

Power and Sample Size Calculation

Presented by

Jim Matthews

Senior Consultant: Statistics

Sydney Informatics Hub

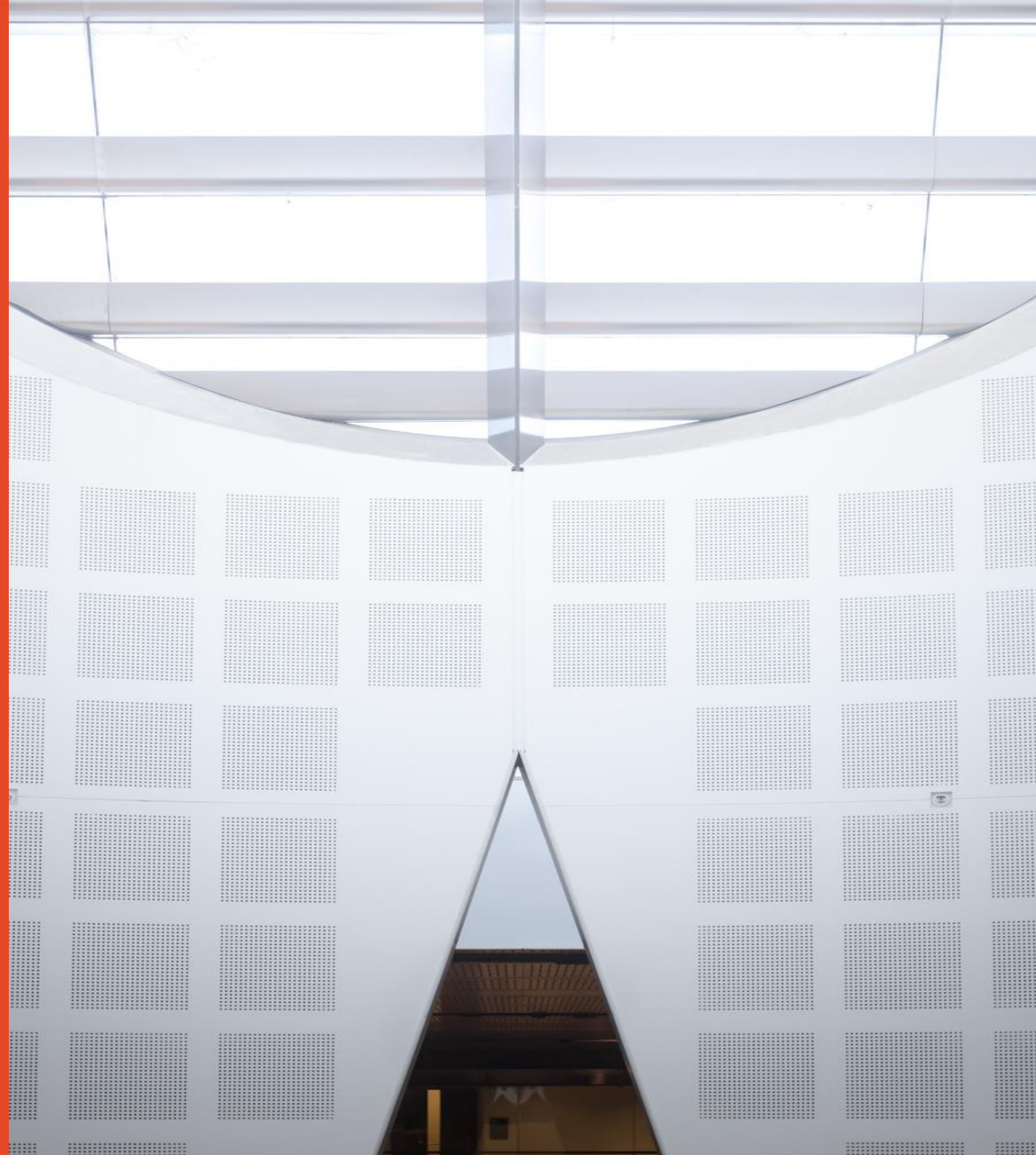
Core Research Facilities

The University of Sydney

August 2022



THE UNIVERSITY OF
SYDNEY



Outline

- Statistical power and sample size calculation - concepts
- Software tools - G*Power
- Example 1: Difference between 2 means (t-test)
- Example 2: Difference between 2 means (Mann-Whitney)
- Example 3: Difference between 2 proportions (z-test)
- Power calculation for other designs
- References

How to use this workshop

- These slides have a dual purpose:
 - To guide our interactive workshops
 - As self-contained reference material and workflows to be used after the workshop
- Some slides are for your reference, not all of the material will be discussed in the workshop. Such slides are marked with this blackboard icon
- Ask short questions or clarifications during the workshop. There will be breaks during the workshop for longer questions. You can email us about the material in these workshops at any time, or request a consultation for more in-depth discussion of the material as it relates to your specific project



Why do we need to calculate power and sample size?

Why do we want to estimate the power of an experiment?

- To know if it is worth doing the experiment
- To plan the time and resources necessary
- To make sure we are not wasting our time
- To get a grant application approved
- To make sure the study design is ethically acceptable

But do I really need to calculate power?

What type of study are you planning?

My study is:

- A pilot study
- Exploratory (no inferences or generalisation planned)
- Qualitative



NO – perhaps not



Sample size may be determined by other considerations, but a power analysis might still help

My study is:

- Confirmatory (pilot study already done)
- Testing a specific hypothesis
- Will make inferences about wider population



YES – Statistical validity is important



Continue with workshop!



What is the power of an experimental design?

The power to know...

Start with the hypothesis that you have generated, for example:

“The means of two groups are different”

In statistics, this is referred to as the alternative hypothesis H_1 .

Classically we test the veracity of the null hypothesis:

H_0 : There is no difference between the means of the two groups

A statistical test of the null hypothesis is always subject to uncertainty, or error. There are two main types of error.

Types of statistical error

Type I error

- Incorrectly rejecting the null hypothesis
- Also called false positive rate
- Referred to as the Significance level, designated by α
- The *convention* is to set the significance level to $\alpha = 0.05$

Type II error

- Incorrectly accepting the null hypothesis
- Also called the false negative rate
- Denoted by β
- Power is the complement of Type II error, denoted by $1 - \beta$
- We want Power to be as high as possible, typically $1 - \beta > 0.8$

Types of statistical error

When we perform a null hypothesis test, we are setting up a binary choice that can result in these types of error.



[This Photo](#) by Unknown Author is licensed under [CC BY-SA](#)

Types of statistical error

Table of error types		Reality Null hypothesis (H_0) is	
		True	False
<u>Decision</u> about null hypothesis (H_0)	Don't reject	Correct inference (true negative) (probability = $1-\alpha$)	Type II error (false negative) (probability = β)
	Reject	Type I error (false positive) (probability = α)	Correct inference (true positive) (probability = $1-\beta$)

Types of hypothesis

The most recognised hypothesis test relates to testing whether two measures are equal or different (a superiority trial).

Other study objectives will lead to other types of hypothesis test. The types below are frequently found in clinical trials:

- Superiority trials
- Equivalence trials
- Non-inferiority trials
- As-good-as-or-better trials
- Bioequivalence trials
- Trials to a given precision

The hypothesis tests that apply will vary depending on the study objective.

