



<https://www.vox.com/future-perfect/21504366/science-replication-crisis-peer-review-statistics>fbclid=IwAR3lIJXfXBVwFWaE5aw4RXHKY

NHST, CI, EFFECT SIZES AND STATISTICAL POWER

Ph.D Programme in Psychology, Linguistics and Cognitive
Neurosciences

QUICK REFRESH ON GLIMs

- From the General Linear Model to the Generalised Linear Model
- Results from first application
- Bridge from GLIM to GLIMM
- Which assumption of the GLIM are violated when GLIMM is needed?
- Are there any ways for checking the need for a GLIMM instead of a GLIM?

PLAN OF THE LESSON

- Part I
 - Icebreakers: NHST and p -values
- Part II
 - Effect sizes, P-values, & Power
 - The language of power analysis
 - Types of power analysis
- Part III
 - Conducting, running and reporting a power analysis
 - Software for Power Analysis:
 - G*Power
 - R (introduction)

PART I

Icebreakers: NHST and p -values

A SHORT PREMISE: SOME HISTORY

- Ronald Fisher in the 20's, described the testing of a null hypothesis and used p values to this purpose. He did not set the 0.05, 0.01 criteria
- Neyman and Pearson, in the 30's, described the extension of Fisher's model by adding the notion of the alternative or research hypothesis
- The use of the p -value as compared with standards of 0.05 and 0.01 followed soon after that
- In early 90's there were suggestions to include effect sizes in reporting (APA Style Manual, 1994).
- In 1999, the Task Force on Statistical Inference (TFSI) was formed to report on the controversy about significance testing and to promote the use of alternative methods The TFSI also described the different effect sizes that could be used.
- Many of the methods for effect size were introduced a long time ago, pre-Fisher, as r and η
- Fisher described eta-squared or the correlation ratio as a measure of variance accounted for
- Cohen's d about 1962 Glass effect size about 1976 Hedges effect size about 1981.

NULL HYPOTHESIS SIGNIFICANCE TESTING (NHST)

- Critical value and corresponding p level of significance (**criterion of significance**):
 - p value the proportion of null experiments, out of all possible experiments, that will turn out significant even when the null hypothesis is true (α usually 05 , more seldom .01)
- It is related to sampling (**false positive**) [Analogously for **type II error, false negative**].
 - We sample (evidence) from all possible values in the population (the sample space). In the long run, sooner or later, we will sample from the 'extremes', which are eccentric and in this case the evidence not the general 'true' behavior.
- **p -value**: probability of observing test statistics greater than the (absolute value of) the critical value, under the null hypothesis (**false positive**)

*

CAN *P*-VALUES BE MEASLEADING? (I WANT TO SEE THE STARS!)

*

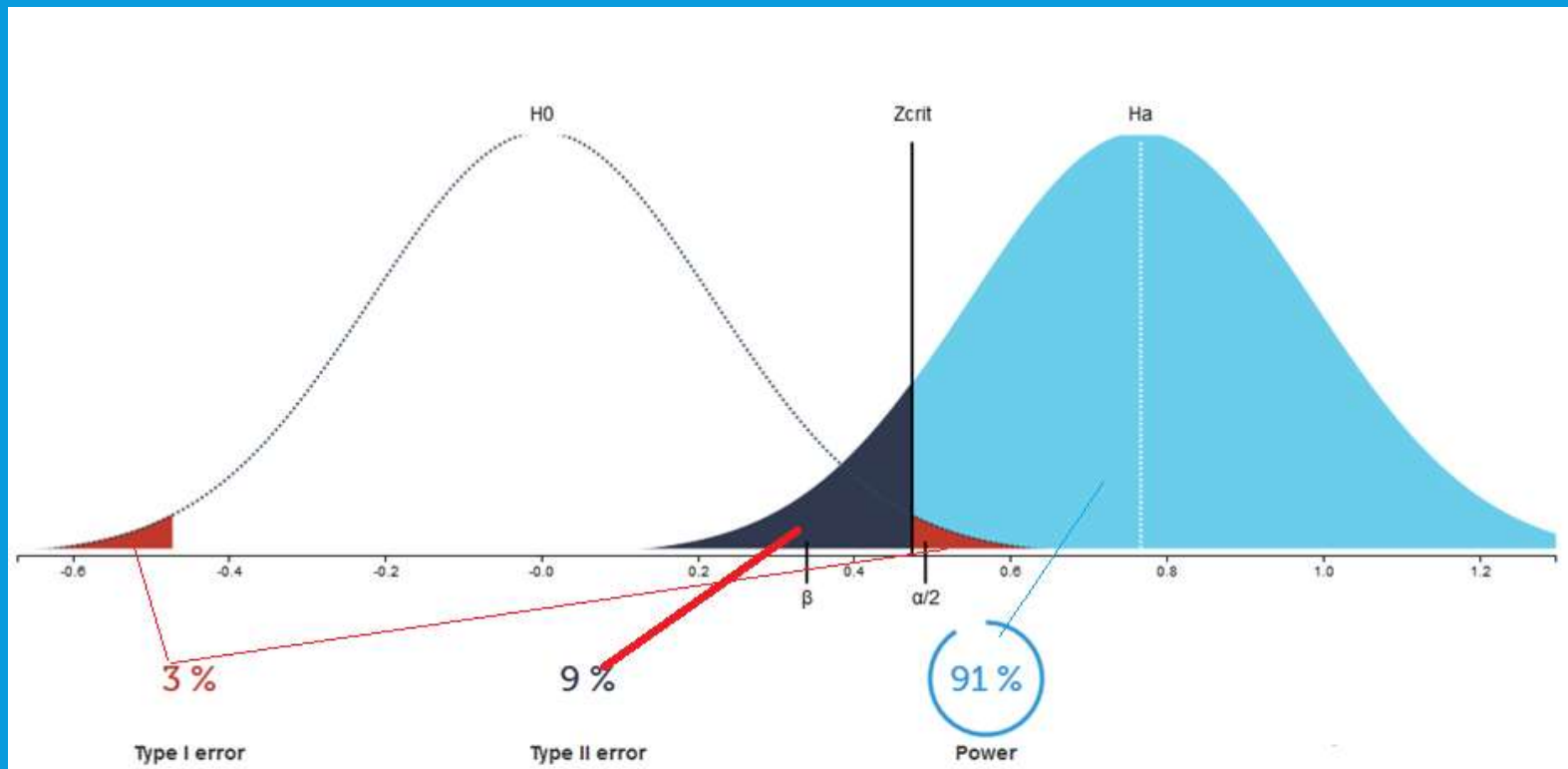
- Statistical significance provides no information about size of effects and other aspects. Many *p*-values are not very meaningful if not read together with the effect size and in a replicability conceptual framework.
- *p*-values can be so specific that they don't tell researchers what they need to know
 - Information on *p* cannot be used to make decisions about how to use the results
- *p*-values indicate only that a difference could be or should not be attributed to chance in extracting that specific sample from all possible sample (this holds under the NHST).

NHST: THE LOGIC REJECTING OR FAILING TO REJECT (FTR)

	Reality: NO EFFECT	Reality: EFFECT EXISTS
Research concludes: FAIL TO REJECT NULL; NO EFFECT	CORRECT FTR	TYPE 2 ERROR (β)
Researcher concludes: REJECT NULL; EFFECT EXISTS	TYPE 1 ERROR (α)	CORRECT REJECT ($1-\beta$)

Type I and II errors
Coeteris paribus, when
type I decreases, type II
increases

NHST: THE GRAPH





TYPE I ERROR RATE

- The first error rate, **the significance level**, is chosen by the experimenter and is conventionally one of 5%, 1% or 0.1% or ($P=0.05$) ($P=0.01$) or ($P=0.001$). Each error rate gives a **threshold value** (on the X axis) that must be exceeded for significance.
- Interpretation: in 5% of experiments in which there is no real treatment effect the t statistic will exceed the threshold value due to chance in sampling
- Choice of Type I error rate. The smaller the rate chosen the stronger the evidence that there is an effect of treatment.

The typical reference t-test for the difference in the means of two population

Example: difference in height between basket players and the general population, treatment vs control.



BEWARE

- nominal alpha: the probability of making a Type I error when all the assumptions are met
- real alpha: the probability of making a Type I error when one or more of the assumptions are violated
- Real alpha higher than nominal alpha: alpha inflation



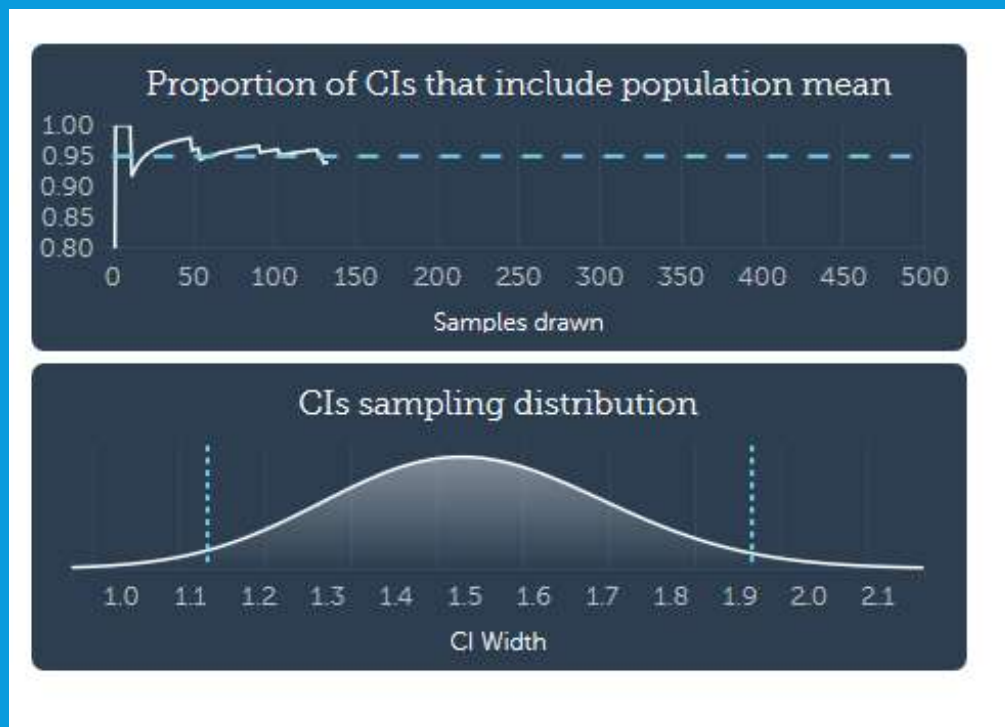
NHST VERSUS CI

t test
mean under the null hypothesis: 120 sentences

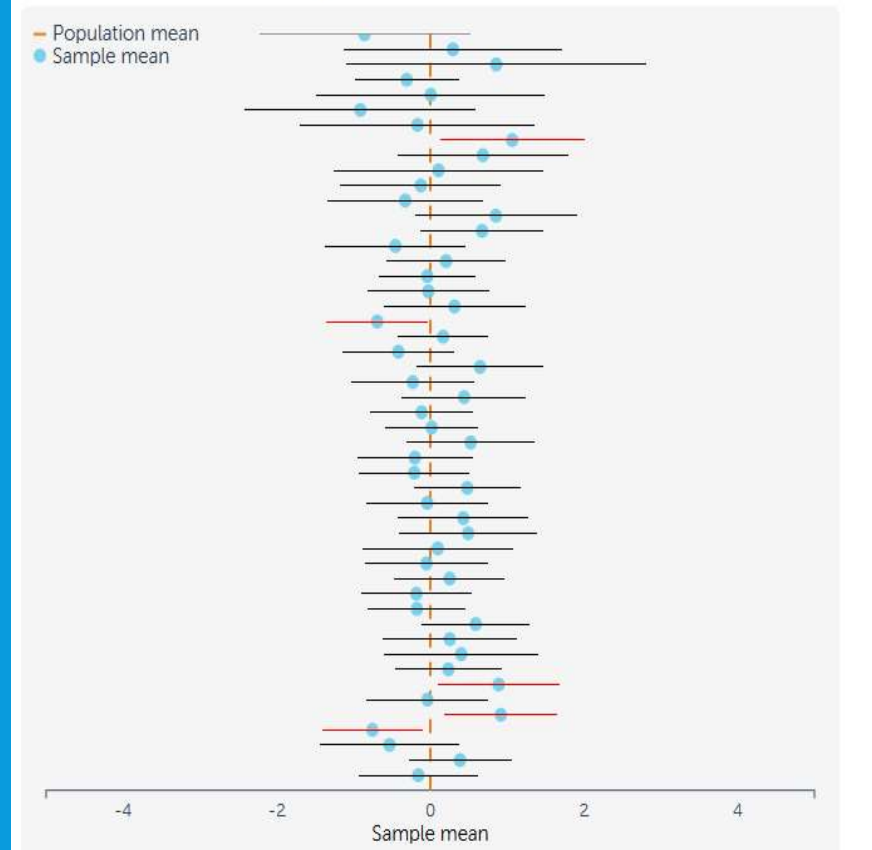
	t	df	Sign. (two tails)	mean diff	95% CI	
					lower	upper
tot_fras	8462673,000	101	,000	67049020,00	51332092,00	82765947,00

What is the main difference between NHST and CI?

