Welcome to the course on

Experimental Design and Statistics: The Good, the Bad and the Ugly

Manuel Berdoy

Make sure you have downloaded the app on your smart phone (IOS or Android) from: get.meetoo.io Whilst we wait to start you can enter the meeting Our Meeting ID is: 107-854-181



By continuing, I accept the <u>Terms & Privacy Policy</u>

Experimental Design and Statistics.

- Aim: understand the Essence of Experimental Design and Statistics
- Ethos for today (important): we are all <u>interacting friends</u>.
- Material for distribution throughout the day
- Group exercises

Experimental Design and Statistics.

Your Handouts:

- 1. From P values to Power Calculation (LECTURE NOTES)
- 2. Analysis of Variance and the Control of Variation. (LECTURE NOTES)
- 3. Workshop exercises
- 3b, 3c papers, as part of workshop exercises
- 4. Assumption of Parametric tests, and data transformation
- 5. + 5b Overview of Several Experimental designs and analyses
- 6. Some relevant websites (mentioned during the lectures)
- 7. Adjusted vs Sequential Sum of Squares (additional material)



Manuel 's Experimental Design

Meeting ID:

107-854-181

Document 1: From P values to Power calculation Why Statistics ?

"Lies, Damned lies and Statistics". This is not "Statistics", but politics.

"All statistics shows is that most of us have more than the average number of legs"

Factually correct, but that is also not what

"statistics" is about.

Statistics is about finding the <u>truth</u> in an <u>imperfect</u> world, and saving time and effort... and (sometimes) lives.

Searching for Truth = Experiment Test of drug doses on cognitive abilities (tested in maze). Exp 1: 3 treatments administered to N=18 male rats, group-housed in 3 cages of 6 rats). Control administered to cage 1 (6 rats), D1 to cage 2,etc.

Conclusion (kind of encouraging) D1 no effect. D2 significant effect t-tests P=0.055 (ns)

Time spent in maze





Searching for Truth = Experiment ?

NEWS

EDISON'S BULBSFAIL TO LIGHT UP AUCTION First all-science collection sells modestly a Christied. www.auturc.com/tews

Animal experiments under fire for poor design

In the contentious world of animal research, one question surfaces time and again: how useful are animal experiments as a way to prepare for trials of medical treatments in humans? The issue is crucial, as public opinion is behind animal research only if it helps develop better drugs. Consequently, scientists defending animal experiments insist they are essential for safe clinical trials, whereas animal-rights activists vehemently maintain that they are useless.

Now a British team has made the first attempt to answer the question in a scientific way, and the result suggests that animal researchers need to raise their game. The team claims that animal experiments are often poorly designed, and so fail to lay the ground properly for subsequent human studies.

The study looked at six treatments that have been evaluated in detail in human trials. The researchers assessed whether animal studies had accurately predicted the outcome of the human work, a task that involved reviewing more than 200 papers. In three of the six cases, the answer was no (P. Perel *et al. Br. Med.*). doi:10.1136/bmj.39048.49728.BE;2006).



Are animals being wasted in badly thought through experiments?

e.g. lack of reporting of:

Randomisation, Double blinding,

External Validity

+ other problems...

The problem: Basic errors are widespread

- <u>experimental unit</u> wrongly identified in 48% of 99 papers surveyed (Hurlbert 1984)
- 54% of 141 articles in *Infection* and *Immunity* had <u>errors of analysis</u>, reporting or both (Olsen 2003)
- 79 of 157 neuroscience papers used <u>incorrect comparisons</u> of results (Neuwenhuis et al 2009 -)

 "<u>Random allocation</u> of animals to experimental groups was reported in only 13% of all the studies in the sample" (Kilkenny et al 2009)
 Only 14% of all papers [susceptible to observer bias] also reported that they used <u>blinding</u>." (Blinding is an effective way of reducing bias)

We know that these errors have an effect:

- 11 publications, 29 experiments, 408 animals
- Improved outcome by 44% (35-53%)



We know that these errors have an effect:

BIAS

- Of 290 studies surveyed those which did not report that they <u>randomised</u> were more likely to report positive findings than those that did.
- Those that reported neither <u>randomising nor blinding</u> were even more likely to report positive findings. (Bebarta et al 2003)

WASTE of ANIMALS:

Many researchers do not use efficient "blocked" and "factorial" designs



Fang & Casadevall 2011

Retraction Watch: Tracking retractions as a window into the scientific process https://retractionwatch.com/

ANNOUNCEMENT

Reducing our irreproducibility

Over the past year, *Nature* has published a string of articles that highlight failures in the reliability and reproducibility of published research (collected and freely available at go.nature.com/ huhbyr). The problems arise in laboratories, but journals such as this one compound them when they fail to exert sufficient scrutiny over the results that they publish, and when they do not publish enough information for other researchers to assess results properly.

