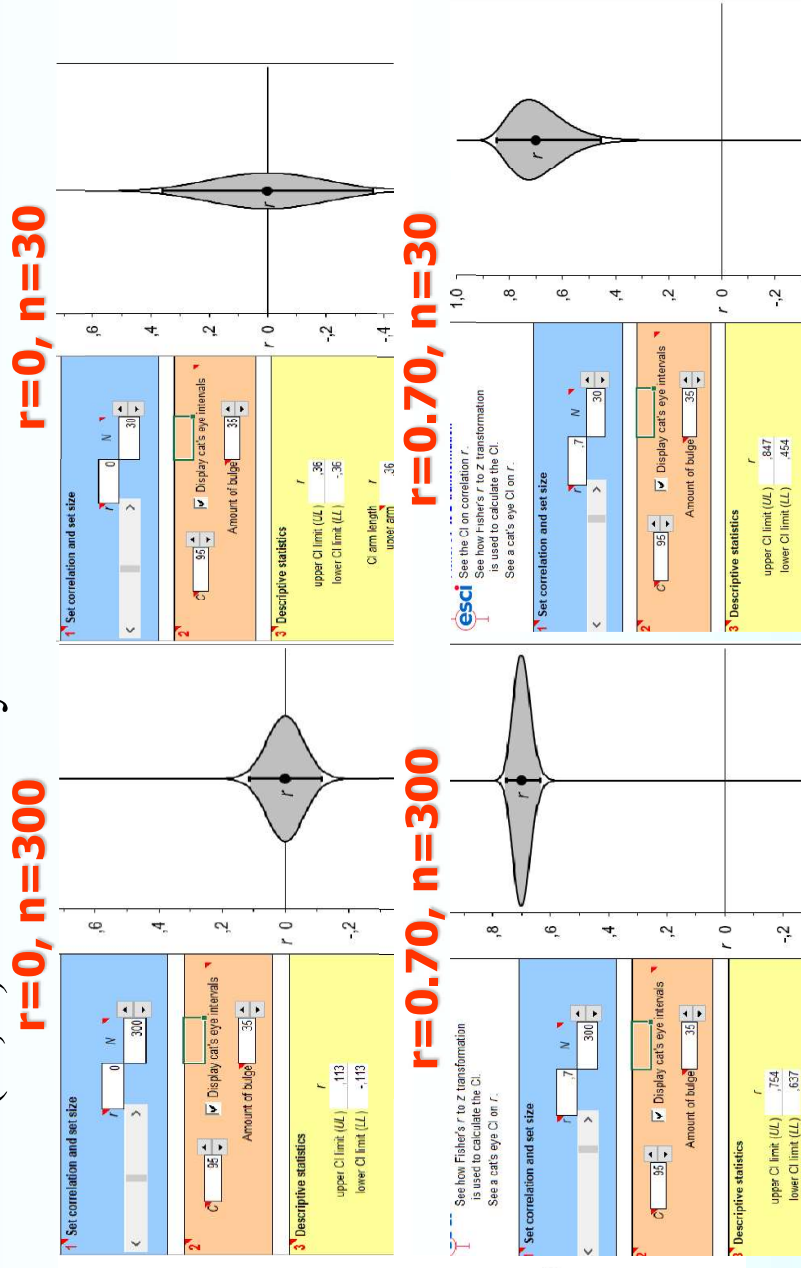- CI for correlations (r) are bounded (-1, 1) and often asymmetrical



r=0, n=30

r=0, n=300

r=0.70, n=30

r=0.70, n=300

Define Fisher Transformation: $z_r = \dfrac{\ln\left(\frac{1+r}{1-r}\right)}{2}$

Define: $L = z_r - \dfrac{z_{1-\frac{\alpha}{2}}}{\sqrt{N-3}}$

$U = z_r + \dfrac{z_{1-\frac{\alpha}{2}}}{\sqrt{N-3}}$

The 100(1-α)% confidence interval is defined as:

$\left(\dfrac{e^{2L}-1}{e^{2L}+1}, \dfrac{e^{2U}-1}{e^{2U}+1}\right)$

# Hypothesis testing I

# Hypothesis testing

- When we have data, we can estimate some parameters from them (e.g., mean, correlation)

- We saw that the estimate of this parameter can be more or less **accurate**

- But we can also make inferences from the estimated parameter

- If the parameter is different from a certain value (e.g., 0)

- If the parameter is different between certain groups (e.g., experimental vs. control, male vs. female)

- This is the realm of **hypothesis testing** (or statistical inferences from data)

# Statistical inference on a parameter

- After computing a parameter, we wish to exclude the possibility that the observed values are not simply due to chance (e.g, their difference from 0; their difference between groups)

  **Statistical Inference**

- Null-Hypothesis Significance Testing (NHST)

# Testing hypotheses on a parameter

- Suppose I find in a sample (N=100) that A (alcohol consumption) and B (aggressive behavior) are correlated r=.253. I want to say that there is a relationship between consuming alcohol and being aggressive.