- NHST: acceptance region is the bilateral CI under the null hypothesis, i.e. when the value of the null hypothesis is considered the TRUE mean in the population
- CI: no hypothesis or knowledge whatsoever about the value of the mean (or means difference) in the population
- Is the CI informative? When its precision is high, i.e. when the CI width is small: (upper limit-lower limit)



95% confidence intervals

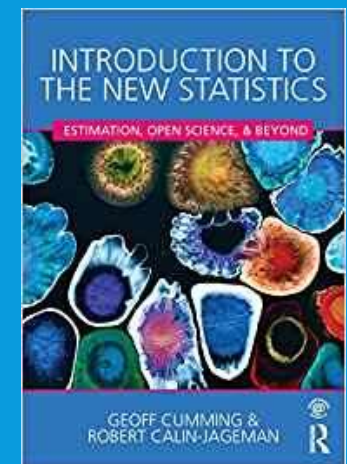
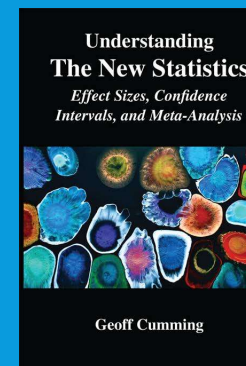


NEW VERSUS OLD (?)

The accent is on replicability. The 'novelty' lies in this perspective.

The need for a larger use of CIs has been underlined by APA and it does not imply NHST neglect.

Contents have been largely known and discussed (see article discussing NHST and CI, written in 1979 and later papers in *Psychological Science*)



PART II

Effect sizes, P-values, sample size and Power

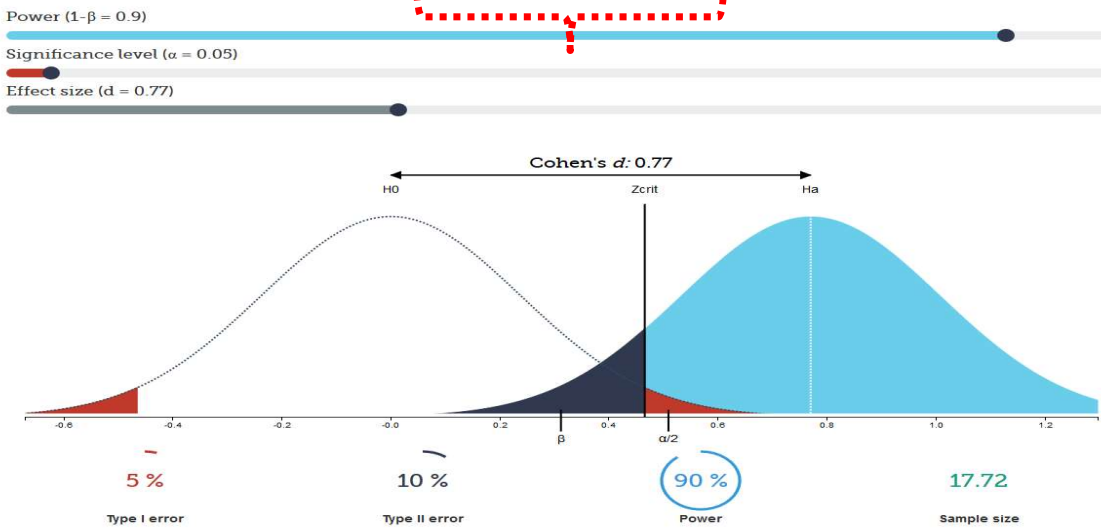
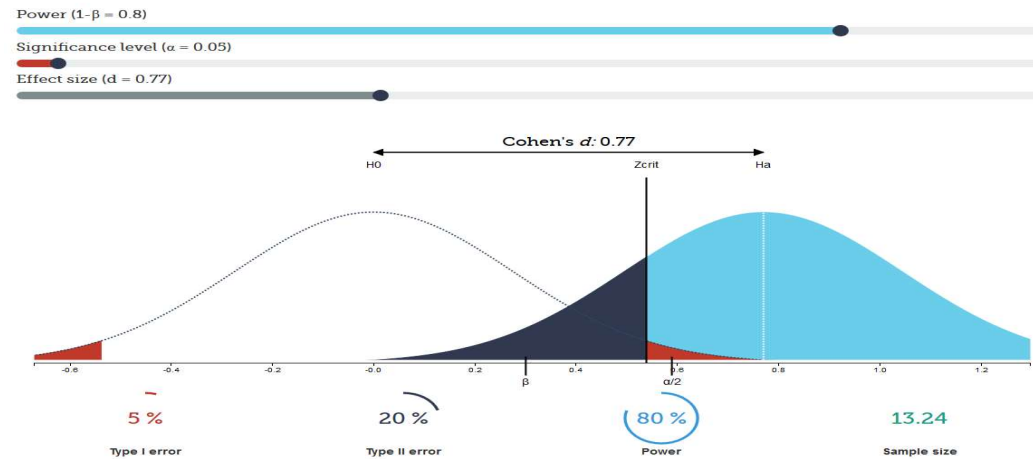
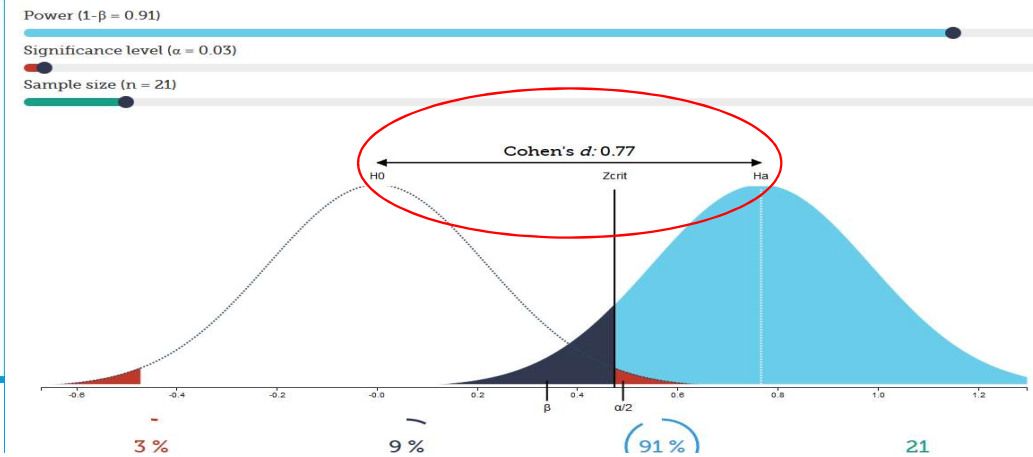
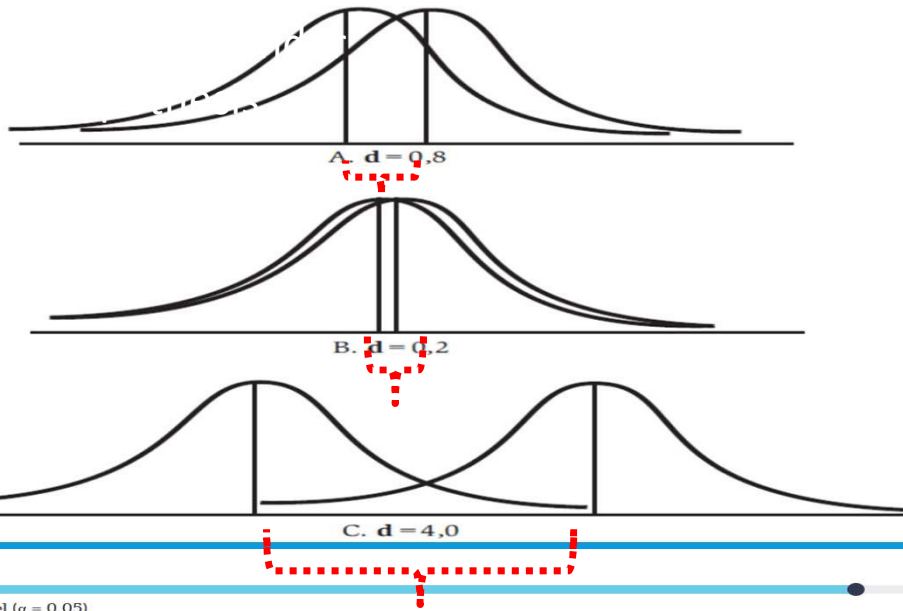
TYPE I ERROR RATE AND THE EFFECT SIZE

- In research, p values are one of the two measures used to report study results. The other one is the effect size (ES).
- An ES is what the result found, e.g. the difference found in the mean scores between two groups. It measures the strength of the result and it is pure, as it does not depend on sample size., unlike p -values.
- $ES = (\text{Meantreatment} - \text{Meancontrol}) / S_{\text{dpooled}}$ *It is a pure number!*
- A general definition of ES is that it is a family of indexes that measure the magnitude of a treatment effect
- What is the probability of detecting an effect when the effect exists in the population?
- Real effects may be very important or very unimportant. Is the ES generalizable?

