

## Sample size calculation with GPower

Wilfried Cools (ICDS) & Sven Van Laere (BiSI)  
<https://www.icds.be/>

# sample size calculation: why & how

- what about my sample size ?
- more is better  $\equiv$  observation  $\rightarrow$  information
  - but increasingly less so
  - but limited resources (time/money)
  - but ethical and practical considerations
- how many is good enough ?
  - depends on what you want
- workshop  $\rightarrow$  how to calculate it ?
  - understand the reasoning
  - apply for simple but common situations
- not one simple formula for all  $\rightarrow$  GPower to the rescue

2/61

# sample size calculation: a design issue

- linked to statistical inference
  - **testing**  $\rightarrow$  power [probability to detect existing **effects**]
  - **estimation**  $\rightarrow$  accuracy [size of **confidence intervals**]
- before data collection, during design of study
  - requires understanding: future data, analysis, inference (effect size, focus, ...)
  - conditional on assumptions & decisions
- not always possible nor meaningful !
  - easier for experiments (control), less for observational studies
  - easier for confirmatory studies, much less for exploratory studies
  - NO retrospective power analyses  $\rightarrow$  OK for future study only

Hoenig, J., & Heisey, D. (2001). The Abuse of Power: The Pervasive Fallacy of Power Calculations for Data Analysis.

*The American Statistician*, 55, 19–24.

# simple example

- evaluation of radiotherapy to reduce a tumor in mice
- comparing treatment group with control (=conditions)
- tumor induced, random assignment treatment or control (equal if no effect)
- after 20 days, measurement of tumor size (=observations)
- analysis:
  - unpaired t-test to compare averages for treatment and control
- goal:
  - if the average size in treatment is at least 20% less than control then we want to detect it (significance)
- the main issue:
  - how to calculate the required sample size to determine the effect aimed for ?

4/61

# overview

- PART I: building blocks in action for t-test
  - sizes: effect size, sample size
  - errors: type I ( $\alpha$ ), type II ( $\beta$ )
  - distributions:  $H_0$ ,  $H_A$
  - criterion: confidence (estimation), power (testing)
- PART II: moving beyond independent t-test
  - dependent groups
  - non-parametric distributions
  - multiple groups (ANOVA: omnibus, pairwise, focused)
  - proportions, correlations, ...

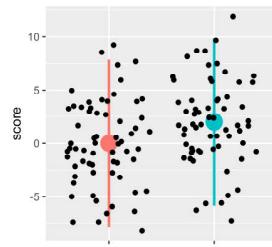
## PART I: building blocks in action for t-test

5/61

6/61

## reference example

- sample sizes easy and meaningful to calculate for well understood problems
- apriori specifications
  - intend to perform a statistical test
  - comparing 2 equally sized groups
  - to detect difference of at least 2
  - assuming an uncertainty of 4 SD on each mean
  - which results in an effect size of .5
  - evaluated on a Student t-distribution
  - allowing for a type I error prob. of .05 ( $\alpha$ )
  - allowing for a type II error prob. of .2 ( $\beta$ )
- sample size conditional on specifications being true



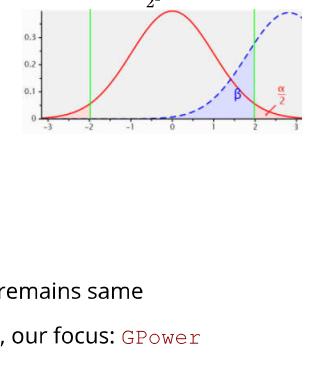
<https://icds.shinyapps.io/shinyt/>

7/61

## a formula you could use

- for this particular case:
  - sample size ( $n \rightarrow ?$ )
  - difference ( $d=\text{signal} \rightarrow 2$ )
  - uncertainty ( $\sigma=\text{noise} \rightarrow 4$ )
  - type I errors ( $\alpha \rightarrow .05$ , so  $Z_{\alpha/2} \rightarrow -1.96$ )
  - type II errors ( $\beta \rightarrow .2$ , so  $Z_{\beta} \rightarrow -0.84$ )
- sample size = 2 groups x 63 observations = 126
- note: formula's are test and statistic specific but logic remains same
- this and other formula's implemented in various tools, our focus: [GPower](#)

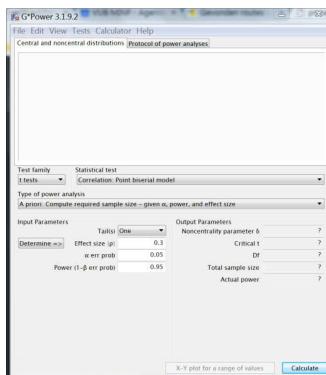
$$n = \frac{(Z_{\alpha/2} + Z_{\beta})^2 * 2 * \sigma^2}{d^2}$$
$$n = \frac{(-1.96 - 0.84)^2 * 2 * 4^2}{2^2} = 62.79$$



8/61

# GPower: a useful tool

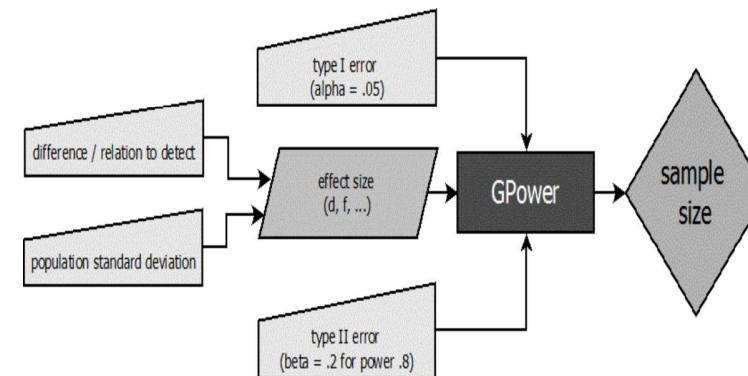
- popular and well established
- free @ <http://www.gpower.hhu.de/>
- implements wide variety of tests
- implements various visualizations
- documented fairly well
- note: not all tests are included !