- You do not believe it, and you think **it was a fluke,** but in reality there is no relationship between the two variables (r=0): one can be aggressive without drinking alcohol or drink alcohol without becoming more aggressive

- I want to show you that it was not a fluke. I use a Sherlock Holmes–like approach: *if this result is by chance, how likely would I be to find it?* Depending on the results, you could become less skeptical

II

THE PRINCIPLES OF EXPERIMENTATION, ILLUSTRATED BY A PSYCHO-PHYSICAL EXPERIMENT

### 5. Statement of Experiment

A LADY declares that by tasting a cup of tea made with milk she can discriminate whether the milk or the tea infusion was first added to the cup. We will consider the problem of designing an experiment by means of which this assertion can be tested. For this purpose
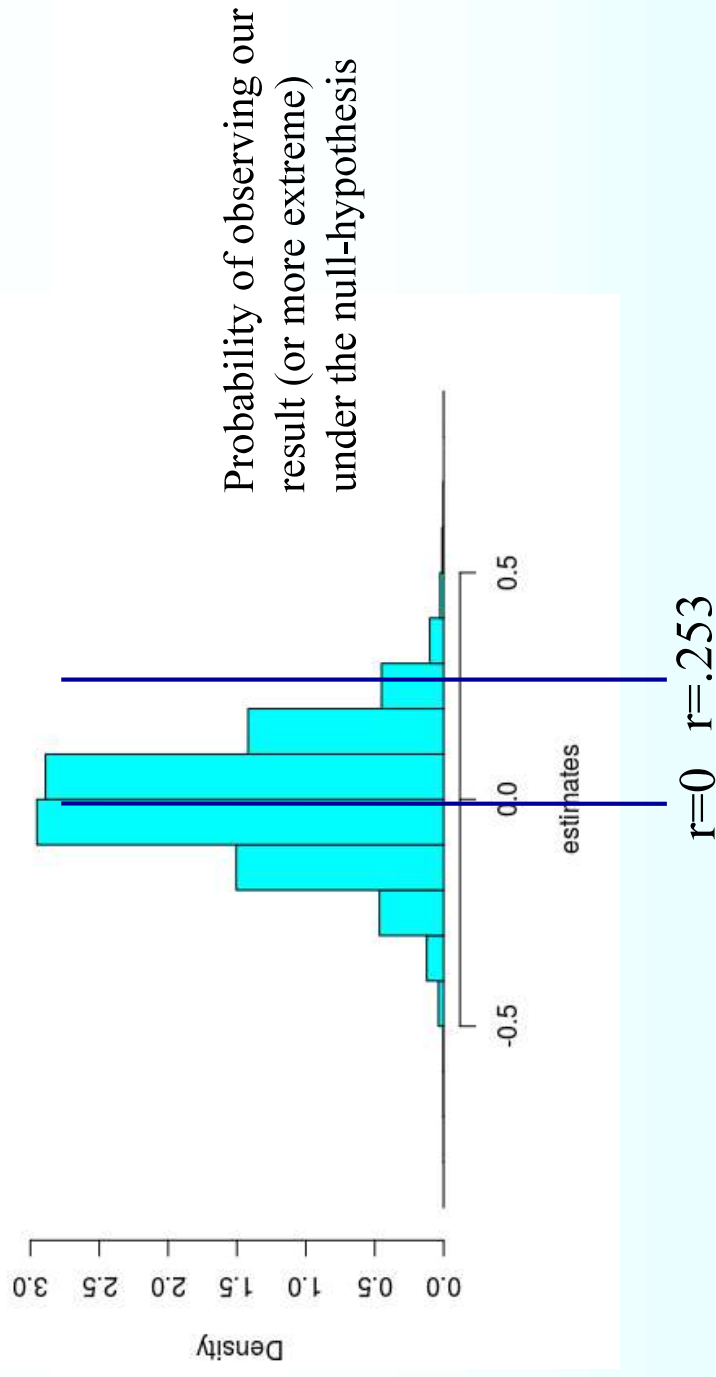
# Null-hypothesis: the what-if scenario

- Let assume you are right (it was a fluke), and the actual relation between hours of studying and grade is null

- This means that in the **population,** r=0. Let's call this population value the null-hypothesis population

- What if I could repeat my study hundreds of time, drawing each time an equivalent sample from the **null-hypothesis population.**