EFFECT SIZE

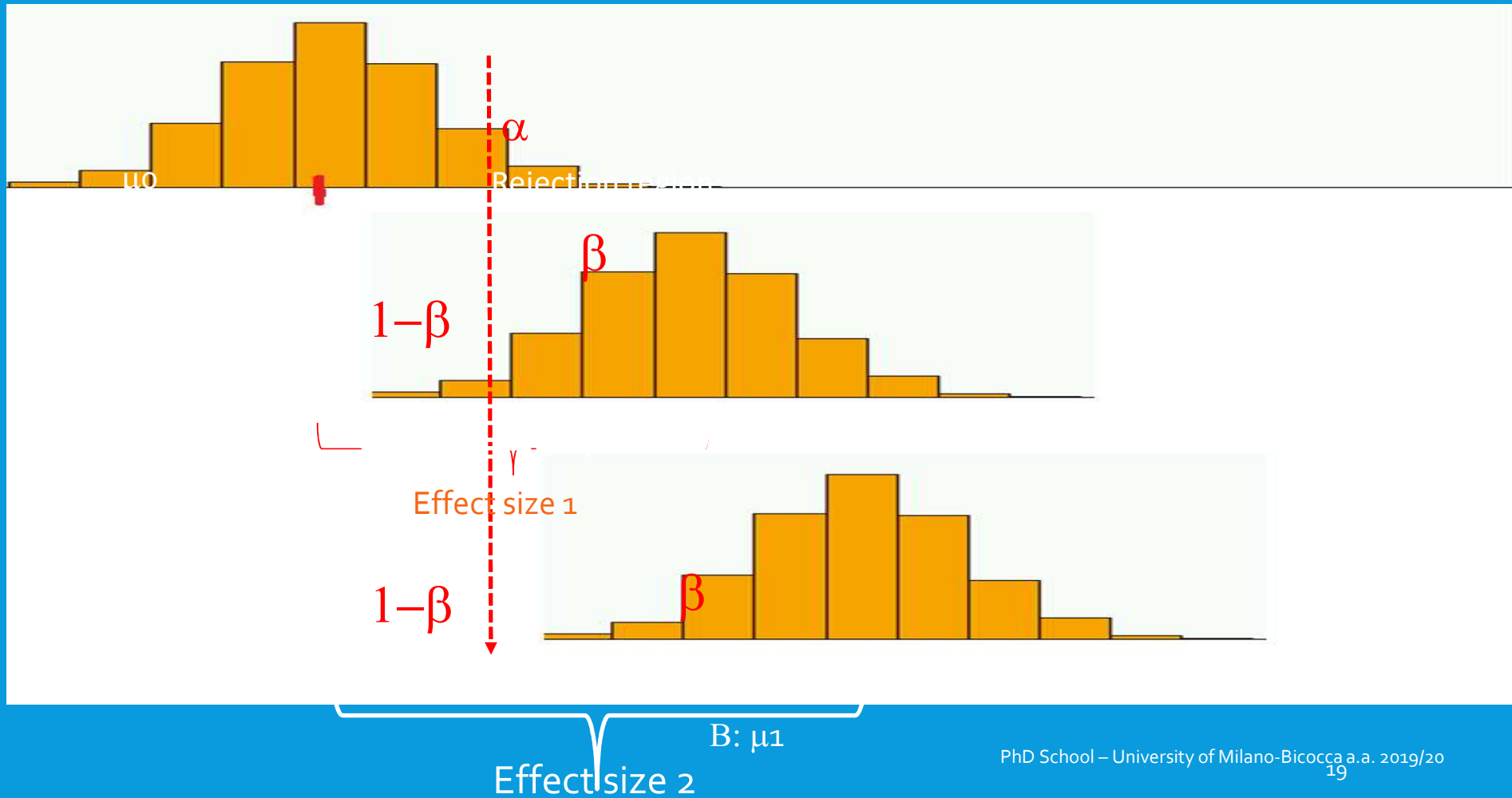
- Knowledge of the magnitude of a treatment effect is qualitatively different than knowing if the effect is real.
- Real ES may be very important or very unimportant (negligible).
- Two types of ES:
 1. standardized units of difference
 2. Variance-accounted-for statistics

UNDER THE NORMALITY ASSUMPTION: THE BIGGER D, THE BETTER



Bigger and bigger effect size d

A bigger effect size requires a smaller power to be detected



ES IN STANDARDIZED UNITS

An ES is the observed difference between means, proportions, etc.

To be useful for comparative purposes, this difference needs to be standardized.

Standardization relies on the *pooled variance*, i.e. on the variance computed on all subjects.

- Cohen's d (most common)
- Glass's g

EFFECT SIZES IN STANDARDIZED UNITS OF DIFFERENCE

- Observed difference on means, proportions, etc. The difference needs to be standardized
- Cohen's d (most common), Glass's g'

Cohen's d (1969)

Population Form

$$d = \frac{\mu_1 - \mu_2}{\sigma}$$

Statistic Form

$$d = \frac{\bar{X}_1 - \bar{X}_2}{s}$$

σ and s represent the total score standard deviation

First recognized effect size.
Mean differences in units of common population standard deviation (called pooled sd).
Cohen used this as the basis his research on power.

Glass's Effect Size (1978)

$$g' = \frac{\bar{X}_{Experimental} - \bar{X}_{Control}}{S_{Control}}$$

Glass proposed a modification of the Cohen d the common standard deviation replaced with the standard deviation of the control group.

VARIANCE-ACCOUNT-FOR STATISTICS

- Very similar to the correlation coefficient (r) and the coefficient of determination (r^2). These provide indications of the proportion of variance that can be attributed to the treatment.
- For mixed models, the intraclass correlation coefficient, estimates the effect size. It measures the proportions of total variance in the second (higher) level of the model.
- In case of nominal variables, Cramer's V is applied, as a transform of the Chi-square
 - Eta-squared, η^2
 - Intraclass Correlation, ICC
 - Cramer's V

EFFECT SIZES AS VARIANCE-ACCOUNTED-FOR STATISTICS

Model	Effect size
Regression, anova	<p>Eta-squared, η^2</p> <p>More often used in meta-analysis</p> <p>A η^2 of 0.25 would indicate that 25% of the total variation is accounted for by the treatment variation.</p>
Mixed models	<p>intraclass Correlation, ICC</p> <p>Interaction can be tricky: Results showed that power varies significantly as a function of model type and whether or not the model is the main model for the study</p>
chi-square applications	<p>Cramér's v where $df^* = \min(r - 1, c - 1)$ and $r =$ number of rows and $c =$ number of columns</p>

$$\eta^2 = \frac{SS_{Treatment}}{SS_{Total}}$$

$$\rho = \frac{\text{population variance between level-two units}}{\text{total variance}}$$

$$v = \sqrt{\frac{\chi^2}{n \cdot df^*}}$$

Mathieu, J

Relating η^2 to Effect Size d (Coehn, 1988)

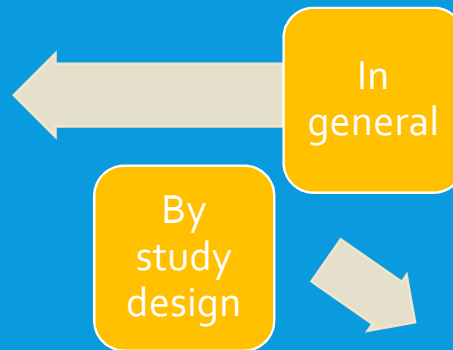
$$r = \frac{d}{\sqrt{d^2 + 4}}$$

EFFECT SIZES THRESHOLDS

Cohen's Standard	<i>d</i>	<i>r</i>	<i>r</i> ²
LARGE	0.8	.371	.138
	0.7	.330	.109
	0.6	.287	.083
MEDIUM	0.5	.243	.059
	0.4	.196	.038
	0.3	.148	.022
SMALL	0.2	.100	.010
	0.1	.050	.002
	0.0	.000	.000

Cramer's V

<i>df</i> *	<i>small</i>	<i>medium</i>	<i>large</i>
1	.10	.30	.50
2	.07	.21	.35
3	.06	.17	.29
4	.05	.15	.25
5	.04	.13	.22



	Effect Size Benchmarks		
Statistic	Small	Medium	Large
Means - Cohen's d	0.2	0.5	0.8
ANOVA - <i>f</i>	0.1	0.25	0.4
ANOVA - eta squared	0.01	0.06	0.14
Regression <i>f</i> -test	0.02	0.15	0.35
Correlation - <i>r</i> or point serial	0.1	0.3	0.5
Correlation - <i>r</i> squared	0.01	0.06	0.14
Association - 2 x 2 table -OR	1.5	3.5	9
Association - Chi-square - <i>w</i> or Phi	0.1	0.3	0.5

INTERPRETING THE EFFECT SIZES/ STANDARDIZED MEASURES

- A d or g of 1.2 indicates that the range of difference among the means is one and two-tenths of the size of the standard deviation.



Cohen chose three values that had been used extensively as standards for effect size.

Beware: Cohen warned about using these standards in practice. The major problem is that effect sizes are influenced by the number of samples and the sample sizes.

POWER

- High error rate II β means that a real difference between treatments is unlikely to be detected.
- A rate of 0.5 means that 50% of all possible experiments will not detect effects.
- Power = $1 - \beta$ = : The probability of rejecting the null hypothesis when it is false, i.e. to detect a real difference.
- In everyday language, power is the probability of concluding that the group means differ on the basis of your sample evidence, when the groups means actually differ in the population.

SAMPLE SIZE, TYPE I ERROR RATE, EFFECT SIZE AND POWER

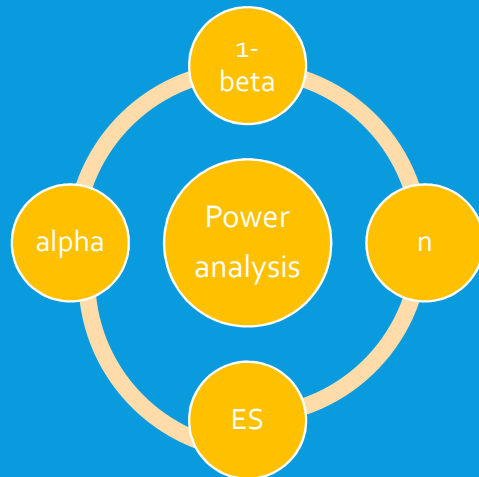
POWER ANALYSIS

- Type II error rate is related to:
 - sample size
 - Effect size
 - SD
 - Type I error rate α (the smaller type I error, the higher type II, *coeteris paribus*).
- These factors – ES, sample size, α , power, form a closed system. Once any three are established, the fourth is completely determined.
- What is a power analysis? A process by where one of several statistical parameters can be calculated given others.
- Usually, a power analysis calculates needed sample size given some expected effect size, alpha, and power.

A POWER ANALYSIS INVOLVES FOUR STATISTICAL MEASURES, 'FIXING' THREE OF THEM RESEARCHERS SOLVE FOR THE FOURTH

Probability of Type I error α

- Probability of finding significance where there is none
- False positive
- Usually set to .05



- n
- The sample size - usually the parameter you are solving for

Power $1-\beta$





- Probability of finding true significance
- True positive
- e beta is :
- Usually set to .80

ES


- The 'expected effect' is ascertained from:
- Pilot study results
- Published findings from a similar study or studies
- Sometimes calculated from results if not reported
- Field defined 'meaningful effect'
- knowledge of the field)


THE LANGUAGE OF POWER ANALYSIS

Factors Affecting Power

1. Size of the effect 
2. Standard deviation of the characteristic 
3. Bigger sample size 
4. Significance level desired 

Other factors in GLIMs

In a more complex model with more parameters or with more complex effects, power because SE gets larger 

Bigger error (unexplained variance and therefore smaller R^2 power) 
larger error means larger standard error of parameter estimates.

THE LANGUAGE OF POWER/2

Sample size in each group (assumes equal sized groups)

desired power (typically .84 for 80% power).

$$n = \frac{2\sigma^2 (Z_{\beta} + Z_{\alpha/2})^2}{\text{difference}^2}$$

Standard deviation of the outcome variable

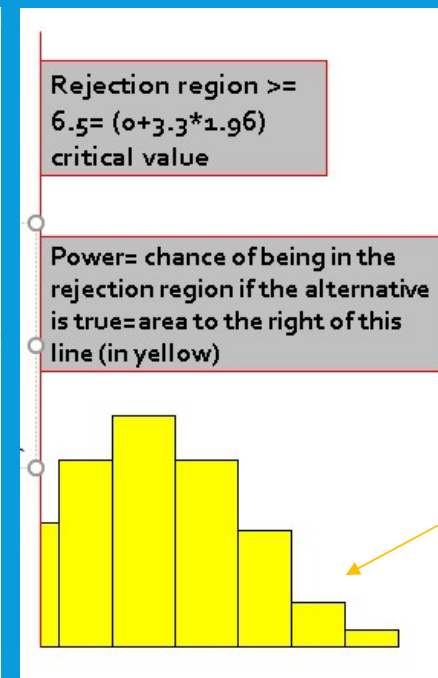
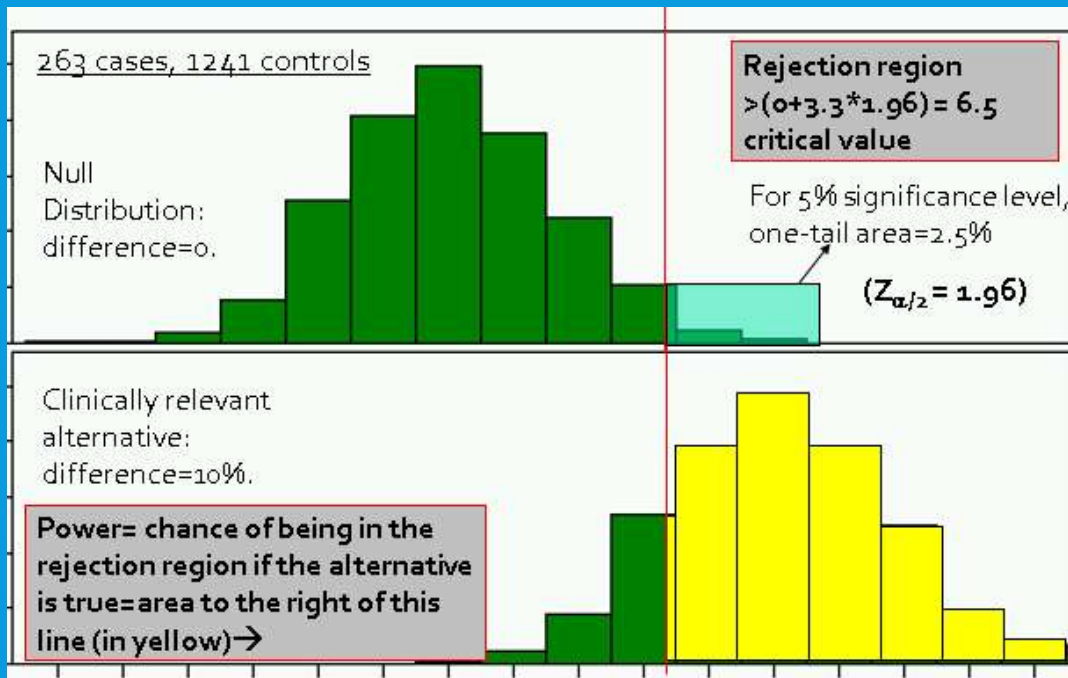
Effect Size (the difference in means)

desired level of statistical significance (typically 1.96).