9/61

# GPower: the building blocks in action

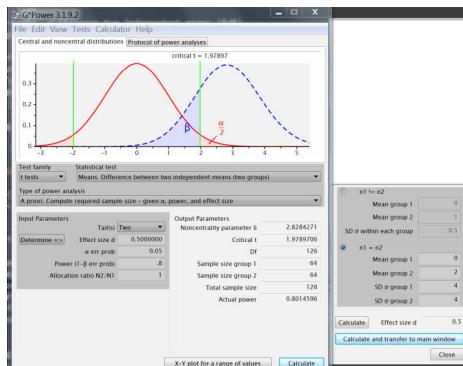
- For this example, calculate sample size based on
  - effect size (difference of interest, scaled on standard deviation)
  - type I and type II error



10/61

## GPower input

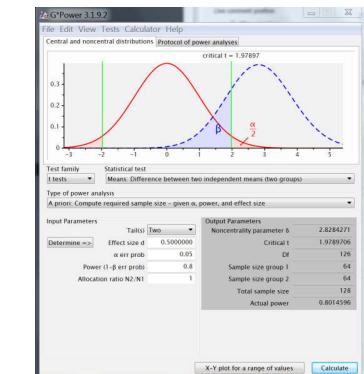
- t-test : difference two indep. means
- apriori: calculate sample size
- effect size = standardized difference
  - Cohen's  $d$
  - $d = |\text{difference}| / \text{SD\_pooled}$
  - $d = |0-2| / 4 = .5$
- $\alpha = .05$ , 2 - tailed ( $\alpha/2 \rightarrow .025 \& .975$ )
- $power = 1 - \beta = .8$
- allocation ratio = 1
- ~ reference example



11/61

## GPower output

- sample size ( $n$ ) =  $64 \times 2 = 128$
- degrees of freedom ( $df$ ) =  $126$  ( $128 - 2$ )
- plot showing null  $H_0$  and alternative  $H_a$  distribution
  - in GPower central and non-central distribution
  - $H_0$  & critical value → decision boundaries
    - critical  $t = 1.979$ ,  $qt(.975, 126)$
  - $H_a$ , shift with non-centrality parameter → truth
    - non centrality parameter ( $\delta$ ) =  $2.8284$   
 $2 / (4 * \sqrt{2}) * \sqrt{64}$
- $power \geq .80 (1 - \beta) = 0.8015$



12/61

# reference example protocol

t tests - Means: Difference between two independent means (two groups)

Analysis: A priori: Compute required sample size

Input: Tail(s) = Two

Effect size  $d = 0.5000000$

$\alpha$  err prob = 0.05

Power ( $1-\beta$  err prob) = .8

Allocation ratio N2/N1 = 1

Output: Noncentrality parameter  $\delta = 2.8284271$

Critical t = 1.9789706

Df = 126

Sample size group 1 = 64

Sample size group 2 = 64

Total sample size = 128

Actual power = 0.8014596

distributions

their cut-offs (type I and II errors)

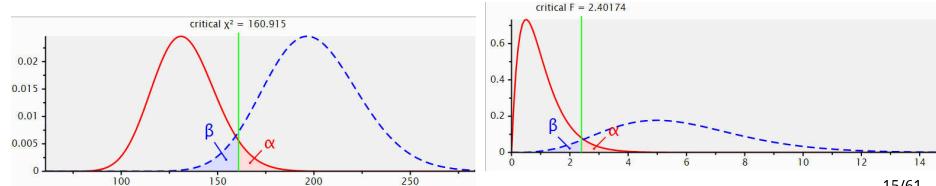
the distance between them (sample and effect sizes)

13/61

14/61

## GPower distributions

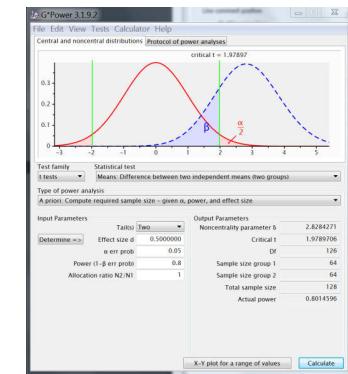
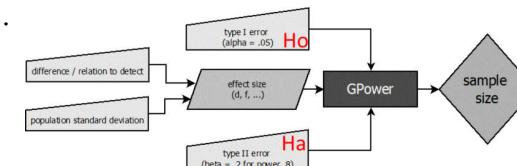
- distribution based test selection
  - Exact Tests (8)
  - $t$ -tests (11) → reference
  - $z$ -tests (2)
  - $\chi^2$ -tests (7)
  - $F$ -tests (16)
- focus on the density functions
- design based test selection
  - correlation & regression (15)
  - means (19) → reference
  - proportions (8)
  - variances (2)
- focus on the type of parameters



15/61

## Ho and Ha distributions

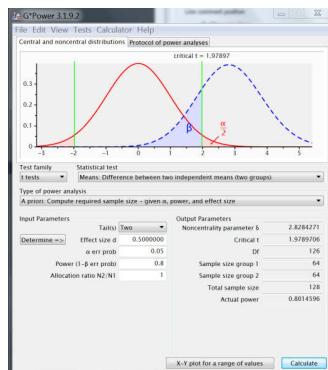
- $H_0$  acts as *benchmark* → eg., no difference
  - $H_0 \sim t(0, df)$  *cutoff* using  $\alpha$ ,
  - reject  $H_0$  if test returns *implausible* value
- $H_a$  acts as *truth* → eg., difference of .5 SD
  - $H_a \sim t(ncp, df)$
  - $ncp$  as violation of  $H_0$  → shift (location/shape)