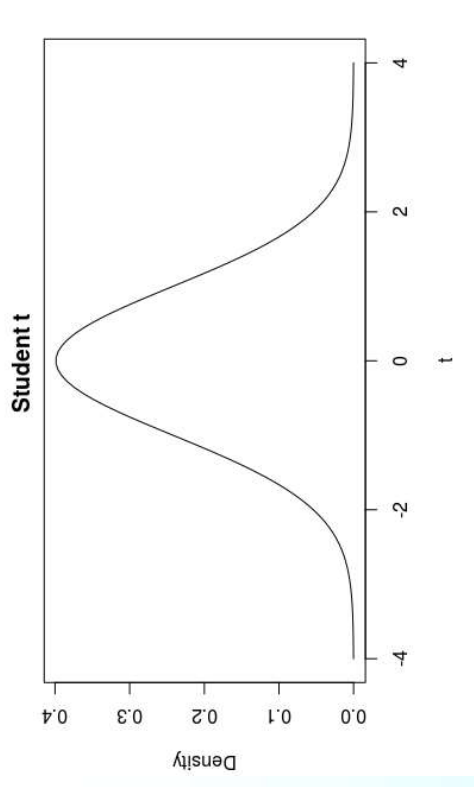- Every time I would have a new estimate of my r in any given sample



SHEER GUESS WORK?

# Null-hypothesis: the what-if scenario, should r be 0

- We could ask how likely is to obtain the result, r=.253, under the null-hypothesis scenario (that is, if r=0 in the population) by extracting thousands of samples and looking at the distribution of the results



Probability of observing our result (or more extreme) under the null-hypothesis

r=0   r=.253

# Null-hypothesis: a priori distribution

- Under the given circumstances (the assumptions of the test) we know a-priori what will be the distribution of all possible estimates under the null-hypothesis

- We know that if we keep resampling from a **normal** distribution where r=0, we can use a Student-t distribution to estimate the probability of values
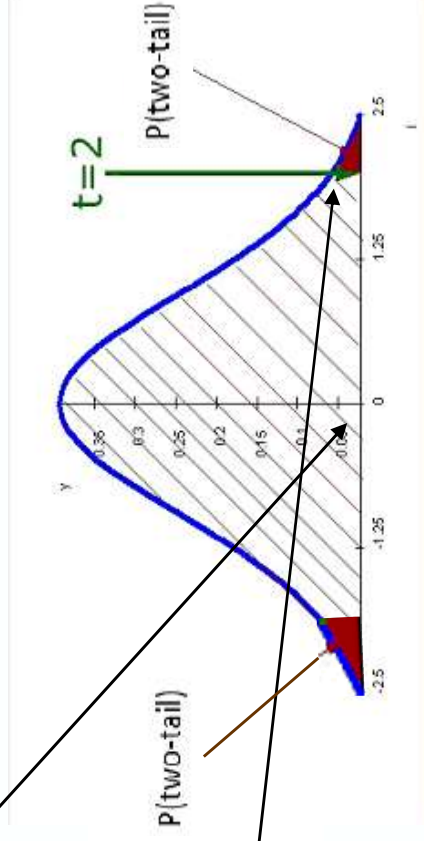
# t-test

● Thus we can compare our result with the expected distribution under the null hypothesis

● Null hypothesis: r=0

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



● This is a t-test: p is the probability of obtaining our observed result (or larger) if the null hypothesis is true. Here t=2.589, and p=.011. What does it mean?

# p-value: How small is small?

- Conventional cut-off values are **0.05** (5% of error), and **0.01** (1% of error). Newer approaches suggest **0.005** in some circumstances (see after).

- These values are almost always used in social science research. It means that if you reject the null hypothesis if and only if the probability (**p**) of your value is less than 0.05 (or 0.01), you will be mistaken at most 5% (or 1%) of the times.

- Although these values represent a widespread standard, they are only a convention. In Physics they have the 5-sigma convention

**2-sigma:** 95.5 percent

**3-sigma:** 99.73 percent

**4-sigma:** 99.993 percent

**5-sigma:** 99.99994 percent

So that means that purely statistical fluctuations will give you a result way out in the 5-sigma range 0.00006 percent of the time.

When physicists announce that they have a 5-sigma result, that means that there's a 1 in 3.5 million chance that it was the result of a statistical fluctuation over the spectrum of experiments they performed. Particle physicists working on the CMS and ATLAS experiments are

# How to interpret a p value

Yeahhhh!! I have a p=.003, which is statistically significant (<.05)!! This should be interpreted considering that:

a) p is the probability that the results are due to chance, the probability that the null hypothesis (H0) is true.

b) p is the probability that the results are not due to chance, the probability that the null hypothesis (H0) is false.

c) p is the probability of observing results as extreme (or more) as observed, if the null hypothesis (H0) is true.

d) p is the probability that the results would be replicated if the experiment was conducted a second time.

e) None of these.