See reference for further details: Julious, Steven A. Sample Sizes for Clinical Trials . Boca Raton: CRC Press/Taylor & Francis, 2010. Print.

Hypothesis test or estimation of effect size?

What if we don't want to perform a hypothesis test?

What if we just want to estimate group means for example?

The same power calculation process can be applied.

We will consider why later.

Power calculation

How do we estimate the power of an experiment?

- It will depend on:
 - Sample size (more samples = more power)
 - Chosen significance level (typically $\alpha = 0.05$)
 - Minimum effect size to detect (larger minimum effect = more power)
 - Variance within groups (larger variance = less power)
 - Experimental design and type of statistical hypothesis test

Decisions regarding the experimental design can be critically important in determining statistical power.

This is covered in the “Experimental Design” workshop.

Sample size calculation - workflow

Often we need a sample size given a required minimum power

Sample size calculation workflow steps

1. Determine experimental design and statistical test
2. Set α and $1 - \beta$
3. Set the smallest effect size of interest
4. Estimate the variance
5. Calculate the minimum sample size
6. Explore scenarios

1. Determine experiment type and statistical test

For example:

Experimental Design	assumptions	proposed statistical test
Comparison of 2 means	independent groups, normality	Student's t-test
Comparison of 2 means	independent groups, no assumption of normality	Mann-Whitney U test
Comparison of 2 proportions	independent groups	z-test
Comparison of means, more than 2 groups	independent groups normality	ANOVA, F-test

2. Set α and $1 - \beta$

Setting values of parameters

- Typically choose $\alpha = 0.05$
- Typically choose $1 - \beta = 0.8$ (or higher)
- Sometimes power ($1 - \beta$) is required at 0.90 or 0.95

3. Set the smallest effect size of interest

What is the smallest effect size of interest?

- Decide on a smallest effect size of interest (sesoi). This should be based on the smallest effect size that is of *scientific interest*.

3. Set the smallest effect size of interest

too small

Effect size chosen is smaller than necessary



- The sample size is larger than necessary
- Possible waste of resources
- Can achieve statistical significance with an effect that is too small to be interesting or useful

just right

Effect size chosen is based on sesoi



- The sample size is just right
- If statistical significance is achieved, then it will align with scientific significance
- Most efficient use of resources

too large

Effect size chosen is larger than necessary

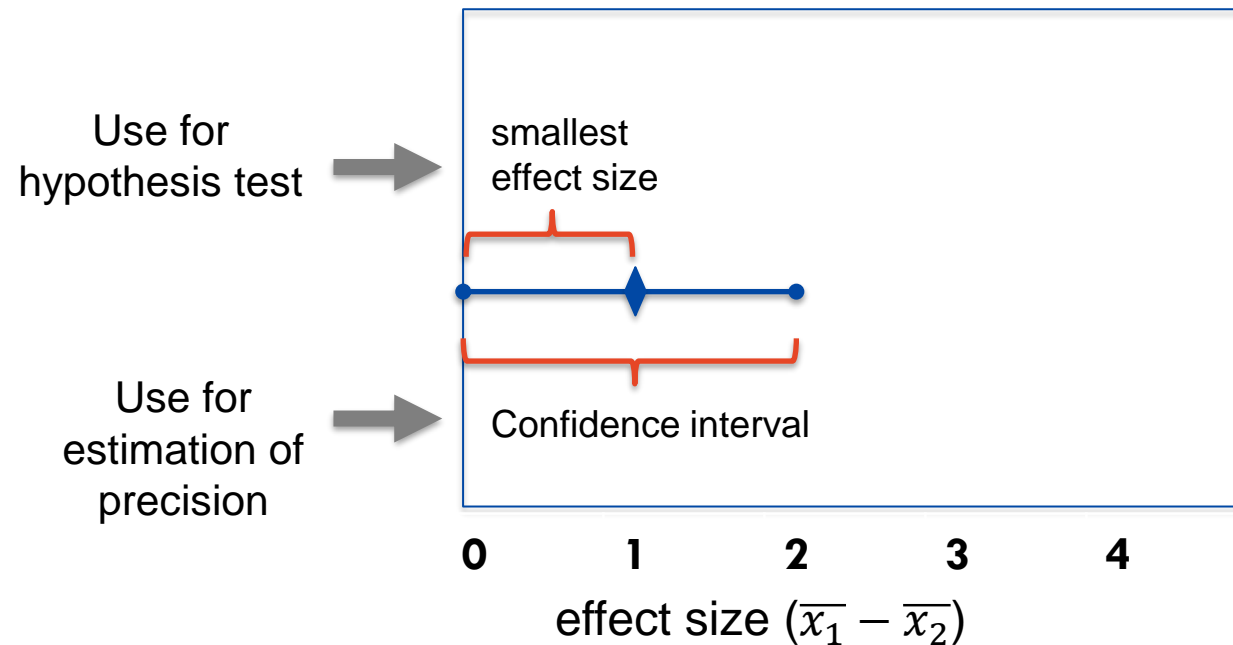


- The sample size is too small
- May detect large effects only
- Not able to achieve statistical significance for small effect sizes of interest
- Could be a waste of resources
- Will lead to a higher Type I error rate over the long run (poor reproducibility)

goldilocks

3. Set the smallest effect size of interest

Effect size – what it means for the hypothesis test and for the estimation of effect size



Further reading on the use of CI for sample size calc: see chapter 3 of **“Determining Sample Size Balancing Power, Precision, and Practicality”** by Dattalo

The minimum confidence interval width is twice the smallest effect size

4. Estimate the variance

Within study variance may be the big unknown in this calculation

How to estimate it?

- Estimate standard deviation (or proportions) from previous experiments?
- Consider theoretical bounds (eg for 5pt scales, proportions)
- Simulate some data and evaluate possible scenarios
- Seek expert knowledge?
- If no idea, may be best to do pilot study

4. Estimate the variance **Standardised Effect Size**

Alternative: Use the Standardised Effect Size

Many effect sizes can be “standardised” by considering the ratio of the effect size to a within group standard deviation.

For example: Cohen’s d is the ratio of the difference in means to the pooled standard deviation

$$d = \frac{\overline{x_1} - \overline{x_2}}{s}$$

Cohen’s d is therefore analogous to the number of standard deviations difference, or the z-score difference.

4. Estimate the variance **Standardised Effect Size**

Alternative: Use the Standardised Effect Size

Instead of deciding on effect size and an estimate of SD, we can choose a value of Cohen's d based on accepted interpretations of relative size.

<i>Effect size</i>	<i>d</i>	<i>Reference</i>
Very small	0.01	Sawilowsky, 2009
Small	0.20	Cohen, 1988
Medium	0.50	Cohen, 1988
Large	0.80	Cohen, 1988
Very large	1.20	Sawilowsky, 2009
Huge	2.0	Sawilowsky, 2009

Other guidelines are published for other standardised effect sizes.