All-purpose power formula

$$Z_{power} = \frac{\text{difference}}{\text{standard error(difference)}} - Z_{\alpha/2}$$

HOW α (ES, ETC.) DETERMINES POWER



10: mean under H1
3.3=SE

Power:

$$P(Z > (6.5 - 10) / 3.3) = P(Z > 1.06) = 85\%$$

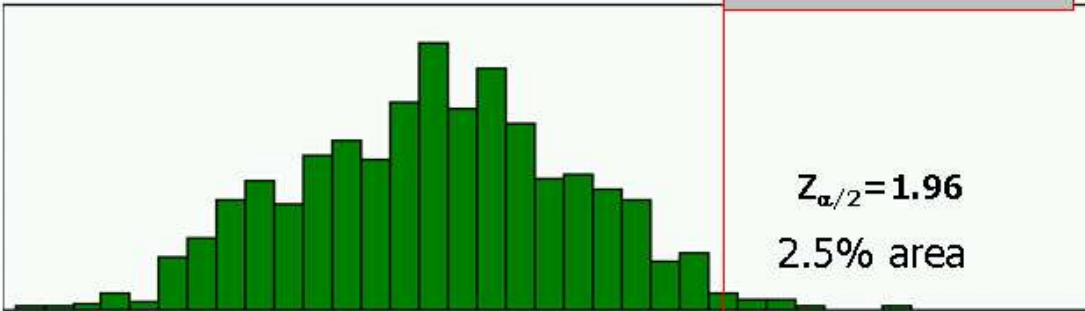
HOW ES AND VARIABILITY DETERMINES POWER

Big SE and small ES

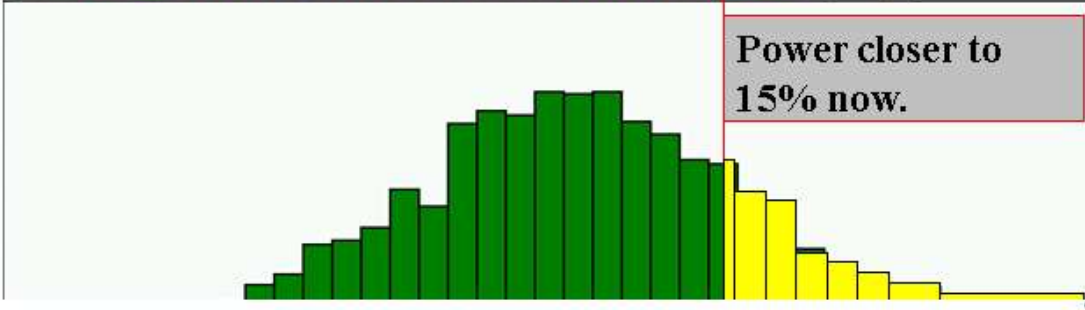
n=1504

Critical value=
 $0+10*1.96=20$

$Z_{\alpha/2}=1.96$
2.5% area



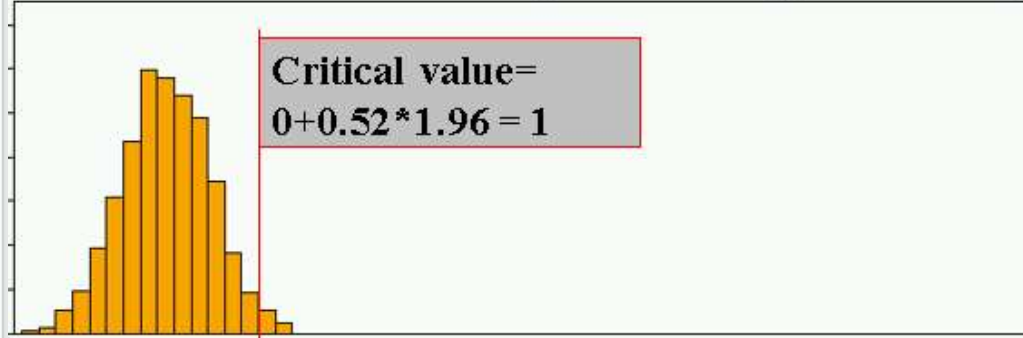
Power closer to
15% now.