16/61

# non-centrality: $H_0 \rightarrow H_a$

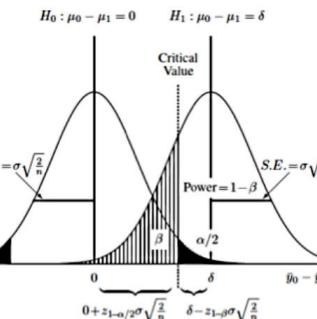
- $ncp$ : non-centrality parameter
  - shift between  $H_0$  and  $H_a$ 
    - assumed effect size (target or signal)
    - conditional on sample size (information)
  - overlap → power or sample size using  $\alpha$  on  $H_0$  and  $\beta$  on  $H_a$
- $H_a$  is NOT interchangeable with  $H_0$ 
  - absence of evidence  $\neq$  evidence of absence
  - equivalence testing ( $H_a$  for 'no effect')



<https://icds.shinyapps.io/shinyt/>

17/61

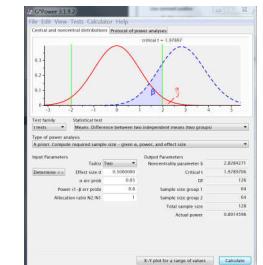
# divide by n perspective on distributions



- remember:
 
$$n = \frac{(Z_{\alpha/2} + Z_{\beta})^2 * 2 * \sigma^2}{d^2}$$

$$n = \frac{(-1.96 - 0.84)^2 * 2 * 4^2}{2^2}$$

$$n = 62.79$$



18/61

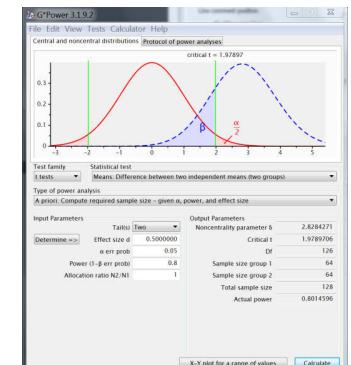
# no $H_a$ for estimation

- focus on estimation, plausible values of effect, not testing or power
- sample size without type II error  $\beta$ , power,  $H_0$  or  $H_a$
- ~ divide by n perspective but shifted to estimate
- precision analysis → set maximum width confidence interval
  - let E = maximum half width of confidence interval to accept
  - for confidence level  $1 - \alpha$
  - $n = z_{\alpha/2}^2 * \sigma^2 * 2 / E^2$  (for 2 groups)
- equivalence with statistical testing
  - if 0 (or other reference) outside confidence bounds → significant
- NOT GPower

19/61

# type I/II error probability

- distribution cut-off's (density → AUC=1)
- decide whether to reject  $H_0$  assuming  $H_a$
- two types of error
  - $P(\text{infer}=\text{Ha} | \text{truth}=\text{H}_0) = \alpha$
  - $P(\text{infer}=\text{H}_0 | \text{truth}=\text{Ha}) = \beta$
- two types of correct inference
  - $P(\text{infer}=\text{H}_0 | \text{truth}=\text{H}_0) = 1 - \alpha$
  - $P(\text{infer}=\text{Ha} | \text{truth}=\text{Ha}) = 1 - \beta \rightarrow \text{power}$
- cut-off 'known'  $H_0$  for statistical test
  - two tailed → both sides informative on  $H_0$
  - one tailed → one side not informative on  $H_0$

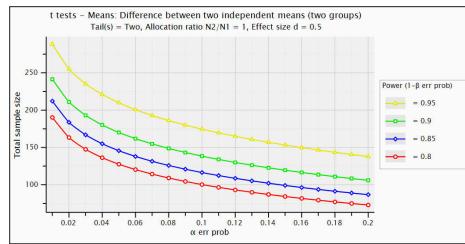


	infer=Ha	infer=Ho	sum
truth=H0	$\alpha$	$1-\alpha$	1
truth=Ha	$1-\beta$	$\beta$	1

20/61

## error exercise : create plot

- create plot  
(X-Y plot for range of values)
- plot sample size by type I error
- set plot to 4 curves
  - for power .8 in steps of .05
- set  $\alpha$  on x-axis
  - from .01 to .2 in steps of .01
- use effect size .5

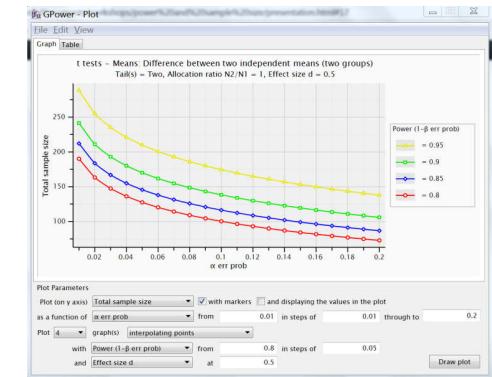


notice Table option

21/61

## error exercise : interpret plot

- where on the red curve (right)  
type II error = 4 \* type I error ?
- when smaller effect size (.25),  
what changes ?
- switch power and sample size  
(32 in step of 32)  
what is relation type I and II  
error ?