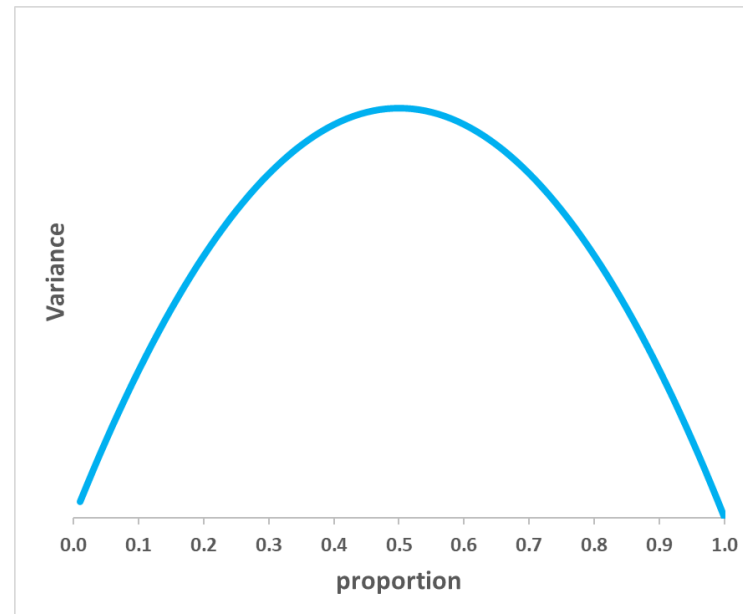
Note however that interpretation can vary across different fields of study.

4. Estimate the variance **Theoretical upper bound**

For proportions the maximum variance occurs when $p = 50\%$ and is at a minimum when $p = 0\%$ and 100% .

So we can use $p = 50\%$ to find a theoretical upper bound.

$$\text{Variance}(p) = p(1 - p)$$



4. Estimate the variance **Theoretical upper bound**

For ordinal responses such as 5pt scales a similar limit applies:

Possible responses are: 1, 2, 3, 4 or 5

Mean=3 Min=1 Max=5

$$\text{Variance}(5\text{pt scale}) = (\max - \text{mean})(\text{mean} - \min)$$

$$\text{Max Variance}(5\text{pt scale}) = (5 - 3)(3 - 1) = 4$$

In practice the actual variance will be smaller than the max. A rule of thumb is shown on StackExchange

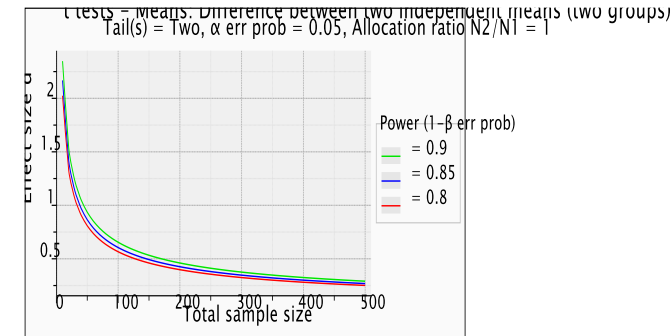
<https://stats.stackexchange.com/questions/23519/how-do-i-evaluate-standard-deviation>

5. Calculate the minimum sample size

- This is typically done using a software package (we will use **G*Power** in this workshop)
- Formulae for the calculation vary with the type of experimental design and the statistical test

6. Explore scenarios

- Don't just calculate a single sample size n !
- Use the software to calculate n for a range of scenarios in order to explore the consequences of uncertainty in the values used in the calculation
- This is called a **Power Analysis**
- *Consider also the shape of the cost curve for sample data collection*



For the above example note:
increasing sample size up to ~ 100
yields big effect size detection benefit,
but increasing sample size beyond
 ~ 100 yields diminishing returns.

Recap

Sample size calculation workflow steps

1. Determine experiment type and statistical test
2. Set α and $1 - \beta$
3. Set the smallest effect size of interest
4. Estimate the variance
5. Calculate the minimum sample size
6. Explore scenarios

Examples using G*Power software

We will work through 3 simple examples

1. Difference between 2 means (continuous response)
2. Difference between 2 means (survey response)
3. Difference between 2 proportions

Power calculation software

G*Power

- Download from website:
- <http://www.psychologie.hhu.de/arbeitsgruppen/allgemeine-psychologie-und-arbeitspsychologie/gpower.html>
- Current release 3.1.9.7 (Windows) 17 March 2020 (and 3.1.9.6 for Mac)
- Program has a simple user interface
- There is also a manual available online:
http://www.psychologie.hhu.de/fileadmin/redaktion/Fakultaeten/Mathematisch-Naturwissenschaftliche_Fakultaet/Psychologie/AAP/gpower/GPowerManual.pdf

G*Power 3.1.9.7

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 |p|: 0.3

α err prob: 0.05

Power ($1 - \beta$ 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

1. Difference between 2 means

Example: Chicken Welfare – Bone density

The bone density of chickens is an important indication of their welfare. We want to test to see if (mineral) bone density can be improved from 120 to at least 130 mg/cm³

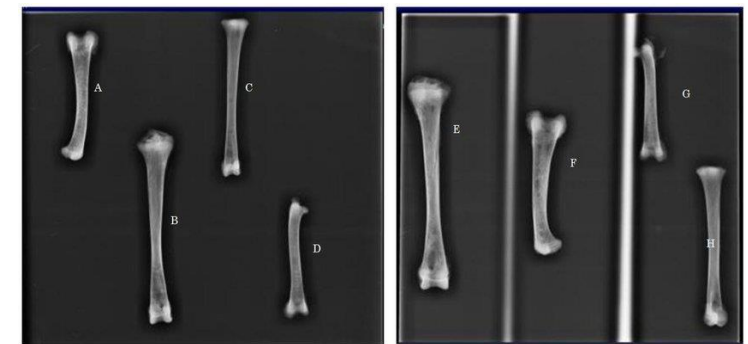


Treatment Group = high mineral diet

Control Group = normal diet

Response variable: Measure the tibia bone density after 6 weeks growth.
How many chickens do I need to detect a difference in bone density of 10 mg/cm³?

What type of statistical test will we perform?



TY - JOUR AU - Mabelebele, Monnye AU - Norris, Dannah AU - Siwendu, Ndyebo AU - Ng'ambi, Jones AU - John, Alabi AU - Mbajjorgu, C.A. PY - 2017/01/01 SP - 1387 EP - 1398 T1 - Bone morphometric parameters of the tibia and femur of indigenous and broiler chickens reared intensively VL - 15 DO - 10.15666/aeer/1504_13871398 JO - Applied Ecology and Environmental Research ER -

1. Difference between 2 means

Example: Chicken Welfare – Bone density

- Step 1: We will use a t-test (assume normality)
- Step 2: $\alpha=0.05$ and $1 - \beta=0.8$
- Step 3: Smallest Effect Size of interest is 10 mg/cm³
- Step 4: Estimate the variance
 - We know from previous studies what the typical variation in bone density is for the control diet. We don't know about the treatment diet. We will use an estimate from the control diet of SD=20 mg/cm³
- Assume we will have equal size groups, $n_1=n_2$

1. Difference between 2 means

Step 5: Calculate the minimum sample size

- Put all the information into G*Power
- Note: G*Power will convert the difference in means with the estimated SD to a standardized effect size called Cohen's d.