Huge ES and small SE

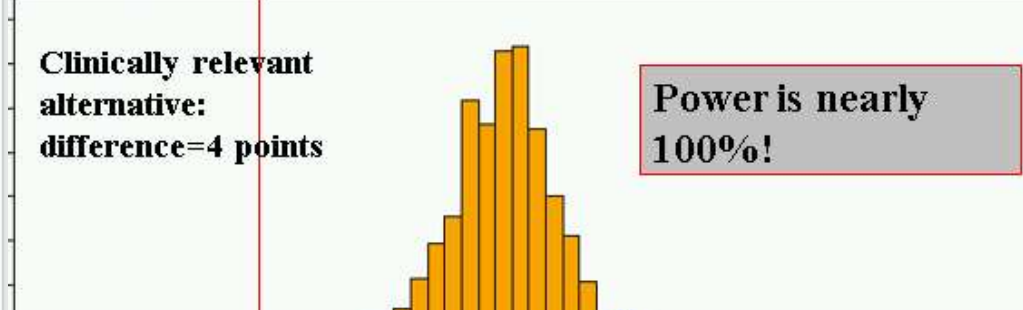
18 treated, 72 controls, sd= 2

Critical value=
 $0+0.52*1.96 = 1$

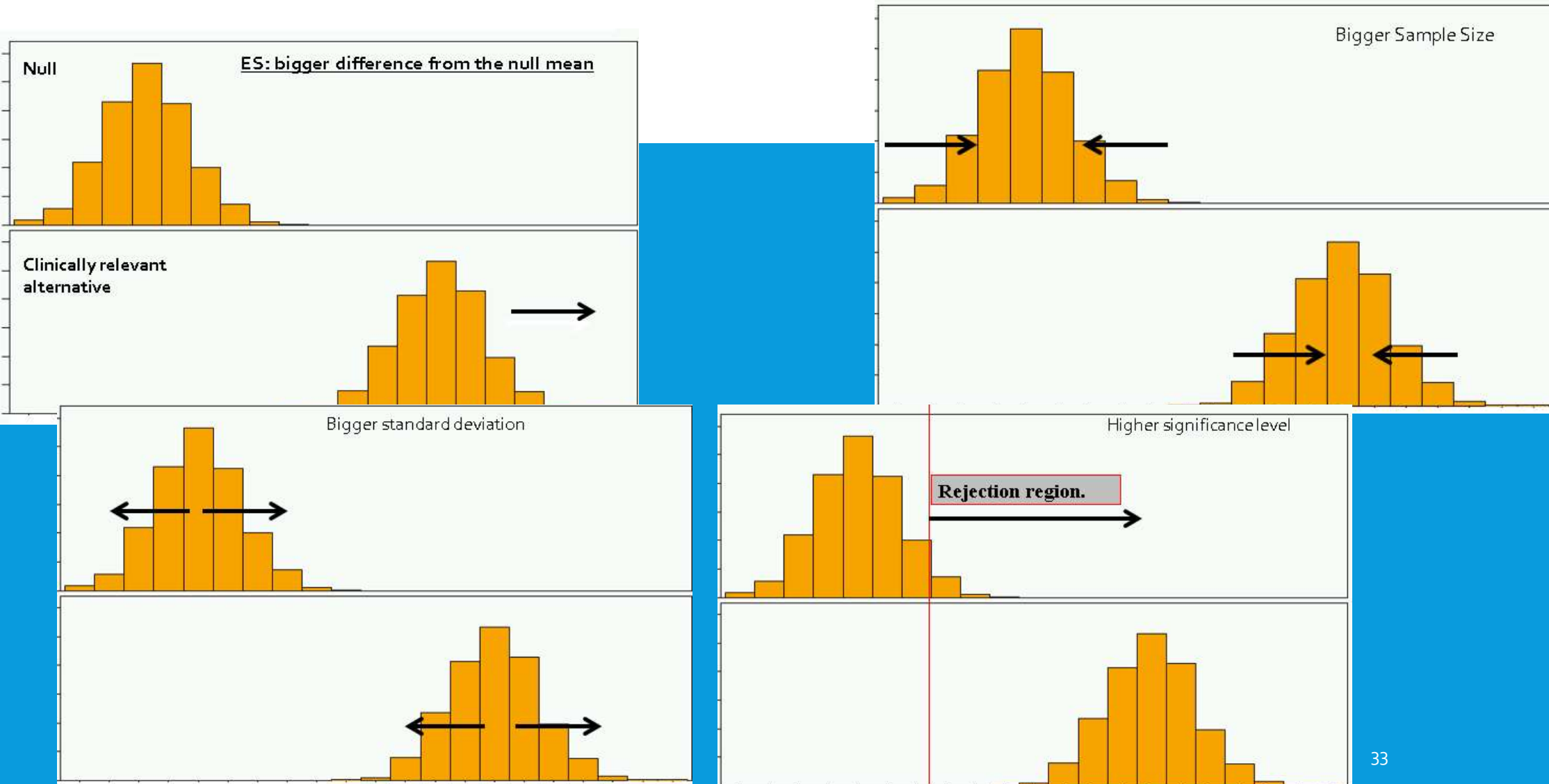


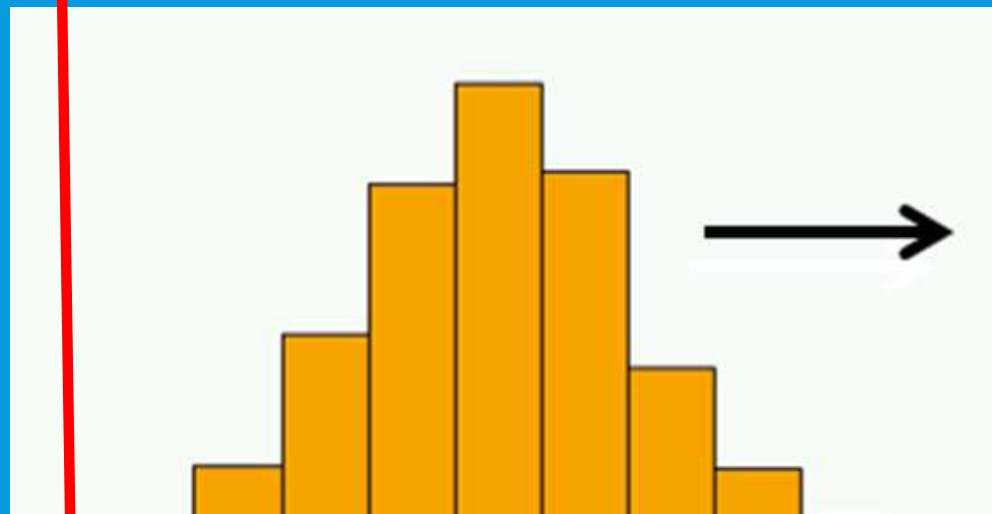
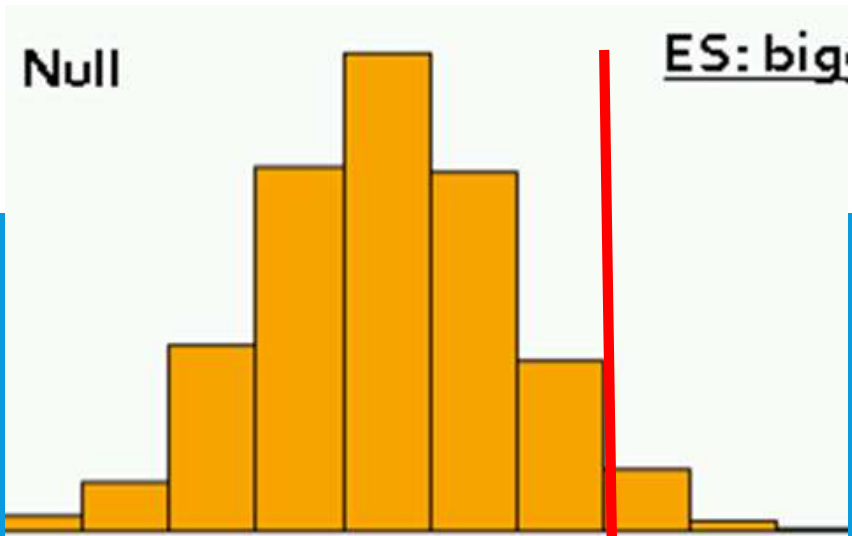
Clinically relevant
alternative:
difference=4 points

Power is nearly
100%!



A BRIEF GRAPHICAL OVERVIEW





TYPES OF POWER ANALYSIS

Power conventions

- Desired level of power: the more the better, value of .80 minimum threshold standard
- Higher power means more precision in estimating ES (tighter CI)

Strategies, determine:

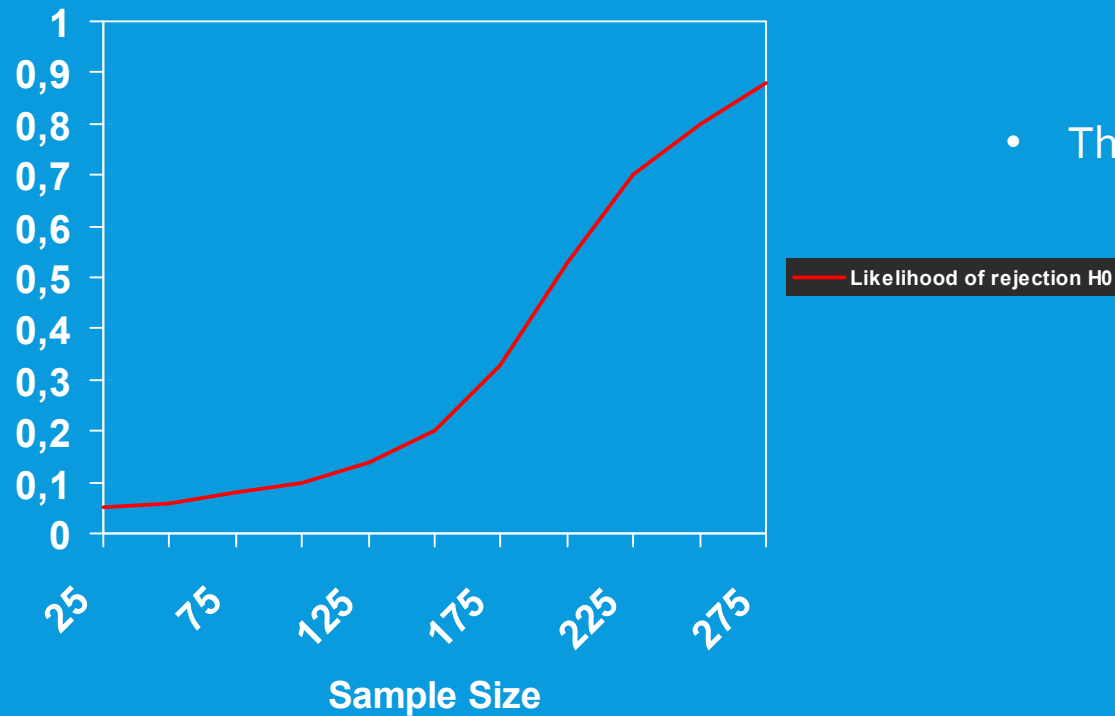
- number of subjects needed (n^*) for given level of power (e.g. .80)
- power for a given design (e.g., completed experiment with fixed n)

Types of power analysis

- A priori: compute n , given alpha, power, ES (advised)
- Post-hoc: compute power, given alpha, n , ES (controversial, in designing a study you need a priori)
 - Criterion: compute alpha, given power, ES, n (not much used)
- Sensitivity: compute ES, given alpha, power, n (useful for Minimal Detectable Effect MDE)

POWER FUNCTIONS

Power ($1-\beta$)



POWER FUNCTION AND F DISTRIBUTION

- F is a ratio of observed effect to error
 - $F = (\text{True Effect} + \text{Error}) / \text{Error}$
- The larger the true treatment effect, the larger F you expect to find
- If the null hypothesis is correct, $E(F) = 1.0$
- The power of most statistical tests in social sciences can be evaluated via the familiar F distribution (D. Lakens)

HOW TO INCREASE POWER

Increase n

- Effects of adding more subjects are not identical to those of adding more observations

Increase ES

- Choose a different research question
- Use stronger treatments or interventions
- Use better measures

Effects of implementing power analysis:

- Stronger studies: larger samples, better measures
- Fewer studies: adequate studies are harder to do than most people realize

CONDUCTING A POWER ANALYSIS

Are all tests the same in the face of power?

- Some statistical tests are more powerful (i.e., better at detecting real/non-zero population effects) than others.
- Parametric tests often are more powerful than non-parametric, because they work with more information from the data.
- GLIM is Minimum Variance (Estimator) when assumptions are met

Power can be calculated for tests of

- Effect for single regressor, subset of regressors controlling for other regressors in model, or all regressors in the model.

CONDUCTING A POWER ANALYSIS

A priori power analysis (sample size planning)

- set α level (max .05), and desired power (min .80)
- Specify (calculate) expected ES(conservative).
- n is a function of the above factors.

Software

For GLIMs G*Power <http://www.gpower.hhu.de/>

- Most popular by far, free download available for both the PC and Mac.
 - It includes an effect size calculator
 - On online tutorial manual

For (few) GLIMMs Optimal Design <http://hlmsoft.net/od/>

- It is somehow related to HLM, the free version is very basic

Commercial scientific software

- Mplus GLIMMs, Latent Models
 - HLM GLIMMs

All models open source R packages

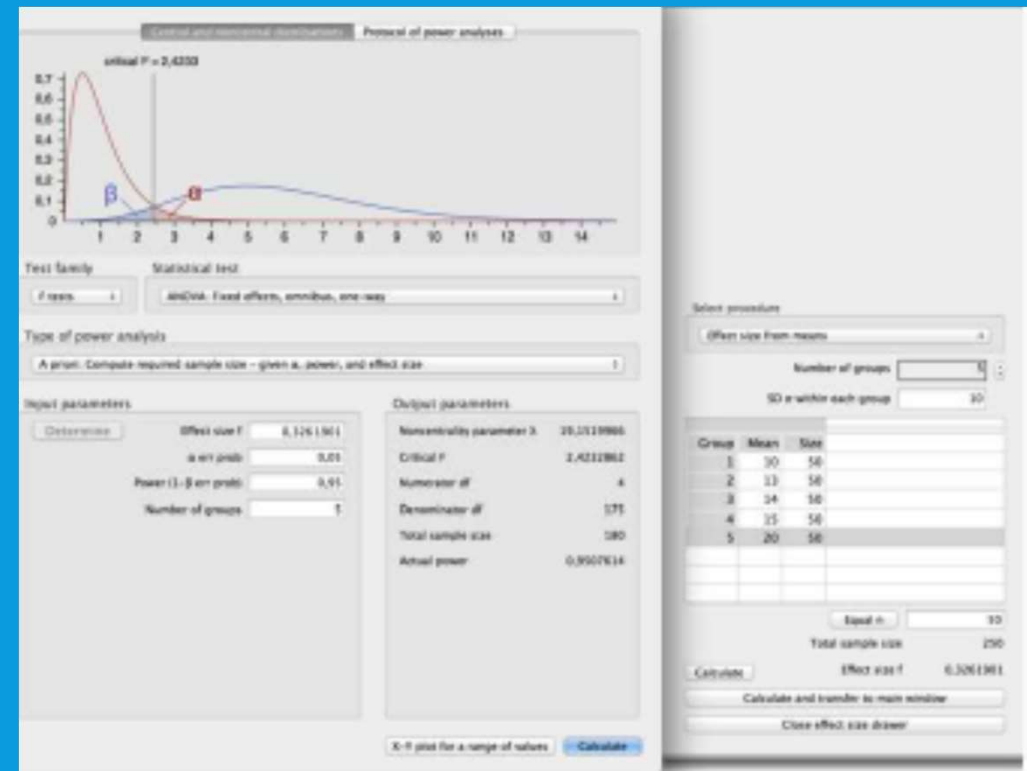
CONDUCTING A POWER ANALYSIS WITH G*POWER



The screenshot shows the UCLA idre website. The header includes the UCLA logo and the text "idre Institute for Digital Research and Education". Below the header is a navigation menu with "HOME", "SOFTWARE", and "RESOURCES". The "G*POWER" section is highlighted, with a sub-header "These pages were developed using G*Power version 3.0.10." and a list of statistical tests: Single-sample t-test, Paired-sample t-test, Independent-sample t-test, Two independent proportions, One-way ANOVA, and Multiple Regression. To the right, there is a video player with the title "Allgemeine Psychologie und Arbeitspsychologie" and a thumbnail image of a person at a computer.

Detailed illustrations can be found on the related UCLA Idre web page: <https://stats.idre.ucla.edu/other/gpower/>

The UCLA idre web page provides also an introductory seminar to power analysis <https://stats.idre.ucla.edu/other/mult-pkg/seminars/intro-power/>



The screenshot shows the G*Power software interface. The main window displays a graph of the power function for a one-way ANOVA. The x-axis represents the effect size f (ranging from 0 to 14), and the y-axis represents the power (ranging from 0 to 0.7). A red curve shows the power function, and a vertical line indicates the critical F value of 2.4233. The graph also shows the distribution of the test statistic under the null hypothesis (blue curve) and the distribution under the alternative hypothesis (red curve). The critical F value is marked on the x-axis, and the corresponding power is marked on the y-axis.