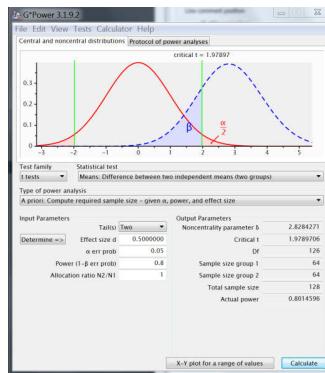


- where on the yellow curve (left)  
type II error = 4 \* type I error ?

22/61

## decide type I/II error probability

- rules of thumb ?
  - $\alpha$  in range .01 - .05 → 1/100 - 1/20
  - $\beta$  in range .1 to .2 → power = 80% to 90%
- $\alpha$  &  $\beta$  inversely related
  - if  $\alpha = 0$  → never reject, no power
  - if power 99% → high  $\alpha$  for same sample size
- determine the balance
  - which error you want to avoid most ?
    - cheap aids test? → avoid type II
    - heavy cancer treatment? → avoid type I
  - $\alpha$  &  $\beta$  often selected in 1/4 ratio  
type I error is 4 times worse !!



23/61

## interim analyses: control type I error

- analyze and proceed ? (peeking)
  - multiple testing as data is collected
  - inflates type I error  $\alpha$
- correct  $\alpha$ 
  - interim analysis specific  $\alpha_i$  with overall  $\alpha$  under control
- suggested technique: alpha spending
  - plan in advance
  - use O'Brien - Fleming bounds
  - NOT GPower
    - determine with simulation tool: <https://icds.shinyapps.io/aspending/>
    - commercial packages eg., PASS and other software (eg., ldbounds in R)

24/61

# for fun: P(effect exists | test says so)

- power → P(test says there is effect | effect exists)
- $P(\text{infer} = Ha | \text{truth} = H_0) = \alpha$
- $P(\text{infer} = H_0 | \text{truth} = Ha) = \beta$
- $P(\text{infer} = Ha | \text{truth} = Ha) = \text{power}$
- $P(\text{truth} = Ha | \text{infer} = Ha) = \frac{P(\text{infer} = Ha | \text{truth} = Ha) * P(\text{truth} = Ha)}{P(\text{infer} = Ha)}$  → Bayes Theorem
- $P(\text{truth} = Ha | \text{infer} = Ha) = \frac{P(\text{infer} = Ha | \text{truth} = Ha) * P(\text{truth} = Ha)}{P(\text{infer} = Ha | \text{truth} = Ha) * P(\text{truth} = Ha) + P(\text{infer} = H_0 | \text{truth} = Ha) * P(\text{truth} = H_0)}$
- $P(\text{truth} = Ha | \text{infer} = Ha) = \frac{\text{power} * P(\text{truth} = Ha)}{\text{power} * P(\text{truth} = Ha) + \alpha * P(\text{truth} = H_0)}$  → depends on prior probabilities
- IF very low probability model is true, eg., .01 ? →  $P(\text{truth} = Ha) = .01$
- THEN probability effect exists if test says so is low, in this case only 14% !!
- $P(\text{truth} = Ha | \text{infer} = Ha) = \frac{.8 * .01}{.8 * .01 + .05 * .99} = .14$

25/61

# effect sizes

- degree with which a certain phenomenon holds ( $\sim H_0$  is false)
  - part of non-centrality (as is sample size) → shift in GPower
  - signal to noise ratio
    - typically not just the signal, to provide scale
    - e.g., difference on scale of pooled standard deviation
  - bigger effect → more easy to detect it (pushing away  $H_0$ )
- 2 main families of effect sizes → test specific
  - differences **d-family** // association **r-family**
  - transformations, e.g.,  $d = .5 \rightarrow r = .243$ 
    - $d = \frac{2r}{\sqrt{1-r^2}}$ ;  $r = \frac{d}{\sqrt{d^2+4}}$ ;  $d = \ln(OR) * \frac{\sqrt{3}}{\pi}$

26/61

# effect sizes of Cohen

Table 1  
ES Indexes and Their Values for Small, Medium, and Large Effects

Test	ES index	Effect size		
		Small	Medium	Large
1. $m_a$ vs. $m_b$ for independent means	$d = \frac{m_a - m_b}{\sigma}$	.20	.50	.80
2. Significance of product-moment $r$	$r$	.10	.30	.50
3. $r_s$ vs. $r_b$ for independent proportions	$q = z_a - z_b$ where $z = \text{Fisher's } z$	.10	.30	.50
4. $P = .5$ and the sign test	$g = P - .50$	.05	.15	.25
5. $\phi$ for independent proportions	$h = \phi_A - \phi_B$ where $\phi = \text{arcsine transformation}$	.20	.50	.80
6. Chi-square for goodness of fit and contingency	$w = \sqrt{\sum_{i=1}^k \frac{(P_{ij} - E_{ij})^2}{E_{ij}}}$	.10	.30	.50
7. One-way analysis of variance	$f = \frac{s_m}{s_e}$	.10	.25	.40
8. Multiple and multiple partial correlation	$f^2 = \frac{R^2}{1-R^2}$	.02	.15	.35

- Cohen, J. (1992). A power primer. Psychological Bulletin, 112, 155–159.