1. Difference between 2 means

Step 5: G*Power

G*Power will use this formula to calculate the sample size:

$$n = 2 \frac{\delta^2}{d^2}$$

where:

n = sample size per group (when $n_1 = n_2$)

δ = non-centrality parameter (of the t statistic, based on α & β)

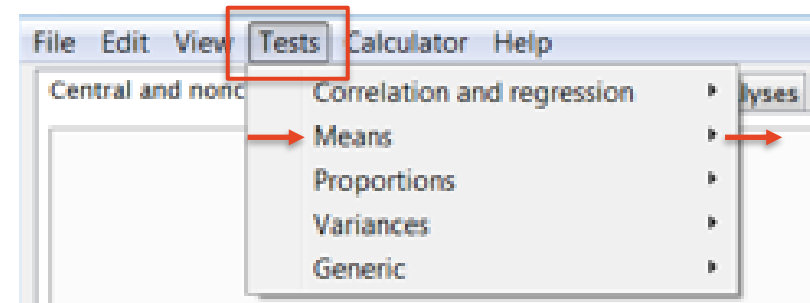
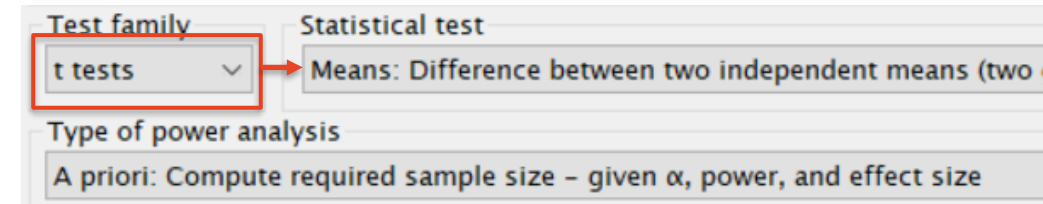
d = standardised effect size (Cohen's d)

1. Difference between 2 means

Step 5: G*Power

There are two ways to find the correct test

- Distribution approach: Select the test family (eg t tests), then the statistical test
- Design based approach: Select the test parameter class (eg means), then the study design
- Select Tests/Means/Two independent groups



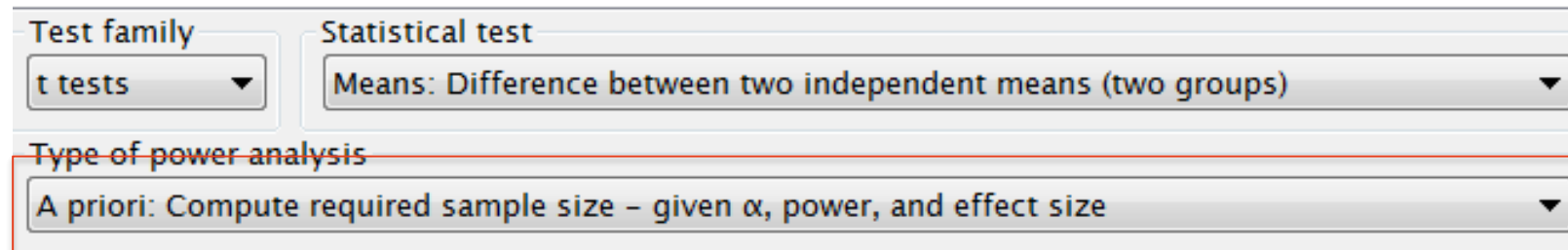
1. Difference between 2 means

G*Power

There are five different types of power analysis

- A priori
- Compromise
- Criterion
- Post Hoc
- Sensitivity

The “A priori” type is suitable for sample size calculation



The screenshot shows the G*Power software interface. It features three dropdown menus. The first, 'Test family', is set to 't tests'. The second, 'Statistical test', is set to 'Means: Difference between two independent means (two groups)'. The third, 'Type of power analysis', is set to 'A priori: Compute required sample size - given α , power, and effect size'. This third dropdown menu is highlighted with a red rectangular border.

1. Difference between 2 means

Example: Chicken Welfare – Bone density

Enter the values for the chick experiment

- Use $\alpha=0.05$ and $1 - \beta=0.8$
- Allocation ratio $N2/N1=1$
- Open the “determine” window to calculate the effect size d. Use means $M1=120$, $M2=130$, $SD1=SD2=20$, “calculate and transfer”
- Effect size is now shown $d=0.5$, select “two” tails, “Calculate”

The screenshot shows a software interface for calculating effect size and sample size. It is divided into two main sections: 'Input Parameters' and 'Output Parameters'.

Input Parameters:

- Tail(s): Two (selected from a dropdown)
- Effect size d: 0.5000000
- α err prob: 0.05
- Power ($1 - \beta$ err prob): 0.8
- Allocation ratio $N2/N1$: 1

Output Parameters:

- Noncentrality parameter δ : 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

Buttons and Actions:

- A red circle highlights the 'Determine =>' button in the Input Parameters section.
- A red circle highlights the 'Calculate and transfer to main window' button in the 'determine' window.
- Red arrows point from the 'Calculated results' text to the output parameters.

1. Difference between 2 means

Example: Chicken Welfare – Bone density

- Group sample sizes are $N1=64$, $N2=64$
- Actual power = 0.8015
- G*Power rounds up the sample size to the nearest integer, so actual power is slightly higher than the minimum requested.

Protocol of the power analysis

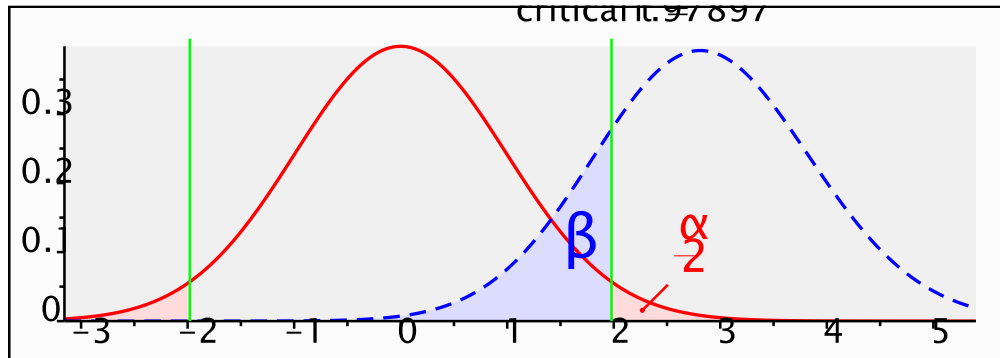
- You may want to save a copy of the calculation from this window (at the top right)

Central and non central distributions

- You may be interested to check the visual display of the test statistics in this window (at the top left)

1. Difference between 2 means

Central and non central test statistic distribution



The central distribution of a test statistic (in red) describes how a test statistic is distributed when the null hypothesis is true.

The non central distribution (blue dashed line) describes how the test statistic is distributed when the null hypothesis is false (alternate hypothesis is true).

Shows the distribution with the minimum effect size threshold.