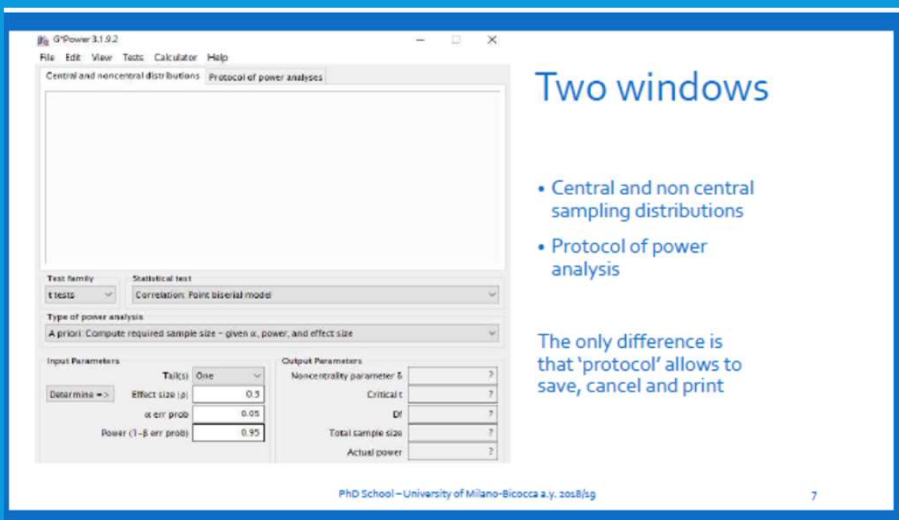
The interface includes several input fields and buttons:

- Test family:** Statistical test: ANOVA: Fixed effects, omnibus, one-way
- Type of power analysis:** A priori: Compute required sample size - given α , power, and effect size
- Input parameters:** Effect size f : 0.1251961, α err prob: 0.05, Power (1 - β err prob): 0.91, Number of groups: 5
- Output parameters:** Noncentrality parameter λ : 28.0219965, Critical F : 2.4232982, Numerator df : 4, Denominator df : 171, Total sample size: 180, Actual power: 0.9007614
- Select procedure:** Effect size from means, Number of groups: 5, SD σ within each group: 30
- Group Mean Size table:**

Group	Mean	Size
1	30	50
2	13	50
3	14	50
4	15	50
5	20	50

Additional controls include "Equal n" (set to 30), "Total sample size" (set to 180), "Calculate" button, "Effect size f " (set to 0.1251961), "Calculate and transfer to main window", and "Close effect size drawer".

CONDUCTING A POWER ANALYSIS WITH G*POWER



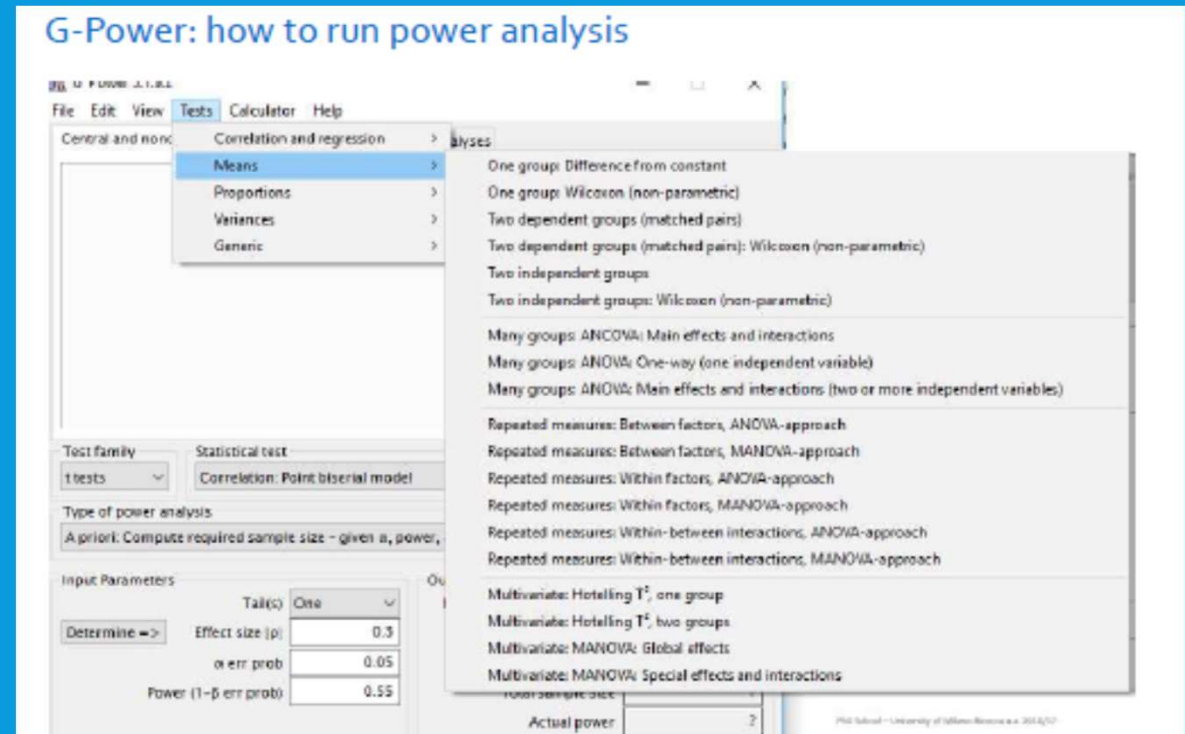
Two windows

- Central and non central sampling distributions
- Protocol of power analysis

The only difference is that 'protocol' allows to save, cancel and print

PHD School – University of Milano-Bicocca a.y. 2018/19 7

G-Power: how to run power analysis



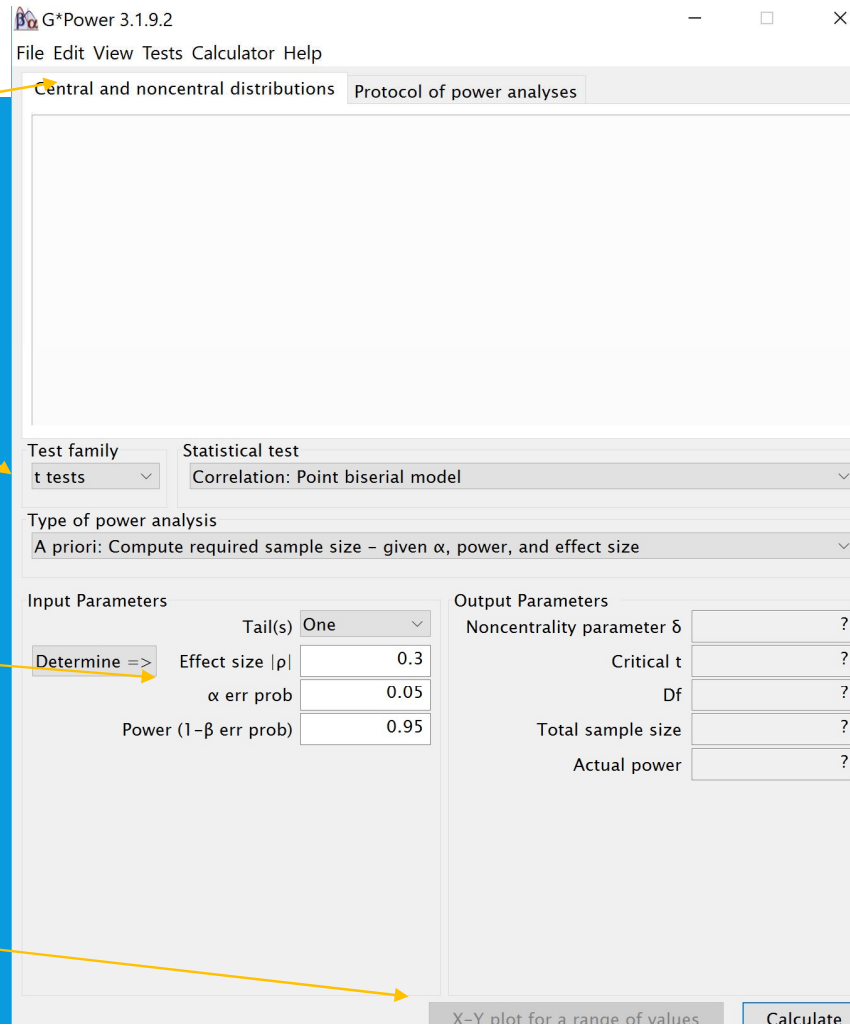
Tests

- Correlation and regression
- Means
- Proportions
- Variances
- Generic

- One group: Difference from constant
- One group: Wilcoxon (non-parametric)
- Two dependent groups (matched pairs)
- Two dependent groups (matched pairs): Wilcoxon (non-parametric)
- Two independent groups
- Two independent groups: Wilcoxon (non-parametric)
- Many groups: ANCOVA: Main effects and interactions
- Many groups: ANOVA: One-way (one independent variable)
- Many groups: ANOVA: Main effects and interactions (two or more independent variables)
- Repeated measures: Between factors, ANOVA-approach
- Repeated measures: Between factors, MANOVA-approach
- Repeated measures: Within factors, ANOVA-approach
- Repeated measures: Within factors, MANOVA-approach
- Repeated measures: Within-between interactions, ANOVA-approach
- Repeated measures: Within-between interactions, MANOVA-approach
- Multivariate: Hotelling T^2 , one group
- Multivariate: Hotelling T^2 , two groups
- Multivariate: MANOVA: Global effects
- Multivariate: MANOVA: Special effects and interactions

PHD School – University of Milano-Bicocca a.y. 2018/19

G*Power



Analysis

Type of power analysis

Inputs

Action buttons

Output

Two sample means

Analysis

Type of power analysis

Inputs

Test family	Statistical test
t tests	Means: Difference between two independent means (two groups)
Type of power analysis	
A priori: Compute required sample size - given α , power, and effect size	
Input Parameters	
Determine =>	Tail(s) One
Effect size d	0.5
α err prob	0.05
Power (1- β err prob)	0.80
Allocation ratio N2/N1	1
Output Parameters	
Noncentrality parameter δ	2.5248762
Critical t	1.6602343
Df	100
Sample size group 1	51
Sample size group 2	51
Total sample size	102
Actual power	0.8058986
X-Y plot for a range of values	
Calculate	

Output

Example 2: Two repeated means

Analysis → Test family: t tests

Type of power analysis → Statistical test: Means: Difference between two dependent means (matched pairs)

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

Inputs → Input Parameters

Parameter	Value
Tail(s)	One
Effect size dz	0.5
α err prob	0.05
Power ($1-\beta$ err prob)	0.80