27/61

# effect sizes of d and r family

Measures of group differences (the <b>d</b> family)		Measures of association (the <b>r</b> family)		Measures of association (the <b>r</b> family)
(a) Groups compared on dichotomous outcomes	(a) Correlation indexes	(b) Proportion of variance indexes	(b) Coefficient of determination	
RD The risk difference in probabilities:	$r$ The Pearson product moment correlation coefficient used when both variables are measured on an interval or ratio scale.	$r^2$ The coefficient of determination: analysis of variance expression	$R^2$ R squared, or the uncorrected coefficient of multiple regression analysis	
The probability of an event or outcome occurring in one group is higher than in another group.	$r$ (or $r_s$ ) Spearman's rho or the rank correlation coefficient used when both variables are measured on an ordinal or rank scale.	$R^2$ Adjusted R squared, or the coefficient of multiple regression analysis adjusted for sample size and the number of predictors.	$R^2_{\text{adj}}$ Adjusted R squared, or the coefficient of multiple regression analysis adjusted for sample size and the number of predictors.	
OR The odds ratio: the odds of an event or outcome occurring in one group with the odds of it occurring in another group.	$r$ The odds ratio: the odds of an event or outcome occurring in one group with the odds of it occurring in another group.	$f$ Cohen's $f$ quantifies the dispersion of means in three or more groups commonly used in ANOVA.	$f$ Cohen's $f$ quantifies the dispersion of means in three or more groups commonly used in ANOVA.	
(b) Groups compared on continuous outcomes	(b) Partial correlations	(b) Proportion of variance indexes	(b) Coefficients of correlation	
Cohen's $d$ , the difference between two groups based on the mean difference between two groups.	$r_{pb}$ The point-biserial correlation coefficient used when one variable is continuous and the other variable is dichotomous.	$r^2$ The proportion of variance explained by the independent variable(s) in the dependent variable(s).	$r^2$ The proportion of variance explained by the independent variable(s) in the dependent variable(s).	
Glass' delta or the uncorrected standardized mean difference between two groups based on the standard deviation of the control group.	$\psi$ The phi coefficient used when variables are dichotomous can be arranged in a 2x2 contingency table.	$\eta^2$ Eta squared or the (uncorrected) coefficient of multiple regression commonly used in ANOVA.	$\epsilon^2$ Epsilon squared: an unbiased estimate of error variance.	
Hedges' $\delta$ or the uncorrected standardized mean difference between two groups based on the pooled weighted standard deviation.	$C$ Peason's $c$ coefficient used when variables and effects can be measured on an interval or ratio scale.	$\omega^2$ Omega squared: an unbiased estimate of error variance.	$\Omega^2$ Omega squared: an unbiased estimate of error variance.	
Probability of superadding: the probability that a random value drawn from one group will be greater than a random value drawn from another.	$V$ Cramér's V: $C = \sqrt{\chi^2 / (N * k)}$ is a measure of effect size that can be used for tables of any size.	$R^2_C$ The coefficient of correlation squared: used for canonical correlation analysis.		
$\lambda$	Goodman and Kruskal's lambda: used for variables that are measured on nominal (or categorical) scales.			

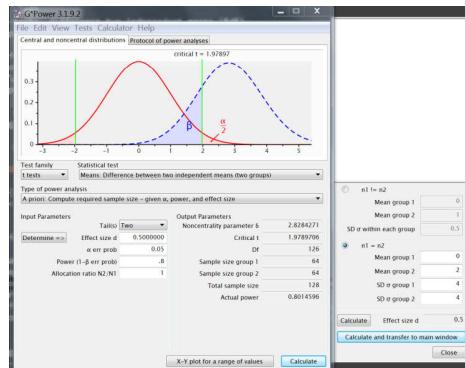
- most important effect size of
  - d: dichotomous - continuous
  - r: correlation - proportion variance

28/61

- Ellis, P. D. (2010). The essential guide to effect sizes: statistical power, meta-analysis, and the interpretation of research results.
- more than 70 different effect sizes... most of them related to each other
- NOT p-value ~ partly effect size, partly sample size or power
  - do not simply compare p-values !

# effect sizes in GPower (Determine)

- often very difficult to specify
- GPower offers help with **Determine**
  - difference group means  
0-2 → signal ~ minimally relevant (or expected)
  - standard deviations (sd)  
4 each group → expected noise ~ natural diversity
  - written to Effect Size  $d$   
.5 → difference in sd
- Reminder: effect size statistic depends on statistical test

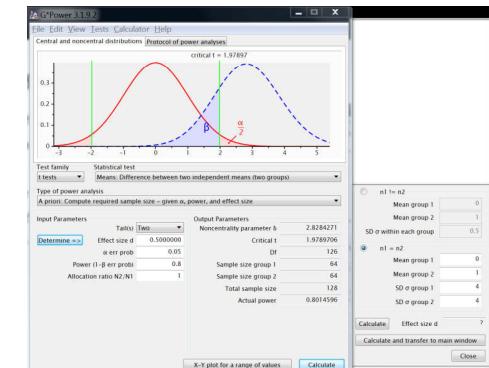


29/61

# effect size exercise : ingredients cohen d

For the **reference example**:

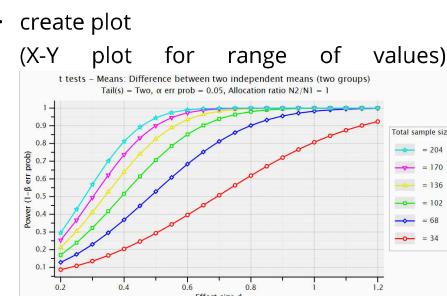
- change mean values from 0 and 2 to 4 and 6, what changes ?
- change sd values to 2 for each, what changes ?
  - effect size ?
  - total sample size ?
  - non-centrality ?
  - critical t ?
- change sd values to 6 for each, what changes ?



30/61

# effect size exercise : plot

- plot power by effect size
- set plot to 6 curves
  - for sample sizes, 34 in steps of 34
- set effect sizes on x-axis
  - from .2 to 1.2 in steps of .05
- use  $\alpha$  equal to .05
- determine (approximately) the three situations from previous slide on the plot
- how does power change when doubling the effect size, eg., from .5 to 1 ?