# Recap of terms

- The hypothesis which describes the effect of chance, is called the **null hypothesis (H0)**

- The probability of obtain our result (or even more extreme) if the null hypothesis (**H0**) is true is called **p**

- The cut-off (0.05 or 0.01) which we use to reject or not the null hypothesis is called **critical alpha**

- The operation which leads us to the decision is called **test of significance**

- If we reject the null hypothesis (**H1**), we say that our result is **statistically significant at level p (or below p)**

- If we do not reject the null hypothesis, we say that our result is **not significant**

# Interpretation

**Statistically significant** does not mean scientifically significant, interesting or even important!

- Statistically significant means that we can exclude (with an error of **p**) that our result is equivalent to a purely random effect

- Statistically significance can often be considered a necessary condition for a result to be of interest, but it is not sufficient

- The importance of a result depends on its practical and theoretical relevance, its **strength** and direction

# Alternative Interpretation

Subgroup analysis

https://xkcd.com/882/

| P-VALUE | INTERPRETATION |
|---|---|
| 0.001 | |
| 0.01 | HIGHLY SIGNIFICANT |
| 0.02 | |
| 0.03 | |
| 0.04 | SIGNIFICANT |
| 0.049 | |
| 0.050 | OH CRAP. REDO CALCULATIONS. |
| 0.051 | ON THE EDGE OF SIGNIFICANCE |
| 0.06 | |
| 0.07 | HIGHLY SUGGESTIVE, SIGNIFICANT AT THE P<0.10 LEVEL |
| 0.08 | |
| 0.09 | |
| 0.099 | HEY, LOOK AT THIS INTERESTING SUBGROUP ANALYSIS |
| ≥0.1 | |

# Some suggested recent readings

Lakens, D. (2021). The practical alternative to the p value is the correctly used p value. *Perspectives on psychological science, 16*(3), 639-648.

Wasserstein, R-L., & Lazar, N.A (2016) The ASA's statement on p-values: Context, process, and purpose. *The American Statistician, 70*, 129-133.

Wasserstein, R. L., Schirm, A. L., & Lazar, N. A. (Eds.). (2019). Statistical inference in the 21st century: A world beyond p < 0.05 [Special issue]. *The American Statistician, 73, 1-19*.

EDITORIAL

Taylor & Francis
Taylor & Francis Group

## The Practical Alternative to the *p* Value Is the Correctly Used *p* Value

Daniël Lakens
Department of Industrial Engineering and Innovation Sciences, Eindhoven University of Technology

## The ASA's Statement on *p*-Values: Context, Process, and Purpose

EDITORIAL

Moving to a World Beyond "*p* < 0.05"

PSYCHOLOGICAL SCIENCE

$SAGE

Taylor & Francis
Taylor & Francis Group

∂ OPEN ACCESS

Check for updates

*Effect sizes*

# But literally H0 is never true...

- Given an infinite sample size, two parameters (e.g., means) will always be significantly different unless they are exactly identical, or one parameter will always be different from zero unless it is exactly zero (cf. **standard error**)

*r = .01 with N=40000 is significantly different from 0 with p<.05 (p=.0456)*

- It is thus important to understand the **effect size** (even if significant, some effects can be of a trivial quantity)

- Different effect size estimators

- Most common: **Cohen's d** and Pearson's **r** (correlation coefficient)

$$Cohen's\ d = \frac{M_1 - M_2}{SD_{pooled}}$$

$$SD_{pooled} = \sqrt{\frac{(SD_1^2 + SD_2^2)}{2}}$$

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

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

$$r = \sqrt{\frac{F(1,-)}{F(1,-) + df_R}}$$

# Effect size: examples

➤ A: ad product

➤ B: control group (irrelevant ad)

• VD: Product evaluation (from 0 to 10)

• A (n=60) : M= 6.50, SD=1.20

• B (n=60) : M= 5.50, SD=1.30

• $SD_{pool} = 1.25$

• $Cohen's\ d = \dfrac{6.50-5.50}{1.25} = 0.80 \qquad r=0.37$

**If** A: M= 7.50, DS=1.20; B: M= 5.50, DS=1.30

$SD_{pool} = 1.25,\ d = \dfrac{7.50-5.50}{1.25} = \mathbf{1.60} \quad \mathbf{r = 0.62}$

**If** A: M= 6.50, DS=2.20; B: M= 5.50, DS=2.30

$SD_{pool} = 2.25,\ d = \dfrac{6.50-5.50}{2.25} = \mathbf{0.44} \quad \mathbf{r = 0.22}$

**Rough guidelines (ES should be understood within research context)**

$r = .1$, $d = 0.2$ (small effect): the effect explains 1% of the total variance.
$r = .3$, $d = 0.5$ (medium effect): the effect explains 9% of the total variance.
$r = .5$, $d = 0.8$ (large effect): the effect explains 25% of the variance.