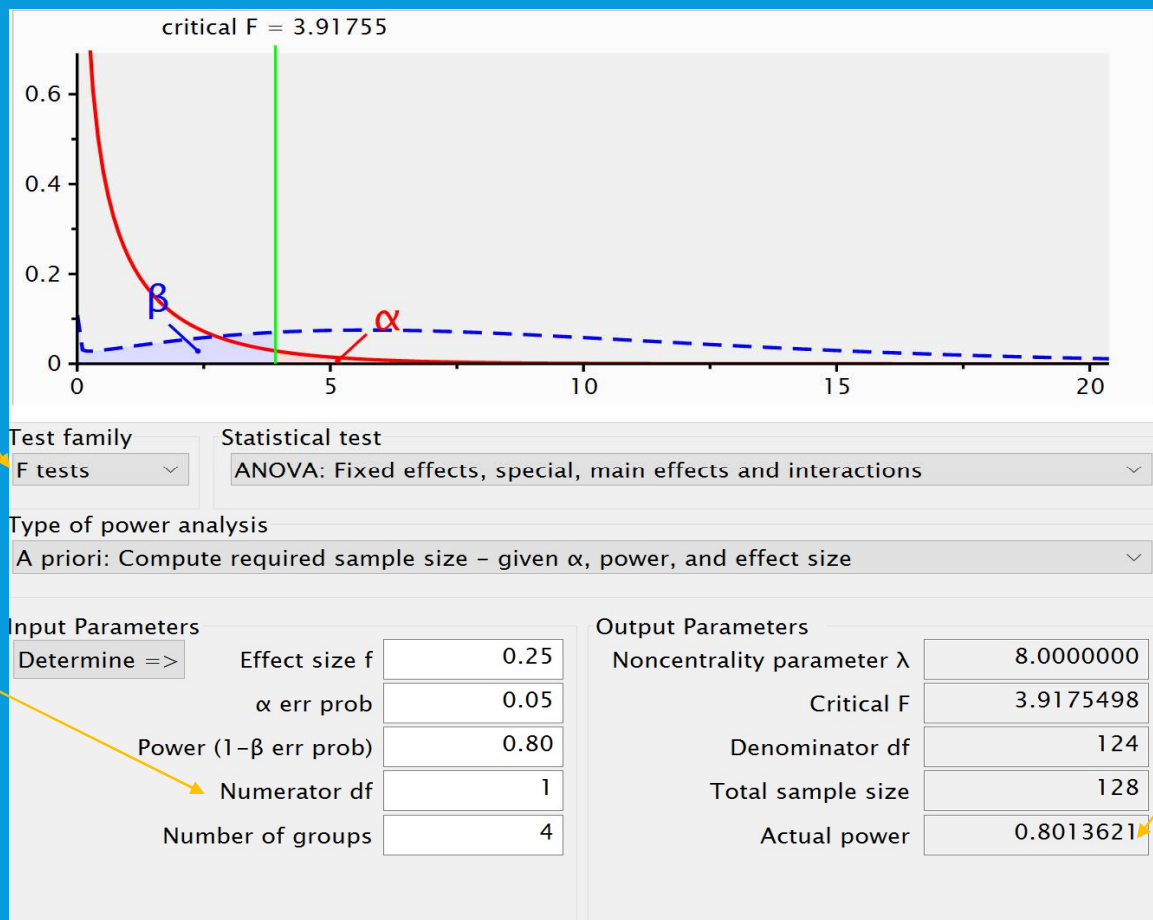
Output → Output Parameters

Parameter	Value
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

ANOVA 2 x 2

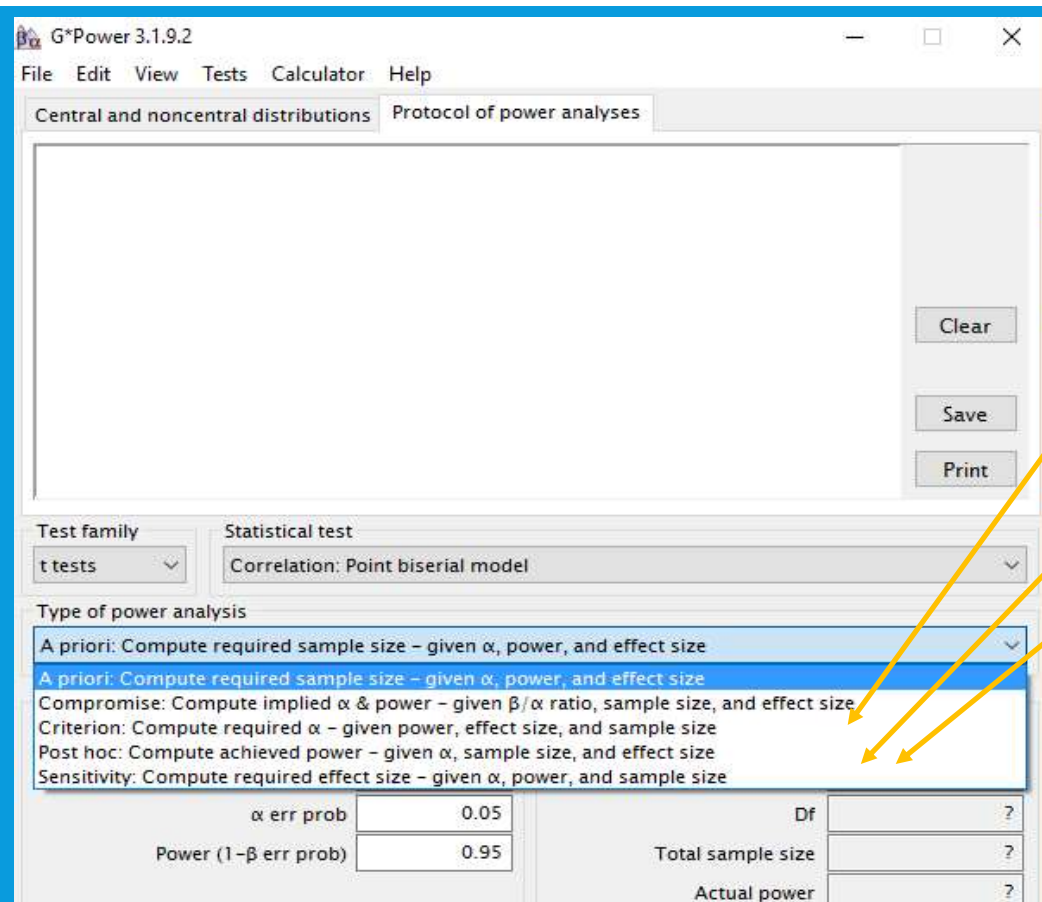
Analysis



Inputs

Output

G-POWER: PROTOCOL OF POWER ANALYSIS. HAVE WE FOUND WHAT WE ARE LOOKING FOR?



A priori: the 'ideal' way of determining power

Post hoc: the 'I could not help it' way of determining power

Sensitivity analysis: I have a feeling that there is an effect, but I do not have evidence so far.

How much do I have to increase my sample size to find (if I am wright) what I am looking for?

G-POWER: HOW TO CHOOSE THE GLIM

The screenshot displays the G-POWER software interface. The 'Tests' menu is open, showing a list of statistical tests. The 'Means' option is selected, which has opened a sub-menu with various test options. Below the menu, the 'Input Parameters' section is visible, showing the 'Determine =>' button and a table of parameters.

Tests Menu:

- Central and non-central distributions
- Correlation and regression analyses
- Means**
 - One group: Difference from constant
 - One group: Wilcoxon (non-parametric)
 - Two dependent groups (matched pairs)
 - Two dependent groups (matched pairs): Wilcoxon (non-parametric)
 - Two independent groups
 - Two independent groups: Wilcoxon (non-parametric)
- Proportions
- Variations
- Generic

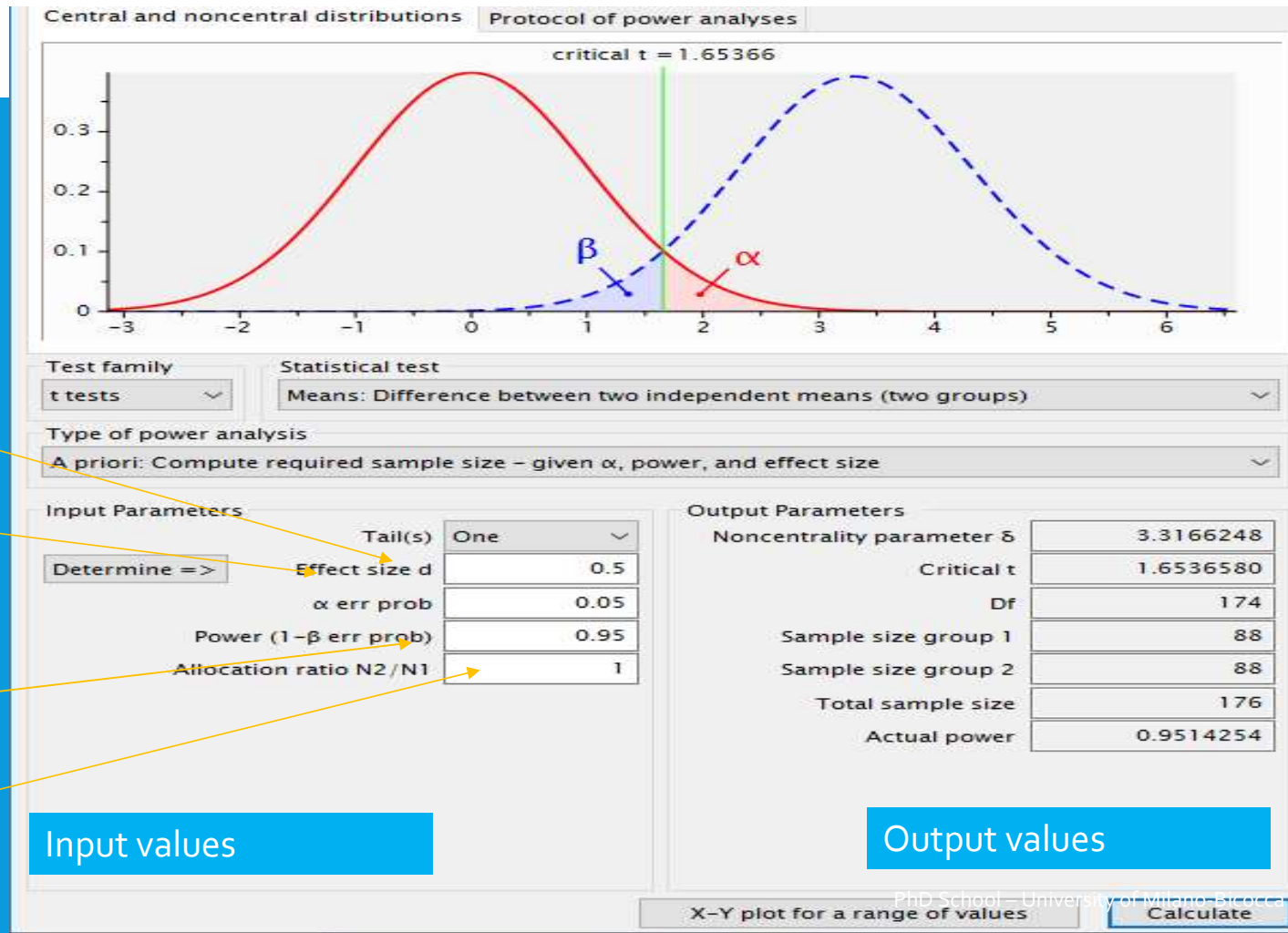
Input Parameters:

Parameter	Value
Effect size $ \rho $	0.3
α err prob	0.05
Power ($1 - \beta$ err prob)	0.55

Actual power: ?

PhD School – University of Milano-Bicocca a.s. 2016/17

A COMPLETED EXERCISE : POWER ANALYSIS FOR TEST FOR TWO INDEPENDENT MEANS



Desired effect size $d=0.5$

Type I error, 0.05

Power wanted (in this case, very high, 0.95)