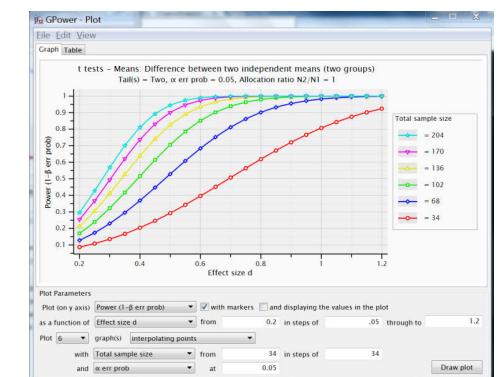
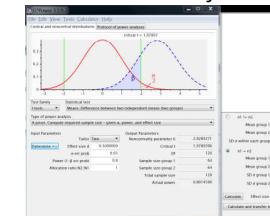


31/61

# effect size exercise : imbalance

For the **reference example**:

- change allocation ratio from 1
  - to 2, .5, 3 and 4, what to conclude ?
    - ratio 2 and .5 ?
    - imbalance + 1 or \* 2 ?
- ? no idea why  $n_1 \neq n_2$



32/61

# effect sizes, how to determine them in theory

- choice of effect size matters → justify choice
- choice of effect size
  - NOT significant → meaningless, dependent on sample size
  - realistic (eg., previously observed effect) → replicate
  - important (eg., minimally relevant effect)
- use **Determine** to get started (check the manual)
  - for independent t-test → means and standard deviations
  - possible alternative is to use variance explained, eg., 1 versus 16
    - with one-way ANOVA ( $f=.25$  instead of  $d=.5$ )
    - with linear regression ( $R^2=.0625$  instead of  $d=.5$ )
  - [https://www.psychometrica.de/effect\\_size.html#transform](https://www.psychometrica.de/effect_size.html#transform)

33/61

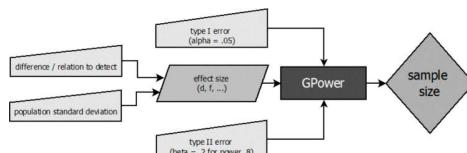
# effect sizes, how to determine them in practice

- experts / patients → use if possible → importance
- literature (earlier study / systematic review) → realistic
- pilot → guestimate dispersion estimate, but very small sample size
- internal pilot → stopping rule (sequential/conditional)
- turn to Cohen → use if everything else fails (rules of thumb)
- guesstimate the input parameters, what can you do ?
  - sd from assumed range / 6 assuming normal distribution
  - sd for proportions (& percentages) at conservative .5
  - sd from control, assume treatment the same
  -

34/61

# relation samples & effect size, errors I & II

- building blocks:
    - sample size ( $n$ )
    - effect size ( $\Delta$ )
    - alpha ( $\alpha$ )
    - power ( $1 - \beta$ )
  - each parameter conditional on others
- GPower → type of power analysis
    - Apriori:  $n \sim \alpha, \text{power}, \Delta$
    - Post Hoc:  $\text{power} \sim \alpha, n, \Delta$
    - Compromise:  $\text{power}, \alpha \sim \beta / \alpha, \Delta, n$
    - Criterion:  $\alpha \sim \text{power}, \Delta, n$
    - Sensitivity:  $\Delta \sim \alpha, \text{power}, n$



35/61

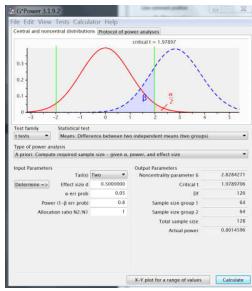
# type of power analysis exercise

- for given example, step through...
  - retrieve power given  $n, \alpha$  and  $\Delta$
  - [1] for power .8, take half the sample size, how does  $\Delta$  change ?
  - [2] set  $\beta/\alpha$  ratio to 4, what is  $\alpha & \beta$ ? what is the critical value ?
  - [3] keep  $\beta/\alpha$  ratio to 4 for effect size .5, what is  $\alpha & \beta$ ? critical value ?

36/61

# getting your hands dirty

# PART II: moving beyond independent t-test



- in R, assuming normality
  - `qt` → get quantile on  $H_0$  ( $Z_{1-\alpha/2}$ )
  - `pt` → get probability on  $H_a$  (non-central)

```
# calculator
m1=0;m2=2;s1=4;s2=4
alpha=.025;N=128
var=.5*s1^2+.5*s2^2
d=abs(m1-m2)/sqrt(2*var)
d=d*sqrt(N/2)
tc=tinv(1-alpha,N-1)
power=1-nctcdf(tc,N-1,d)
```

```
.n <- 64
.deg <- 2*.n-2
.ncp <- 2 / (4 * sqrt(2)) * sqrt(.n)
.power <- 1 -
  pt(
    qt(.975,df=.df),
    df=.df, ncp=.ncp
  ) -
  pt( qt(.025,df=.df), df=.df, ncp=.ncp)
round(.power,4)
```

```
## [1] 0.8015
```

37/61

38/61

## GPower examples beyond the independent t-test

- so far, comparing two independent means
- selected topics with small exercises
  - dependent instead of independent
  - non-parametric instead of assuming normality
  - relations instead of groups (regression)
  - correlations
  - proportions, dependent and independent
  - more than 2 groups (compare jointly, pairwise, focused)
  - more than 1 predictor
  - repeated measures
- GPower [manual](#) 27 tests: effect size, non-centrality parameter and example !!