# Other effect size indexes (from Ellis, 2010)

Table 1.1 *Common effect size indexes*

| Measures of group differences (the d family) | | Measures of association (the r family) |
|---|---|---|
| **(a) Groups compared on dichotomous outcomes** | | |
| RD | The risk difference in probabilities: the difference between the probability of an event or outcome occurring in two groups | |
| RR | The risk or rate ratio or relative risk: compares the probability of an event or outcome occurring in one group with the probability of it occurring in another | |
| OR | The odds ratio: compares the odds of an event or outcome occurring in one group with the odds of it occurring in another | |
| **(b) Groups compared on continuous outcomes** | | |
| d | Cohen's d: the uncorrected standardized mean difference between two groups based on the pooled standard deviation | |
| Δ | Glass's delta (or d): the uncorrected standardized mean difference between two groups based on the standard deviation of the control group | |
| g | Hedges' g: the corrected standardized mean difference between two groups based on the pooled, weighted standard deviation | |
| PS | Probability of superiority: the probability that a random value from one group will be greater than a random value drawn from another | |

Table 1.1 (*cont.*)

| Measures of group differences (the d family) | Measures of association (the r family) | |
|---|---|---|
| | **(a) Correlation indexes** | |
| | $r$ | The Pearson product moment correlation coefficient: used when both variables are measured on an interval or ratio (metric) scale |
| | $\rho$ (or $r_s$) | Spearman's rho or the rank correlation coefficient: used when both variables are measured on an ordinal or ranked (non-metric) scale |
| | $\tau$ | Kendall's tau: like rho, used when both variables are measured on an ordinal or ranked scale; tau-b is used for square-shaped tables; tau-c is used for rectangular tables |
| | $r_{pb}$ | The point-biserial correlation coefficient: used when one variable (the predictor) is measured on a binary scale and the other variable is continuous |
| | $\varphi$ | The phi coefficient: used when variables and effects can be arranged in a $2 \times 2$ contingency table |
| | $C$ | Pearson's contingency coefficient: used when variables and effects can be arranged in a contingency table of any size |
| | $V$ | Cramér's V: like C, V is an adjusted version of phi that can be used for tables of any size |
| | $\lambda$ | Goodman and Kruskal's lambda: used when both variables are measured on nominal (or categorical) scales |
| | **(b) Proportion of variance indexes** | |
| | $r^2$ | The coefficient of determination: used in bivariate regression analysis |
| | $R^2$ | R squared, or the (uncorrected) coefficient of multiple determination: commonly used in multiple regression analysis |
| | $_{adj}R^2$ | Adjusted R squared, or the coefficient of multiple determination adjusted for sample size and the number of predictor variables |
| | $f$ | Cohen's f: quantifies the dispersion of means in three or more groups; commonly used in ANOVA |
| | $f^2$ | Cohen's f squared: an alternative to $R^2$ in multiple regression analysis and $\Delta R^2$ in hierarchical regression analysis |
| | $\eta^2$ | Eta squared or the (uncorrected) correlation ratio: commonly used in ANOVA |
| | $\varepsilon^2$ | Epsilon squared: an unbiased alternative to $\eta^2$ |
| | $\omega^2$ | Omega squared: an unbiased alternative to $\eta^2$ |
| | $R^2_C$ | The squared canonical correlation coefficient: used for canonical correlation analysis |

(*cont.*)

# General logic behind ES

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

$$r = \frac{Covariance\ (x,y)}{S.D.(x)S.D.(y)}$$

$$r = \frac{n\sum xy - (\sum x)(\sum y)}{\sqrt{[n\sum x^2 - (\sum x)^2][n\sum y^2 - (\sum y)^2]}}$$

$$d = \frac{M_1 - M_2}{pooled\ SD}$$

➤ Effect sizes go up when "**signal**" (*numerator*) increases relative to "**noise**" (*denominator*)

# Effect size: useful tools

Read these

**Evaluating Effect Size in Psychological Research: Sense and Nonsense**

**David C. Funder and Daniel J. Ozer**, University of California, Riverside

Use this: https://www.psychometrica.de/effect_size.html (give a look also here http://www.stat-help.com/spreadsheets.html)

Check (or ask) your analysis output (SPSS, R) for effect sizes

Effect size can be calculated starting from different bits of information and can be transformed (e.g., from r to d)

# Some other readings and tools

Some bibliographic 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.

➢ Ellis (2010). *The essential guide to effect sizes.* Cambridge University Press.

➢ Cohen (1992). A power primer. *Psychological Bulletin, 112,* 155-159.

➢ Cohen (1994). The earth is round (ρ < .05). *American Psychologist, 49,* 997-1003.

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

Some online calculators

➢ https://www.psychometrica.de/effect_size.html

➢ https://sites.google.com/site/lakens2/effect-sizes

➢ https://www.campbellcollaboration.org/this-is-a-web-based-effect-size-calculator/explore/this-is-a-web-based-effect-size-calculator

➢ http://www.stat-help.com/spreadsheets.html

# *Hypothesis testing II*

# Testing for Significance

In practice, each test of significance requires:

- We should know what is the **null hypothesis**

- The variability of the parameter we are testing: standard error

- The conventional cut-off is 0.05 (or 0.01)

And software produces for us:

- The calculation of the test

- The probability of our result in the population described by the null hypothesis (**p**)

# Significance of Correlation Coefficient

- Does the correlation coefficient differ from zero?