1. Difference between 2 means

Example: Chicken Welfare – Bone density

Step 6: Explore scenarios

Power Analysis

- It is advisable to explore some different scenarios for different experimental settings.
- Consider how much your within study standard deviation could vary from your point estimate
 - Our estimate is $SD = 20$
 - Possible min value = 15 (optimistic)
 - Possible max value = 30 (pessimistic, conservative)

1. Difference between 2 means

Example: Chicken Welfare – Bone density

For G*Power we will use Cohen's d values to match the possible range of SD values

Min	SD = 15	$d = 10/15 = 0.67$
Expected	SD = 20	$d = 10/20 = 0.5$
Max	SD = 30	$d = 10/30 = 0.33$

1. Difference between 2 means

Example: Chicken Welfare – Bone density

X-Y Plot for a range of values

- Plot (on y axis) change to “power”
- Sample size from 10 to 400 in steps of 5
- Plot “3” graphs with $d = 0.33$ in steps of 0.17

Plot Parameters

Plot (on y axis) Power ($1 - \beta$ err prob) ☐ with markers

as a function of Total sample size from 10 in steps of 5 through to 400

Plot 3 graph(s) interpolating points

with Effect size d from 0.33 in steps of 0.17

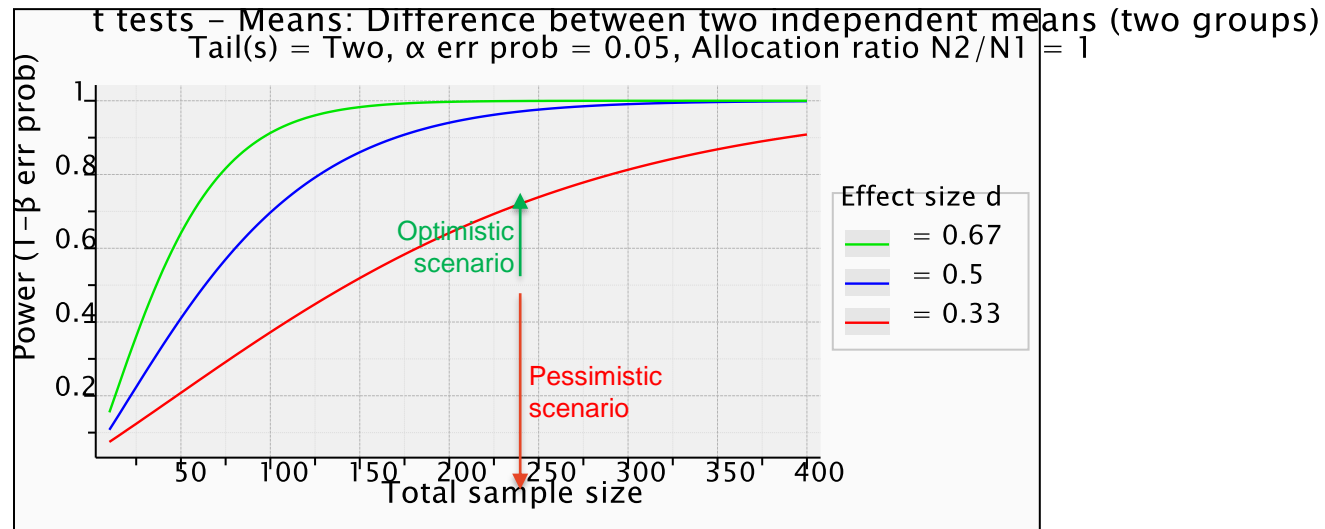
and α err prob at 0.05

Draw plot

1. Difference between 2 means

Example: Chicken Welfare – Bone density

X-Y Plot: sample size vs power



= 10/15

= 10/20

= 10/30

1. Difference between 2 means

Example: Chicken Welfare – Bone density

Remember: the accepted meaning of $d=0.5$ is that this is a “medium” standardised effect size, so our value of d is roughly in the right ballpark.

The sensitivity plot is another visualisation we can use in our power analysis. This plots effect size vs sample size.

1. Difference between 2 means

Example: Chicken Welfare – Bone density

Sensitivity Plot:

We want to look at a wide range of effect sizes. To do this, we will plot a sample size range from 10 up to 400 (as before) with 3 power curves for power = 0.8, 0.85, 0.90.

Plot Parameters

Plot (on y axis) ☐ with markers

as a function of from in steps of through to

Plot graph(s)

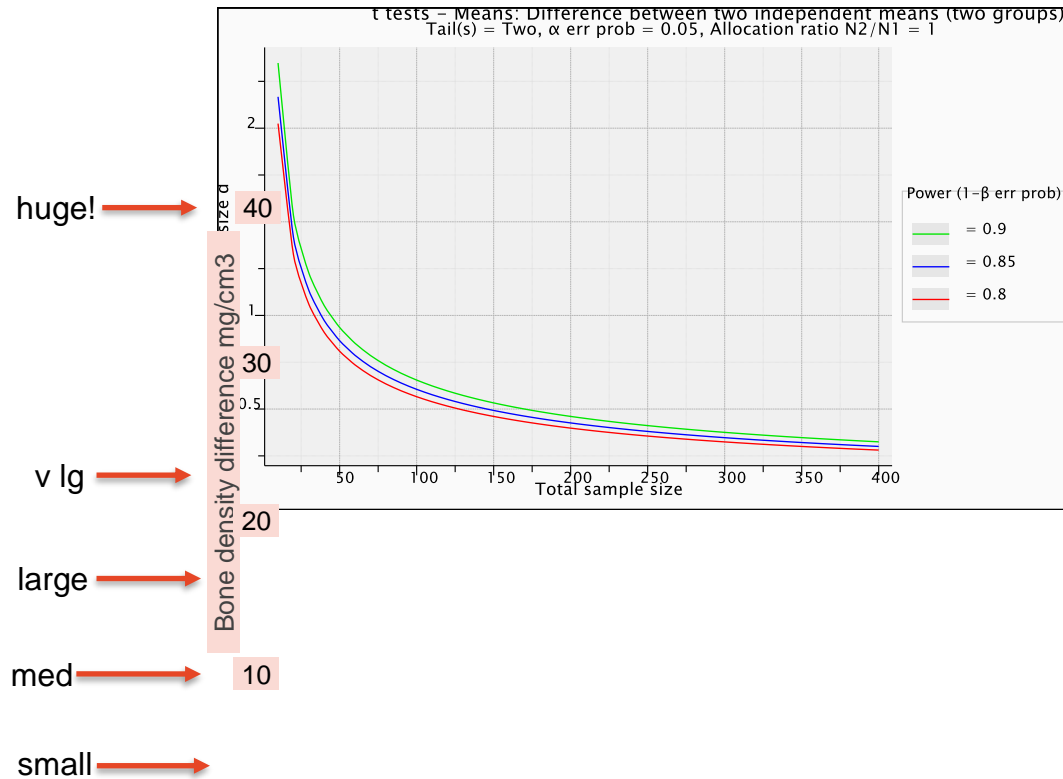
with from in steps of

and at

1. Difference between 2 means

Example: Chicken Welfare – Bone density

X-Y Plot: sample size vs effect size (sensitivity)



Effect size shown
assuming SD = 20

1. Difference between 2 means

Example: Chicken Welfare – Bone density

X-Y Plot: sample size vs effect size (sensitivity)

Customise plot in EXCEL

If you aren't happy with the G*Power plot, select the data from the Table tab and paste it into Excel (or your favourite plotting program).