39/61

## dependence between groups

- if 2 dependent groups (eg., before/after treatment) → account for correlations
- matched pairs (t-test / means, difference 2 dependent means)
- use [reference example](#)
  - [1] use correlation .5 to compare (effect size, ncp, n)
  - [2] how many observations if no correlation exists ([reference example](#))?
  - [3] difference in sample size for correlation .875 ?
  - [4] set original sample size (n=64\*2) and effect size (dz=.5), compare ?

40/61

## non-parametric distribution

- expect non-normally distributed residuals, avoid normality assumption
- only considers ranks or uses permutations → price is efficiency
- avoid when possible, eg., transformations
- two groups → Wilcoxon-Mann-Whitney (t-test / means, diff. 2 indep. means)
- use [reference example](#)
  - [1] how about n ? compared to parametric → what is % loss efficiency ?
  - [2] change parent distribution to 'min ARE' ? what now ?

## relations instead of group differences

- differences between groups → relation observations & categorization
- example →  $d = .5 \rightarrow r = .243$
- note: slope  $\beta = r * \sigma_y / \sigma_x$
- regression coefficient (t-test / regression, one group size of slope)
- sample size for comparing the slope Ha with 0 (=Ho)
  - [1] determine slope ( $\beta$ , with  $\sigma_y = 4.12$  and  $\sigma_x = .5$ )
  - [2] calculate sample size
  - [3] what if  $\sigma_x$  (predictor values) or  $\sigma_y$  (effect and error) increase ?

41/61

42/61

## relations: a variance perspective

- between and within group variance → relation observations & categorization
- regression coefficient (t-test / regression, fixed model single regression coef)
- use [reference example](#), regression style
  - variance within  $4^2$  and between  $1^2$ , totaling  $\sigma_y^2 = 17$
  - [1] calculate sample size, compare effect sizes ?
  - [2] what if also other predictors in the model ?
- note:  $f^2 = R^2 / (1 - R^2)$

## more groups to compare, 4 cases

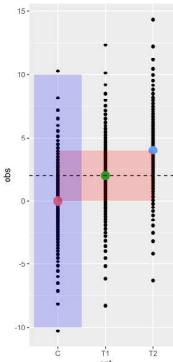
- simple example: assume one control, and two treatments
- if more than two groups, several options
  - test whether at least one differs → omnibus F-test (variances)
  - test whether all differ from each other → pairwise comparisons
  - test whether selected pairs differ → contrast (t-test)
  - test whether linear combinations of pairs differ → contrasts (t-tests)  
eg., control versus each of the average of treatments
- if multiple tests → inflation of type I error ( $\alpha$ )
  - correct  $\alpha$  or p-value, eg., using Bonferroni
  - make more tentative inferences

43/61

44/61

## F-test statistic

- multiple groups → not one effect size  $\text{d}$
- F-test statistic & effect size  $f$
- $f$  is ratio of variances  $\sigma_{\text{between}}^2 / \sigma_{\text{within}}^2$
- example: one control and two treatments
  - reference example + 1 group
  - within group observations normally distributed
  - means  $C=0, T1=2$  and  $T2=4$
  - sd for all groups ( $C, T1, T2$ ) = 4



## more groups: pairwise

- assume one control, and two treatments
  - interested in all three pairwise comparisons → maybe Tukey
    - typically run aposteriori, after omnibus shows effect
  - use t-test with correction of  $\alpha$  for multiple testing
- apply Bonferroni correction for original 3 group example
  - [1] resulting sample size for three tests ?
  - [2] what if biggest difference is ignored, sample size ?
  - [3] with original 64 sized groups, what is the power ?

45/61

## more groups: omnibus

- for one control and two treatments → test that at least one differs
- one-way Anova (F-test / Means, ANOVA - fixed effects, omnibus, one way)
- effect size  $f$ , with numerator/denominator df (derived from  $\eta^2$ )
- start from reference example,
  - [1] what is the sample size ? ncp ? critical F ? does size matter ?
  - [2] set extra group, either mean 1 or 4, what are the effect / sample sizes ? and with the mean of middle group away from global mean ?
  - [4] derive effect size with variance between .5 and within 3, 1 and 6 ?

46/61

## sample size calculation benefit from focus

- better to focus during the design on specific questions
  - only consider the main comparisons in focus (eg. primary endpoints)
    - only interested in comparing two treatments → t-test
  - only consider smallest of relevant effects, largest sample size
  - set up contrasts (next slide)
- sample size calculations (design)
  - not necessarily equivalent to statistics
  - requires justification to convince yourself and/or reviewers
- example:
  - statistics: group difference evolution 4 repeated measurements → mixed model
  - power: difference treatment and control last time point → t-test

47/61

48/61

## more groups: contrasts

- assume one control and two treatments
  - set up 2 contrasts for T1 - C and T2 - C
  - set up 1 contrast for average(T1,T2) - C
- each contrast requires 1 degree of freedom
- each contrast combines a specific number of levels
- effect sizes for planned comparisons must be calculated !!
  - contrasts (linear combination)
  - standard deviation of contrasts

$$\sigma_{contrast} = \frac{|\sum \mu_i * c_i|}{\sqrt{N \sum_i^k c_i^2 / n_i}}$$

with group means  $\mu_i$ , pre-specified coefficients  $c_i$ , sample sizes  $n_i$  and total sample size  $N$