Correlation Matrix

| | | v1 | v19 |
|---|---|---|---|
| v1 | Pearson's r | — | |
| | p-value | — | |
| | 95% CI Upper | — | |
| | 95% CI Lower | — | |
| | N | — | |
| v19 | Pearson's r | 0.370*** | — |
| | p-value | <.001 | — |
| | 95% CI Upper | 0.459 | — |
| | 95% CI Lower | 0.273 | — |
| | N | 334 | — |

*Note.* * p < .05, ** p < .01, *** p < .001

Because p<.05, we can reject the null hypothesis that r is equal to zero

# One-tailed vs two-tailed t-test

- One-tailed test: we test that r=0 against r>0

- Two-tailed test: we test that r=0 against r≠0

# Hypothesis testing among groups

- Theoretical hypothesis: Frustration increases aggression

- Empirical test: we need a situation that produces frustration (**independent variable, IV**) and a variable that measures aggression (**dependent variable, DV**)

- **IV**: Experimental manipulation. Ss are given a memory test and told that they have done poorly (**frustration**) or well (**control**)

- **DV**: Ss are asked to evaluate the confederate for suitability to an assistant research post

# Defining and testing hypotheses

- **H1**: Frustration increases aggression

- **H0**: Frustration does not affect aggression

- To test it, we analyze sources of variation in the scores (DV)

- Two sources of variation

a) due to the experimental effect (systematic variance, effect variance, variation **between** groups)

b) due to other sources (non-systematic random variance, error variance, variation **within** groups)

- We need a test statistics with a known probability distribution

$$test\ statistic = \frac{variance\ explained\ by\ the\ effect}{error\ variance} = \frac{variance\ between\ groups}{variance\ within\ groups}$$

# The two sources of variance



| Group A | Group B |
|---------|---------|
| 10 | 6 |
| 8 | 5 |
| 7 | 9 |
| 6 | 4 |
| 11 | 5 |
| 9 | 7 |
| 6 | 3 |
| 4 | 8 |
| 10 | 6 |
| **7.89** | **5.89** |

# Partitioning of variance



Between groups

Total variance

**=**

Within groups

$$test\ statistic = \frac{variance\ explained\ by\ the\ effect}{error\ variance} = \frac{variance\ between\ groups}{variance\ within\ groups}$$

The greater this value, the better; its significance is evaluated against a probability distribution of the the test statistic considering a certain threshold (e.g., p=.05)

# *Errors of statistical inference*

# Errors of inference

- Frequentist approach

- There are three types of errors

- NHST*: Type I error (False positives)
  Type II error (False negatives)

- CI (aka "The New Statistics":
  Estimate error (imprecision)

NHST= Null Hypothesis Significance Testing

(what you have been 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) |
|---|---|---|
| **Conclusion of the significance test (SAMPLE)** — Null is false | Type I error ($\alpha$) | Correct decision ($1-\beta$) |
| Null is true | Correct decision ($1-\alpha$) | Type II error ($\beta$) |

# Errors of inference in NHST

- **Type I error:** *Erroneously rejecting the null hypothesis (<u>F</u>alse <u>positive</u>).* 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>F</u>alse <u>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.

# How to control Type I errors?

- 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 (*p value threshold*).

- Conventionally, researchers use an alpha rate ($\alpha$) 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)

UNIVERSITA' DEGLI STUDI DI MILANO BICOCCA

# How to control Type II errors?

- 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 Analysis*

# What is power?



critical t = 2.03452

1-α

β

1-β
(Power)

α

POPULATION

| | Null is true (H0 is correct) | Null is false (H1 is correct) |
|---|---|---|
| Null is true | Correct decision (1-α) | Type II error (β) |
| Null is false | Type I error (α) | Correct decision (1-β) |

SAMPLE

# The key determinants of power

- Power is determined by four elements

1) Decision criterion ($\alpha$)
2) Sample size ($n$)
3) Effect size ($\delta$)
4) Desired power ($1 - \beta$)

- Fixing one of the elements one can derive the others

# A simple example

- Fix $\alpha=.05$ and $(1-\beta)=.80$

- Plot sample size and effect size for a two sample t-test

# What affects power?

- 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

$$\eta^2 = \frac{SS_{Effect}}{SS_T}$$

$$d = \frac{M_1 - M_2}{pooled\ SD}$$

$$r(v,x) = \frac{cov(v,x)}{sd(v) * sd(x)}$$

# Power as a function of ES and N

t tests – Means: Difference between two independent means (two groups)
Tail(s) = One, α err prob = 0.05, Allocation ratio N2/N1 = 1