GPower - Plot

File Edit View

Graph Table

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

Tail(s) = Two, α err prob = 0.05, Allocation ratio = 1:1

#	Total sample size	Power (1- β err prob) = 0.8	Power (1- β err prob) = 0.85	Power (1- β err prob) = 0.9
		Effect size d	Effect size d	Effect size d
1	10.0000	2.02444	2.16752	2.34795
2	20.0000	1.32495	1.41736	1.53369
3	30.0000	1.05980	1.13359	1.22644
4	40.0000	0.909129	0.972389	1.05199
5	50.0000	0.808708	0.864966	0.935757
6	60.0000	0.735621	0.786789	0.851171
7	70.0000	0.679351	0.726601	0.786054
8	80.0000	0.634299	0.678413	0.733919
9	90.0000	0.597169	0.638700	0.690955
10	100.000	0.565882	0.605236	0.654752
11	110.000	0.539050	0.576537	0.623705
12	120.000	0.515707	0.551570	0.596694
13	130.000	0.495156	0.529589	0.572915
14	140.000	0.476881	0.510044	0.551770
15	150.000	0.460492	0.492514	0.532806
16	160.000	0.445684	0.476677	0.515673
17	170.000	0.432219	0.462275	0.500093
18	180.000	0.419905	0.449105	0.485845
19	190.000	0.408587	0.437000	0.472750
20	200.000	0.398138	0.425824	0.460660
21	210.000	0.388452	0.415464	0.449452
22	220.000	0.379440	0.405825	0.439024
23	230.000	0.371027	0.396828	0.429291
24	240.000	0.363150	0.388403	0.420177
25	250.000	0.355755	0.380493	0.411620
26	260.000	0.348794	0.373048	0.403566
27	270.000	0.342226	0.366023	0.395966
28	280.000	0.336015	0.359381	0.388781
29	290.000	0.330131	0.353088	0.381973
30	300.000	0.324546	0.347114	0.375510
31	310.000	0.319235	0.341434	0.369365
32	320.000	0.314176	0.336023	0.363512
33	330.000	0.309351	0.330862	0.357929
34	340.000	0.304741	0.325932	0.352595
35	350.000	0.300331	0.321216	0.347493
36	360.000	0.296108	0.316698	0.342606
37	370.000	0.292058	0.312367	0.337920
38	380.000	0.288169	0.308208	0.333421
39	390.000	0.284432	0.304211	0.329097
40	400.000	0.280836	0.300365	0.324937

2. Difference between 2 means (Mann-Whitney)

The Mann-Whitney U test is a non-parametric version of the t-test for a difference in means. It is based on ranks. (also called Wilcoxon rank sum)

This is used when the data are not approximately normally distributed, or the underlying distribution is not normal.

Often used for ordinal data from surveys.

The values of the two groups are combined and ranked. The values are then divided back into the groups and the mean of the assigned ranks for each group is calculated and compared.

The test doesn't use the information about the size of the effect.

2. Difference between 2 means (Mann-Whitney)

Example: Happiness Survey



2

3



5

6



You want to measure happiness using the Lyubomirsky & Lepper scale. Each response ranges from 1 (unhappy) to 7 (happy). The score is the sum of 4 items, so the range is 4~28.

A pilot study on two groups produced the following results that can be used for the power calculation.

	Values		Ranks	
	Single	Married	Single	Married
	12	20	3	1
	11	15	4	2
	10	9	5	6
	6	8	8	7
Avg	9.8	13.0	5	4
SD	2.6	5.6		

2. Difference between 2 means (Mann-Whitney)

Example: Happiness Survey

You want to apply it to different groups of people (eg single vs married) to see if there is a difference in scores.

What is a meaningful difference?

Let's suppose that a minimum difference of 4 points (average of 1 pt difference per item) is the smallest effect size of interest.

2. Difference between 2 means (Mann-Whitney)

Example: Happiness Survey

So, what are our first 4 steps?

Step 1:	Determine experiment type and statistical test	Mann-Whitney
Step 2:	Set α and $1 - \beta$	0.05 and 0.8
Step 3:	Set the smallest effect size of interest	4 points
Step 4:	Estimate the variance	SD1=2.6, SD2=5.6

2. Difference between 2 means (Mann-Whitney)

Example: Happiness Survey

Sample size calculation

Heuristic method

“Do the calculations as if performing the corresponding parametric test (i.e. the t-test), then add 15% to the sample size.

2. Difference between 2 means (Mann-Whitney)

Example: Happiness Survey

- Tests>Means>Two Independent Groups
- Click “Determine” (different SDs so use $n1 = n2$)
- Enter expected means (use 9.5 and 13.5, equates to 4pt diff)
- Enter SDs from pilot study ($SD1 = 2.6$, $SD2 = 5.6$)

Test family: t tests
Statistical test: 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:
Tail(s): One
Effect size d: 0.9162174
 α err prob: 0.05
Power ($1 - \beta$ err prob): 0.95
Allocation ratio $N2/N1$: 1
Determine =>
Output Parameters:
Noncentrality parameter δ : ?
Critical t: ?
Df: ?
Sample size group 1: ?
Sample size group 2: ?
Total sample size: ?
Actual power: ?

☐ $n1 \neq n2$
Mean group 1: 0
Mean group 2: 1
SD σ within each group: 0.5
☒ $n1 = n2$
Mean group 1: 9.5
Mean group 2: 13.5
SD σ group 1: 2.6
SD σ group 2: 5.6
Calculate
Effect size d: 0.9162174
Calculate and transfer to main window
Close

2. Difference between 2 means (Mann-Whitney)

Example: Happiness Survey

- Check α , $1 - \beta$, two tails, allocation ratio=1.
- Calculate sample size. $N=20$ per group
- Add 15%. $N=20 \times 1.15 = 23$

Input Parameters		Output Parameters	
Determine =>	Tail(s) Two	Noncentrality parameter δ	2.8973338
Effect size d	0.9162174	Critical t	2.0243942
α err prob	0.05	Df	38
Power ($1 - \beta$ err prob)	0.8	Sample size group 1	20
Allocation ratio N2/N1	1	Sample size group 2	20
		Total sample size	40
		Actual power	0.8060552

2. Difference between 2 means (Mann-Whitney)

Theoretical approach

Statistical procedures can be compared according to their efficiency.

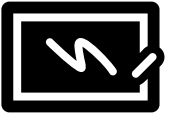
One test is more efficient than another if it requires fewer observations to obtain a given result.

The relative efficiency of two tests is the ratio of their efficiencies.

With smaller sample numbers, parametric tests are often more efficient than non-parametric tests although they approach equal efficiency with larger sample sizes.

The Asymptotic Relative Efficiency (ARE) is the limit of the relative efficiencies as the sample size increases. It can be calculated or set and is used in the sample size calculation, along with the effect size.

It can be shown that the minimum ARE for these two tests is 0.864.



2. Difference between 2 means (Mann-Whitney)

Example: Happiness Survey

- Under “Tests” select “Means” and then the option:
- “Two independent groups: Wilcoxon (non-parametric)”
- Use the same values as before:
- Two tails, $\alpha=0.05$ and Power=0.80, group means and SDs.
- Select Parent distribution = “min ARE”
- Calculate sample size >> N=23 per group

Input Parameters		Output Parameters	
Tail(s)	Two	Noncentrality parameter δ	2.8880475
Parent distribution	min ARE	Critical t	2.0248452
Determine =>	Effect size d	Df	37.7440000
	α err prob	Sample size group 1	23
	Power (1- β err prob)	Sample size group 2	23
	Allocation ratio N2/N1	Total sample size	46
		Actual power	0.8034207

3. Difference between 2 proportions

Example: Happiness survey

The survey scores could also be analysed as proportions by considering how many report a value above a threshold (say >14 means “happy”)

Singles group $P1$ = proportion of subjects respond “happy”

Married group $P2$ = proportion of subjects respond “happy”

Effect size: Say we want to find a minimum difference in proportions of $P1 - P2 = 0.1$ What sample size is required?