49/61

50/61

## more groups: contrasts exercise

- one-way ANOVA (F-test / Means, ANOVA-fixed effects,special,main,interaction)
- obtain effect sizes for contrasts (assume equally sized for convenience)
  - $\sigma_{contrast}$  T1-C:  $\frac{(-1*0+1*2+0*4)}{\sqrt{2*((-1)^2+1^2+0^2)}} = 1$ ; T2-C:  $= 2$ ;  $(T1+T2)/2-C: = 1.4142$
  - with  $\sigma = 4 \rightarrow$  ratio of variances for effect sizes f .25, .5, .3536
- sample size for each contrast, each 1 df and 2 groups
  - [1] contrasts nrs. 1 or 2
  - [2] contrasts nrs. 1 AND 2
  - [3] contrasts nr. 3

## multiple factors

- multiway ANOVA (F-test / Means, ANOVA-fixed effects,special,main,interaction)
- multiple main effects and interaction effects
  - interaction: group specific difference between groups
    - degrees of freedom (#A-1)\*(#B-1)
  - main effects: if no interaction (#X-1)
  - get effect sizes for two way anova  
<https://icds.shinyapps.io/effectsizes/>
- sample size for **reference example**, assume a second predictor is trivial
  - [1] what is partial  $\eta^2$  ?
  - [2] sample size ?

51/61

52/61

## dependence within groups (repeated)

- if repeated measures  $\rightarrow$  account for correlations
- repeated measures (F-test / Means, repeated measures...)
- 3 main types
  - within: like dependent t-test for 2 or more measurements
  - between: use of multiple measurements per group
  - interaction: difference of change over groups
- correlation within subject (unit)
  - informative on group differences within subject
  - redundancy for between group differences

## repeated measures within

- possible to have only 1 group (within subject comparison)
- use effect size  $f = .25$  (1/16 explained versus unexplained)
  - [1] use zero correlation to compare with sample size independent t-test
  - [2] for one group use correlation .5, compare sample size dependent t-test
  - [3] double number of groups to 2
  - [4] double number of measurements to 4 (correlation 0 and .5), impact ?

## repeated measures between

- use effect size  $f = .25$  (1/16 for variance or 2/4 for means)
  - [1] use correlation 0 and .5 with 2 groups and 2 measurements, sample size ?
  - [2] for correlation .5, compare 2 or 4 measurements, sample size ?
  - [3] double number of groups to 2

53/61

54/61

## repeated measures within x between

- SPSS idiosyncrasies: <https://www.youtube.com/watch?v=CEQUNYg80Y0>
- use effect size  $f = .25$  (1/16 for variance)
  - [1] use correlation 0, compare 2 groups 2 measurements with rep. between ?
  - [2] use correlation 0.5, compare 2 groups 2 measurements with rep. within ?
  - [3] use correlation .5, compare 2 groups and 4 measurements, sample size ?
  - [4] repeat with 4 groups and 4 measurements, sample size ?

## correlations

- when comparing two independent correlations
- z-tests / correlation & regressions: 2 indep. Pearson r's
- makes use of Fisher Z transformations  $\rightarrow z = .5 * \log(\frac{1+r}{1-r}) \rightarrow q = z_1 - z_2$
- [1] assume correlation coefficients .7844 and .5 effect size & sample size ?
- [2] assume .9844 and .7, effect size & sample size ?
- [3] assume .1 and .3844 effect size & sample size ?

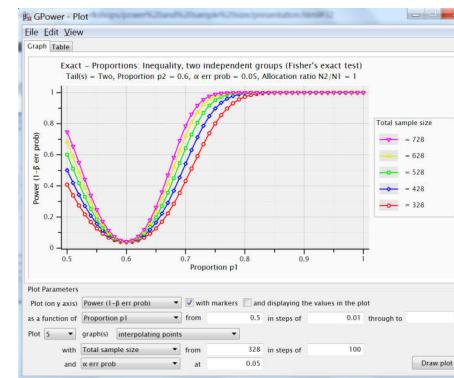
55/61

56/61

# proportions

- comparing two independent proportions → bounded between 0 and 1
- Fisher Exact Test (exact / proportions, difference 2 independent proportions)
- effect sizes in odds ratio, relative risk, difference proportion
  - [1] for odds ratio 2,  $p_2 = .60$ , what is  $p_1$  ?
  - [2] sample size for equal sized, and type I and II .05 and .2 ?
  - [3] sample size when .95 and .8 (difference of .15) and .05 and .2 ?
- treat as if unbounded, ok within .2 - .8, variance is  $p*(1-p)$  → maximally .25 !!
  - [4] use t-test for difference of .15

# proportions exercise



- Fisher Exact Test
  - power over proportions .5 to 1
  - 5 curves, sample sizes 328, 428, 528...
  - type I error .05
  - [1] what happens ?
  - [2] repeat for one-tailed, what is different ?

57/61

58/61

# dependent proportions

- when comparing two dependent proportions
- McNemar test (exact / proportions, difference 2 dependent proportions)
  - include correlations implicitly, discordant pairs → change
  - effect size as odds ratio → ratio of discordance ?!
- assume odds ratio equal to 2, equal sized, type I and II errors .05 and .2, two-way
  - [1] what is the sample size for .25 proportion discordant, and [2] .5, and [3] 1
  - [4] for odds ratio 4 or .25, how the proportion  $p_{12}$  and  $p_{21}$  change ?
  - [5] repeat for third alpha option, and consider total sample size, what happens ?

# conclusion: keep it simple, keep it real

- sample size calculation is a design issue, not a statistical one
- building blocks: sample & effect sizes, type I & II errors, each conditional on rest
- effect sizes express the amount of signal compared to the background noise
- complex models imply complex sample size calculations, if at all possible
- GPower deals with not too complex models
  - simplify using a focus, if justifiable → then GPower can get you a long way
  - use more complex specification for more complex sample size calculations
  - leave GPower, simulation is always an option

59/61

60/61