Effect size d

- ▷ = 0.7
- ◁ = 0.6
- □ = 0.5
- ◇ = 0.4
- ○ = 0.3

Power (1−β err prob)

Total sample size

# How to increase power?

## Power is affected by

- *Sample size*

- *Construct-related (i.e., SIGNAL) variance*

- *Construct-unrelated (i.e., NOISE) variance*

# What is affected by power?

## Higher power means

- *Less False Negatives*
- *Lower overall errors of inference (crucial error rates)*

## Lower power means

- *with multiple outcomes and HARKing: body of conflicting evidence in the literature*
- *with publication bias: presence of many false-positives in the literature*

# Is low power a real problem for Psychology?

## Researchers' Intuitions About Power in Psychological Research

Marjan Bakker[1], Chris H. J. Hartgerink[1], Jelte M. Wicherts[1], and Han L. J. van der Maas[2]
[1]Department of Methodology and Statistics, Tilburg School of Social and Behavioral Sciences, Tilburg University, and [2]Department of Psychology, Psychological Methods, University of Amsterdam

1990; Maxwell, 2004). Specifically, given the typical effect sizes (ESs) and sample sizes reported in the psychological literature, the statistical power of a typical two-group between-subjects design has been estimated to be less than .50 (Cohen, 1990) or even .35 (Bakker et al., 2012). These low power estimates appear to con-

**Yes!**

META-RESEARCH ARTICLE

## Empirical assessment of published effect sizes and power in the recent cognitive neuroscience and psychology literature

Denes Szucs[1]*, John P. A. Ioannidis[2]

1 Department of Psychology, University of Cambridge, Cambridge, United Kingdom, 2 Meta-Research Innovation Center at Stanford (METRICS) and Department of Medicine, Department of Health Research and Policy, and Department of Statistics, Stanford University, Stanford, California, United States of America

PLOS | BIOLOGY

We have empirically assessed the distribution of published effect sizes and estimated power by analyzing 26,841 statistical records from 3,801 cognitive neuroscience and psychology papers published recently. The reported median effect size was D = 0.93 (interquartile range: 0.64–1.46) for nominally statistically significant results and D = 0.24 (0.11–0.42) for nonsignificant results. Median power to detect small, medium, and large effects was 0.12, 0.44, and 0.73, reflecting no improvement through the past half-century. This is so because sample sizes have remained small. Assuming similar true effect sizes in both disci-

# Publication bias and low power



PB

0.30 0.60

The ES will be overestimated. How much depends on the extent of PB and on the prevalence of small samples.

A reader will think that Cohen's d=0.60 but in fact is d=0.30

# Publication bias, Effect Sizes, underpowered studies

**ES:** Cohen's d=0.60 (vs. d= 0.30)

**N for power:**

**80%**                    **90%**

72 Ss (vs. 278)        98 Ss (vs. 382)

Suppose we run a study with 98 Ss.
Expected power is 0.90 but **real power will be 0.43**

Vicious cycle: PB leads to overestimated ES leading to underpowered studies leading to non replicated effects, **even assuming that the effects are true and the researchers do not "cheat"**

# Why power analysis to plan studies?

- Without logistical constraints (infinite resources and no costs), only accuracy in estimating parameters should matter (e.g., AIPE, Maxwell, 2008 Ann Rew Psych)

- In an accuracy (precision) 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)

- Increased accuracy means less inference errors (both Type I and Type II)

- *If you want to get it right, increase sample size*

# Precision vs. Power

- They have different aims

INTRODUCTION TO THE NEW STATISTICS
ESTIMATION, OPEN SCIENCE, & BEYOND

GEOFF CUMMING &
ROBERT CALIN-JAGEMAN

Study 1

Study 2

Study 3

Study 4

-1        -0.5        0        0.5        1

**Value of Population Beta Weight of Interest**

*Figure 1.* Illustration of possible scenarios in which planned sample size was considered a "success" or "failure" according to the accuracy in parameter estimation and the power analysis frameworks. Parentheses are used to indicate the width of the confidence interval.

- Precision is valuable no matter everything else

UNIVERSITÀ DEGLI STUDI DI MILANO BICOCCA

# a MINOR practical problem…

- **Big** sample sizes are needed for precise estimates no matter the effect size



AIPE FOR THE STANDARDIZED MEAN DIFFERENCE

# How to calculate power

- Different software and routines (e.g., in R)
- A free comprehensive package is G*Power

**http://www.gpower.hhu.de/**

## G*Power: Statistical Power Analyses for Windows and Mac

G*Power is a tool to compute statistical power analyses for many different $t$ tests, $F$ tests, $\chi2$ tests, $z$ tests and some exact tests. G*Power can also be used to compute effect sizes and to display graphically the results of power analyses.