- Set $\alpha = 0.05$ and $1 - \beta = 0.8$, two tails
- Allocation ratio $N2/N1 = 1$
- We also need to estimate the two proportions. Let's first assume that there will be maximum variance ($p = 0.50$)
- Try using $P1 = 0.55$ and $P2 = 0.45$

3. Difference between 2 proportions

Example: Happiness survey

What are our first 4 steps this time?

Step 1:	Determine experiment type and statistical test	z-test for proportions
Step 2:	Set α and $1 - \beta$	0.05 and 0.8
Step 3:	Set the smallest effect size of interest	0.10
Step 4:	Estimate the variance	P1=0.55, P2=0.45

Note: The variance estimate comes from the proportion estimates.

$$\text{Variance} = p(1-p)$$

3. Difference between 2 proportions

Example: Happiness survey

We need 392 subjects per group to achieve Power=0.80

That's a lot of happy/unhappy people!

Test family		Statistical test	
z tests		Proportions: Difference between two independent proportions	
Type of power analysis			
A priori: Compute required sample size – given α , power, and effect size			
Input Parameters		Output Parameters	
Tail(s)	Two	Critical z	
Proportion p2	0.45	Sample size group 1	392
Proportion p1	0.55	Sample size group 2	392
α err prob	0.05	Total sample size	784
Power (1 – β err prob)	0.8	Actual power	0.8007410
Allocation ratio N2/N1	1		

3. Difference between 2 proportions

Example: Happiness survey

Step 6: Suppose the proportion of subjects responding “happy” is expected to be higher, around 90%.

Try using $P_1=0.85$
and $P_2=0.95$

The screenshot shows a statistical software interface for a two-proportion z-test. The 'Test family' is set to 'z tests'. The 'Statistical test' is 'Proportions: Difference between two independent proportions'. The 'Type of power analysis' is 'A priori: Compute required sample size - given α , power, and effect size'. The 'Input Parameters' section includes: 'Tail(s)' set to 'Two', 'Proportion p2' set to 0.95, 'Proportion p1' set to 0.85, ' α err prob' set to 0.05, 'Power ($1-\beta$ err prob)' set to 0.8, and 'Allocation ratio N2/N1' set to 1. The 'Output Parameters' section includes: 'Critical z' set to 1.9599640, 'Sample size group 1' set to 141, 'Sample size group 2' set to 141, 'Total sample size' set to 282, and 'Actual power' set to 0.8025450.

Input Parameters		Output Parameters	
Tail(s)	Two	Critical z	1.9599640
Proportion p2	0.95	Sample size group 1	141
Proportion p1	0.85	Sample size group 2	141
α err prob	0.05	Total sample size	282
Power ($1-\beta$ err prob)	0.8	Actual power	0.8025450
Allocation ratio N2/N1	1		

Now we only need 141 subjects per group

Note the difference in sample sizes corresponding to the different proportion estimates. Remember the variance of the proportion parameter [$\text{var}=p(1-p)$] is at a maximum at 0.5 and gets smaller close to zero and one.

3. Difference between 2 proportions

G*Power provides a total of 4 options for power calculations for proportions with independent groups:

- Inequality, z-test (used in Happiness intervention example)
- Inequality, Fisher's Exact test
- Inequality, Unconditional exact
- Inequality with offset, Unconditional exact

The Fisher's Exact test should be used when sample sizes are going to be small (say $n_1p_1 \leq 5$ or $n_2p_2 \leq 5$)

– Let's try the Fisher's Exact for the Happiness example

3. Difference between 2 proportions

Example: Happiness survey

Step 6: Use the Fisher's Exact test to get the sample size with

$P1=0.85$ and $P2=0.95$

Fisher's Exact suggests 151 subjects per group.

Not quite the same result as the z-test, but note that the actual alpha is 0.024 rather than 0.05.

Test family		Statistical test	
Exact		Proportions: Inequality, two independent groups (Fisher's exact test)	
Type of power analysis			
A priori: Compute required sample size - given α , power, and effect size			
Input Parameters			
Determine =>	Tail(s)	Two	
	Proportion p1	0.85	
	Proportion p2	0.95	
	α err prob	0.05	
	Power ($1-\beta$ err prob)	0.8	
	Allocation ratio N2/N1	1	
Output Parameters			
Sample size group 1		151	
Sample size group 2		151	
Total sample size		302	
Actual power		0.8005824	
Actual α		0.0243675	

3. Difference between 2 proportions

Example: Happiness survey

Step 6: Explore scenarios

- When considering various scenarios, look for value estimates that provide a conservative power estimate.
- In this example proportions centred around 0.5 represent the most conservative estimate. This gives the largest sample size estimate.
- This principle may also be applied to the study design as well. For example powering your study for a non-parametric test is conservative (Mann-Whitney instead of t-test).

Power Analysis for other designs

G*Power scope

G*Power includes methods to calculate power and sample size for a wide variety of design scenarios, eg

- ANOVA
- Correlation
- Linear Regression
- Logistic Regression

Refer to the manual for details

Effect Sizes for other designs

Effect size for ANOVA

G*Power uses the standardised effect size; Cohen's f is related to the partial eta squared

$$\eta^2 = \frac{f^2}{(1 + f^2)}$$

Partial eta squared is often reported in the ANOVA table output

Effect size for other designs: Use a wide variety of effect size measures

Power Analysis for other designs

From simple designs to complex designs

So far we have considered power analysis for simple designs where the mathematical calculations are tractable and rely on a limited set of assumptions regarding the data to be obtained.

As design complexity increases, it becomes more difficult or perhaps impossible to find an analytical solution to calculate power.

When no formula exists:

- First option – determine sample size for a simplified version of the study design and extrapolate this to the more complex design
- Second option – Use a simulation method (that does not rely on formulae)

Power Analysis for other designs

G*Power limitations

G*Power does not do everything!

Use G*Power for simple to moderately complex designs including where simplification of the design can yield an approximate solution.

Switch to simulation methods for complex study designs where analysis of a simplified design is not sufficiently rigorous.