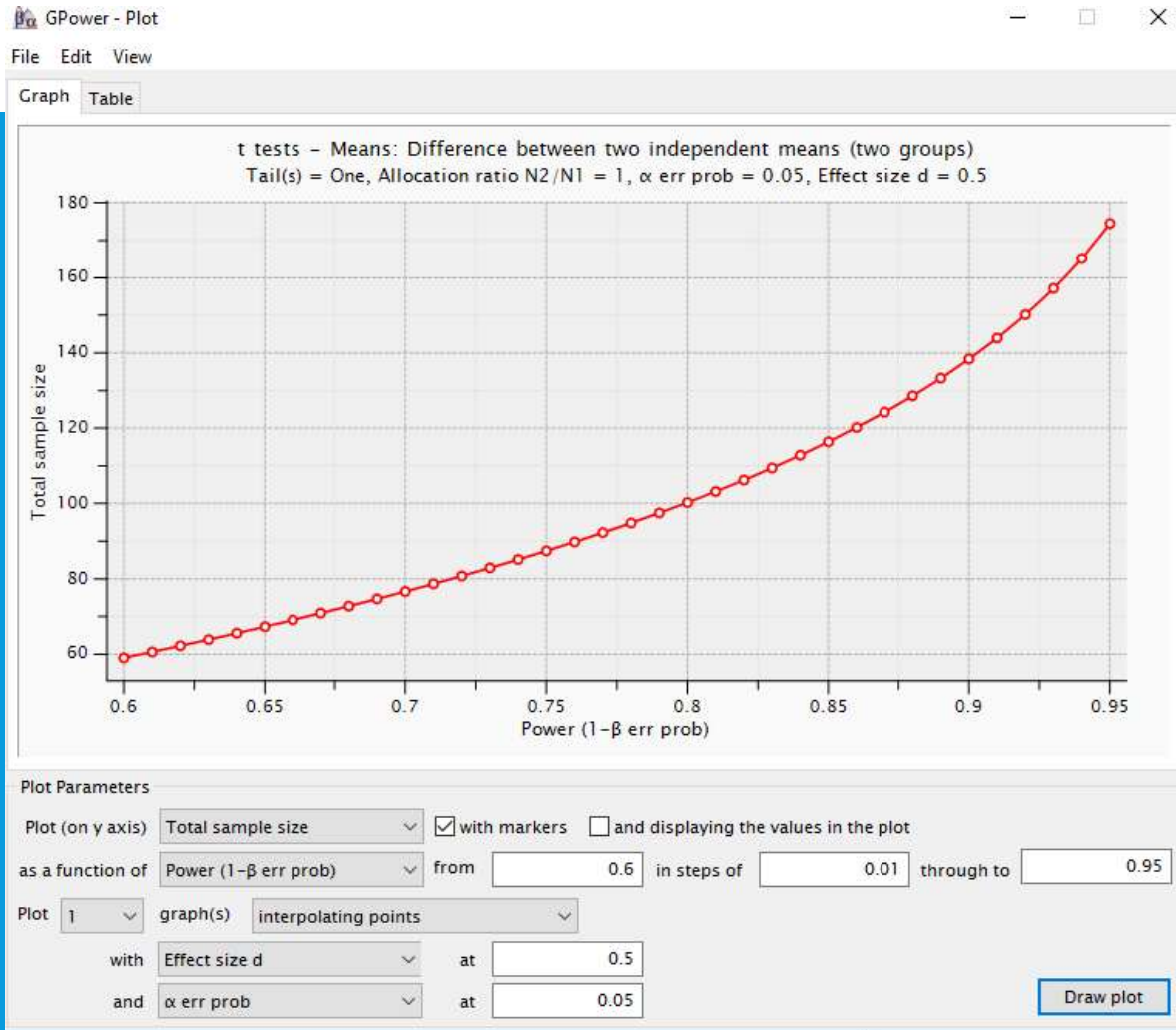
Equal sample sizes,

Sample sizes 88, in total 176 cases

Input values

Output values

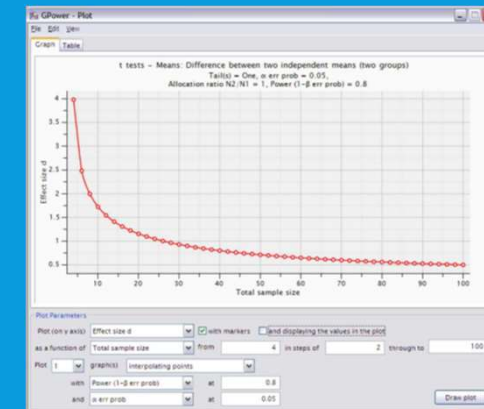
A COMPLETED EXERCISE: SENSITIVITY ANALYSIS FOR T TEST FOR TWO INDEPENDENT MEANS



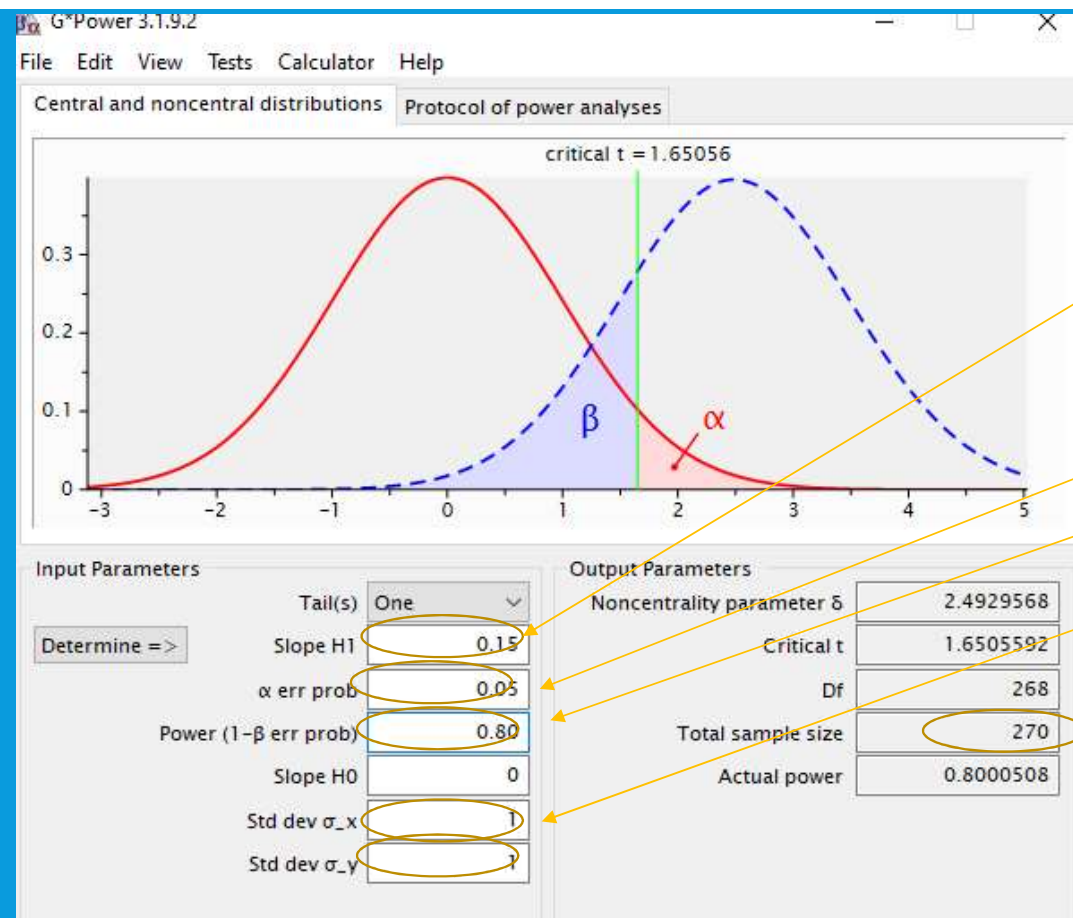
Same Gpower session as in the previous same analysis

Total sample size (y axis)for different power levels (x axis)

Below an example of curve for estimated ES as a function of the sample size n



POWER ANALYSIS: WHAT DO WE EXPECT IN A REGRESSION?



- Sample size in multiple regression
 - Test for single coefficient (slope=0.15): bilateral t-test
 - Alpha=.05
 - Power .80 (minimum)
 - SE for both variables
-
- Output: total sample size: 270

POWER ANALYSIS IN R

[HTTPS://CRAN.R-PROJECT.ORG/WEB/PACKAGES/PAMM/PAMM.PDF](https://cran.r-project.org/web/packages/pamm/pamm.pdf)

Search Results



The search string was "pwr"

Help pages:

pwr::pwr-package	Basic power calculations pwr
pwr::pwr.2p.test	Power calculation for two proportions (same sample sizes)
pwr::pwr.2p2n.test	Power calculation for two proportions (different sample sizes)
pwr::pwr.anova.test	Power calculations for balanced one-way analysis of variance tests
pwr::pwr.chisq.test	power calculations for chi-squared tests
pwr::pwr.f2.test	Power calculations for the general linear model
pwr::pwr.norm.test	Power calculations for the mean of a normal distribution (known variance)
pwr::pwr.p.test	Power calculations for proportion tests (one sample)
pwr::pwr.r.test	Power calculations for correlation test
pwr::pwr.t.test	Power calculations for t-tests of means (one sample, two samples and paired samples)
pwr::pwr.t2n.test	Power calculations for two samples (different sizes) t-tests of means

Random effects
ICC (effect size) in package
sjstats
Power analysis for random
effects in package pamm (and
others)
effect size : Package 'effsize'

Simulations
Superpower in R
Lakens, D., & Caldwell, A. R.
(2019). "Simulation-Based
Power-Analysis for Factorial
ANOVA Designs"

```
> pwr.t.test(n = 30, d = 0.5, sig.level = 0.05)
```

```
Two-sample t test power calculation
```

```
      n = 30
      d = 0.5
sig.level = 0.05
  power = 0.4778965
alternative = two.sided
```

NOTE: n is number in *each* group

Very low power

'ex post', it is debatable theoretically for some methodologists.

We need at least .8 power (Lakens is requires .9!)

```
> pwr.t.test(d = 0.5, power = 0.80, sig.level = 0.05)
```

```
Two-sample t test power calculation
```

```
      n = 63.76561
      d = 0.5
sig.level = 0.05
  power = 0.8
alternative = two.sided
```

NOTE: n is number in *each* group

This is 'ex ante', we need 64 subjects in each group

d: effect size

n= sample size

We can either specify directly d value or the 'intended' difference between means D and the pooled standard deviation

```
> power.t.test(delta = 0.50, sd = 2.25, sig.level = 0.05, power = 0.8)
```

```
Two-sample t test power calculation
```

```
      n = 318.8428
  delta = 0.5
    sd = 2.25
sig.level = 0.05
  power = 0.8
alternative = two.sided
```

```
NOTE: n is number in *each* group
```

```
>
```

Now, we need a different es . We know that the pooled sd is 2.25, we need to detect a difference between the means equal or smaller than .50. This difference is called δ in library `pwr`

We need at least 638 subjects

Power curves (power functions)



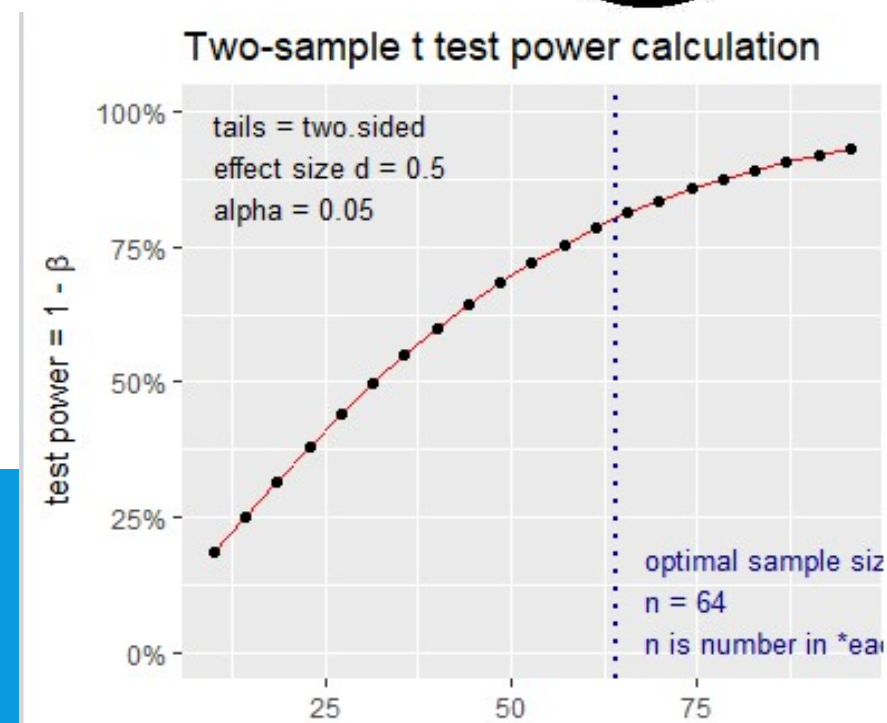
```
> power.t.test(delta = 0.50, sd = 2.25, sig.level = 0.05, power = 0.8)

Two-sample t test power calculation

      n = 318.8428
      delta = 0.5
      sd = 2.25
      sig.level = 0.05
      power = 0.8
      alternative = two.sided

NOTE: n is number in *each* group

> pwr.t.test(d=0.5, power = 0.8, type= "two.sample", alternative = "two.sided")
> plot(pwr.t.test)
> plot(pwr.t.test, xlab="sample size per group")
>
```



pwr - developed by Stéphane Champely- power analysis as outlined by Cohen (1988)

- <https://www.statmethods.net/stats/power.html>

- powerAnalysis: power for experimental design

<https://cran.r-project.org/web/packages/powerAnalysis/powerAnalysis.pdf>

- simr - Power Analysis for Generalised Linear Mixed Models (lme4) by Simulation

<https://cran.r-project.org/web/packages/simr/simr.pdf>

CONSEQUENCES OF LOW POWER

- Low probability of finding true effects: low power means that the chance of discovering effects that are really true is low. Low-powered studies produce more false negatives than high-powered studies.
- Low positive predictive value : the lower the power of a study, the lower the probability that an observed “significant” effect (among of all significant effects) actually reflects a true non-zero effect in the population (vs. a false positive).
- When an underpowered study discovers a true effect, it is likely that the estimate of the magnitude of that effect will be exaggerated. Effect inflation is worst for small, low-powered studies, because they can only detect sample parameter estimates effects when they are large.

ENJOY BEING POWERFUL!