Kontrast

Suchbegriff

### Screenshots (click to enlarge)

Main Window

Main Window (Table)

Power Plot

Power Plot (Table)

# Power analysis calculations

- Examples of calculation of power analysis for some simple designs

- Also based on

INTERNATIONAL REVIEW
OF SOCIAL PSYCHOLOGY

RESEARCH ARTICLE

A Practical Primer To Power Analysis for Simple Experimental Designs

Marco Perugini, Marcello Gallucci and Giulio Costantini

SM and routines available at https://github.com/mcfanda/primerPowerIRSP

# G*Power



G*Power 3.1.9.2

File Edit View Tests Calculator Help

Central and noncentral distributions | Protocol of power analyses

**Test family**
t tests

**Statistical test**
Correlation: Point biserial model

**Type of power analysis**
A priori: Compute required sample size – given α, power, and effect size

**Input Parameters**

| | Tail(s) | One |
| Determine => | Effect size |ρ| | 0.3 |
| | α err prob | 0.05 |
| | Power (1–β err prob) | 0.95 |

**Output Parameters**

| Noncentrality parameter δ | ? |
| Critical t | ? |
| Df | ? |
| Total sample size | ? |
| Actual power | ? |

X–Y plot for a range of values    Calculate

**Analysis**

**Type of power analysis**

**Inputs**

**Action buttons**

**Output**

119

# Example: Two independent groups



Standard pre-study planning approach

Fix ES, α, 1-β

Calculate needed N

Good practice

We consider also sensitivity analysis

# Sensitivity analysis: Starting from N

"Sometimes" resources are fixed

You know that you can collect a certain N

The question becomes what ES can be found with sufficient power

Sensitivity analysis

# Sensitivity plot: N by ES

Figure 2: Sensitivity Plot of G*Power calculating the power of a two independent samples t-test: Lowest detectable effect size as a function of required N.

# Sensitivity plot: N by Power



Figure 3: Sensitivity Plot of G*Power calculating the power of a two independent samples t-test: Power as a function of required N for fixed effect size.

# Inspecting scenarios around effect sizes



t tests – Means: Difference between two independent means (two groups)
Tail(s) = One, Allocation ratio N2/N1 = 1, α err prob = 0.05

# Inspecting scenarios around N

t tests – Means: Difference between two independent means (two groups)
Tail(s) = Two, α err prob = 0.05, Allocation ratio N2/N1 = 1

Total sample size
= 180
= 190
= 200
= 210
= 220

**Plot Parameters**

| Plot (on y axis) | Effect size d | | |
| as a function of | Power (1 – β err prob) | from | 0.6 | in steps of | 0.01 | through to | 0.95 |
| | | | | | | and displaying the values in the plot | ✓ with markers |
| Plot | 5 | graph(s) | interpolating points | | |
| | with | Total sample size | from | 180 | in steps of | 10 |
| | and | α err prob | at | 0.05 | | 0.95 |

Draw plot

# Two repeated measures

**Analysis**

**Type of power analysis**

**Inputs**

**Output**

Test family
t tests

Statistical test
Means: Difference between two dependent means (matched pairs)

Type of power analysis
A priori: Compute required sample size – given α, power, and effect size

Input Parameters

Tail(s) | One

Determine => | Effect size dz | 0.5
α err prob | 0.05
Power (1–β err prob) | 0.80

Output Parameters

Noncentrality parameter δ | 2.5980762
Critical t | 1.7056179
Df | 26
Total sample size | 27
Actual power | 0.8118316

X–Y plot for a range of values | Calculate

# Two repeated measures

The ES for a paired means design is $d_z = \Delta / sd.$

This is not the same as Cohen's d, but calculated with the difference score divided by its SD. To double check with available previous results $d_z = \frac{t}{\sqrt{N}}$

If no previous results, could guess ES as if the two are independent groups and how much the measures are correlated (r/ρ)

$d_z = \frac{d}{\sqrt{2(1-r)}}$, e.g., with d=0.5 and r=.55, $d_z = \frac{0.5}{\sqrt{2(1-0.55)}} = 0.527$

and the other way round

$d = d_z * \sqrt{2(1-\rho)}$, e.g. with $d_z$=0.527 and r=.55,

$d = 0.527 * \sqrt{2(1-.55)} = 0.527 * 0.95 = 0.50$

d = 0.50

# Let's see it together !

1) We want to run a study to replicate a previous research.

   They had two groups and found these results:

   N1=40, M=2.53, SD=0.34

   N2=40, M=2.88, SD=0.42.

   We want to estimate N with power at 80% and at 95%

   ($\alpha$=.05, one-tail, two-tails). **N?**

2) We expect a subtle but theoretically important effect

   ($d$=0.15). How many N with power at 90% or 80%

   ($\alpha$=.05, one-tail, 2 groups). **N?**

3) We can run a study with about 120 Ss. What ES can we

   detect at power 80% ($\alpha$=.05, two-tails) for two-groups

   between design? **ES?**

# Let's see it together !

# Let's see it together !