Power Analysis – by simulation

Simulation based power estimation

- Simulate (many) data sets
- Analyse each data set and test for statistical significance
- Calculate the proportion of significant p values

$$Power = \frac{\text{significant simulations}}{\text{all simulations}}$$

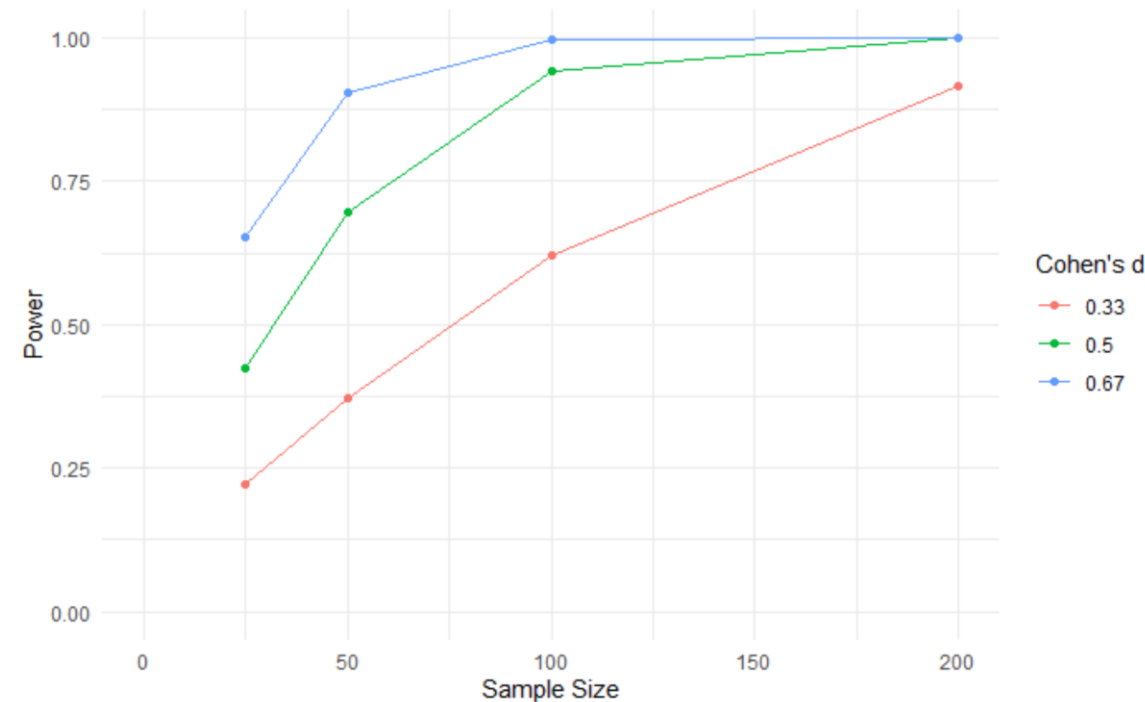
- The ‘trick’ is to set the parameters of the simulation in a sensible, realistic way

<https://link.springer.com/article/10.3758/s13428-021-01546-0>

Power Analysis – by simulation

Example 1: Chicken Welfare - bone density (difference between 2 means)

- Simulation in R using package “paramtest”
- Results for this simple simulation will be very similar to those obtained from G*Power.
- See R Markdown files for details



Software for Power Analysis

Free and Open Source software

- R /R Studio:
 - Base R has functions covering basic proportions, t-tests, etc.
 - Package “pwr” has 9 functions covering proportions, t-tests, ANOVA, chi-square and correlations
 - Package “epiR” has 23 functions covering many statistics including AUC, sensitivity and specificity
 - Package “paramtest” basic power calculations by simulation
 - Package “mixedpower” for generalised linear mixed models
 - Package “simr” simulation based power calculations for mixed models
- Online calculators such as www.powerandsamplesize.com and <http://sample-size.net/>
- G*Power is a dedicated (free) program
- Make your own in Excel! (for example see Lakens, D. (2013). Calculating and reporting effect sizes to facilitate cumulative science: A practical primer for t-tests and ANOVAs. Frontiers in Psychology, 4:863. doi:10.3389/fpsyg.2013.00863)

Software for Power Analysis

Proprietary \$\$ software

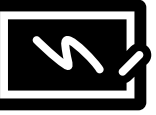
- Packages such as STATA, SPSS and SAS include a calculator
- GraphPad have “StatMate” separate to Prism
- PASS by NCSS dedicated software esp. for medical research



Power calculation references

- **G*Power** <http://www.psychologie.hhu.de/arbeitsgruppen/allgemeine-psychologie-und-arbeitspsychologie/gpower.html>
- **NCSS PASS Statistical software** <https://www.ncss.com/software/pass/>
- **Causal Evaluation** <https://www.causalevaluation.org/power-analysis.html>
- **Epi Tools for disease prevalence (by AUSVET)**
<http://epitools.ausvet.com.au/content.php?page=SampleSize>
- **Demidenko (Dartmouth) for logistic regression**
<https://www.dartmouth.edu/~eugened/power-samplesize.php>
- **National Institutes of Health (NIH – USA) for cluster randomised trials**
<https://researchmethodsresources.nih.gov/SampleSizeCalculator.aspx>
- **UCSF Clinical and Translational science institute** (Survival for clinical research) <http://www.sample-size.net/sample-size-survival-analysis/>
- **Lakens, D. Open Science Framework** <https://osf.io/ixGcd/>

Power Analysis – library references



- **Cohen, Jacob. Statistical Power Analysis for the Behavioral Sciences.**
Burlington: Elsevier Science, 2013. Print.
https://sydney.primo.exlibrisgroup.com/permalink/61USYD_INST/14vvljs/alma991005702359705106
- **Dattalo, Patrick. Determining Sample Size Balancing Power, Precision, and Practicality**
Oxford: Oxford University Press, 2008. Print.
https://sydney.primo.exlibrisgroup.com/permalink/61USYD_INST/14vvljs/alma991015395569705106
- **Julious, Steven A. Sample Sizes for Clinical Trials**
Boca Raton: CRC Press/Taylor & Francis, 2010. Print.
https://sydney.primo.exlibrisgroup.com/permalink/61USYD_INST/14vvljs/alma991000960739705106
- **Ryan, Thomas P., and Thomas P Ryan. Sample Size Determination and Power.**
Somerset: John Wiley & Sons, Incorporated, 2013. Web.
https://sydney.primo.exlibrisgroup.com/permalink/61USYD_INST/1367smt/cdi_askewsholts_vlebooks_9781118439203



Asking for more help

SIH Resources

- Our website

www.sydney.edu.au/research/facilities/sydney-informatics-hub.html

or Google “Sydney Informatics Hub”

- SIH training

- Sign up to our mailing list to be notified of upcoming training:

<https://signup.e2ma.net/signup/1945889/1928048/>

- Hacky Hour www.sydney.edu.au/research/facilities/sydney-informatics-hub/workshops-and-training/hacky-hour.html

or Google “Sydney Hacky Hour”

Other resources

- OLE courses

- Other University Accessible

- LinkedIn Learning [[The Data Science of Experimental Design](#) see Ch 6]

Acknowledging SIH



All University of Sydney resources are available to Sydney researchers **free of charge**. The use of the SIH services including the Artemis HPC and associated support and training warrants acknowledgement in any publications, conference proceedings or posters describing work facilitated by these services.

The continued acknowledgment of the use of SIH facilities ensures the sustainability of our services.

Suggested wording:

General acknowledgement:

"The authors acknowledge the technical assistance provided by the Sydney Informatics Hub, a Core Research Facility of the University of Sydney."

Acknowledging specific staff:

"The authors acknowledge the technical assistance of (name of staff) of the Sydney Informatics Hub, a Core Research Facility of the University of Sydney."

For further information about acknowledging the Sydney Informatics Hub, please contact us at sih.info@sydney.edu.au.

End of Workshop

- Thank you for your interest and attention
- Questions and comments welcome
- We appreciate your feedback via the on-line survey

- **Jim Matthews** BEng MStat | Senior Consultant: Statistics
- The University of Sydney
- Sydney Informatics Hub | Core Research Facilities
- Rm 386 Merewether Building (H04) | The University of Sydney | NSW | 2006
- +61 412 246 271
- Jim.Matthews@sydney.edu.au | sydney.edu.au