From next month, *Nature* and the Nature research journals will introduce editorial measures to address the problem by improving the consistency and quality of reporting in life-sciences articles. To ease the interpretation and improve the reliability of published results we will more systematically ensure that key methodological details are reported, and we will give more space to methods sections. We will examine statistics more closely and encourage authors to be transparent, for example by including their raw data.

Central to this initiative is a checklist intended to prompt authors to disclose technical and statistical information in their submissions, and to encourage referees to consider aspects important for research reproducibility (go.nature.com/oloeip). It was developed after discussions with researchers on the problems that lead to irreproducibility, including workshops organized last year by US National Institutes of Health (NIH) institutes. It also draws on pubWe know that these errors have an effect: LACK OF EXTERNAL VALIDITY (no agreement in 50% of cases)

The cost of of low reproducibility in preclinical research (i.e. not just animal research) in the US alone, is estimated at:... \$ 28 000 000 000

(Freedman, Cockburn & Simcoe PLOS Biology | DOI:10.1371/, 2015)

The Economics of Reproducibility in Preclinical Research

Leonard P. Freedman 🖾, Iain M. Cockburn, Timothy S. Simcoe

Published: June 9, 2015 • DOI: 10.1371/journal.pbio.1002165

Fig 2. Estimated US preclinical research spend and categories of errors that contribute to irreproducibility.



Freedman LP, Cockburn IM, Simcoe TS (2015) The Economics of Reproducibility in Preclinical Research. PLoS Biol 13(6): e1002165. doi:10.1371/journal.pbio.1002165

DATA PIPELINE

The design and analysis of a successful study has many stages, all of which need policing.

Data collection

Experimental design

Rat given no drug (placebo - control) Rat given a drug Does the drug have an effect ? i.e. is there a difference between the 2 groups ?





In a perfect world...2 would be enough

N = 2



WEIGHT

Number of rats

In the real world...: variation (i.e. noise)

N = 16



WEIGHT

N = 90





How confident are we that there is a difference? i.e. How much would you be prepared to bet ?: POLLOPEN





Number of rats



Presentation reminder



Number of rats

Presentation reminder



(no drug)

Treatment (drug) Searching for Truth = Experiment Test of drug doses on cognitive abilities (tested in maze). Exp 1: 3 treatments administered to N=18 male rats, group-housed in 3 cages of 6 rats). Control administered to cage 1 (6 rats), D1 to cage 2,etc.





43.59%

48.72%

When comparing means I prefer to use:

1. Standard Deviation

2. Standard Error

- 3. Whatever looks smaller 2.56%
- 4. Whatever the software gives me 0%
- 5. It does not matter (they are related) 0%
- 6. I don't really know what they are 5.13%

Relationship between SE and SD



You need to write this down

Power Calculation

- Not always realistic
- And we do it to some extent anyway
 - Power calculation is a way of formalising the process
 - Ethical grounds & scientific grounds for it. (Regulatory, Grant Awarding bodies)

Ethical + Scientific Issues

Ethical, using a medical example: If test of new drug will have adequate power with a sample of 100 patients, then inappropriate to use 200. (and with animals, non consenting and possibly painful.

Scientific + ethical:

Conversely,

If test of new drug requires 200 patients to yield adequate power, then inappropriate to use 100. Patients accept to be part of the study on the assumption that it will yield useful results. (animal equivalent is that no result due to lack of power is a waste of animals).

Some Reminders

HO= Nul hypothesis = Nil hypothesis = no effect

"State of Nature"	Reject HO Accept HO (Find effect) (Find no effect)
No effect (HO is true)	Type I error CORRECT Alpha
Effect (HO is false)	p value The prob. that exp. will give a false positive result (e.g.due to random fluctuations) IT DOES HAPPEN

Some Reminders

HO= Nul hypothesis = Nil hypothesis = no effect

"State of Nature"	Reject HO (Find effect)	Accept HO (Find no effect)
No effect (HO is true)	Type I error	CORRECT
	Alpha p value	"Absence of evidence is no evidence of absence"
Effect (HO is false)	CORRECT	Type II error

Some Reminders

HO= Nul hypothesis = Nil hypothesis = no effect

"State of Nature"	Reject HOAccept HO(Find effect)(Find no effect)
No effect (HO is true)	Type I error CORRECT Alpha Prob. of detecting a specified effect
	p value at specified significance level
Effect (HO is false)	CORRECT Type II error (1-Beta – Power) Beta

The 6 variables "determining" the chance of statistical significance

Significance level [arbitrary, set at P= 0.05 min]
Desired Power of experiment [arbitrary, set at 0.80 - 0.90]
Alternative Hypothesis (1 vs 2 tailed)

- Size of the effect of biological interest
- Variation (i.e. Standard Deviation)
- Sample size (N)

Note: this is a closed system, i.e. fix any five and the sixth can be derived

(= SIGNAL) (= NOISE)

= about False Positive

Statistical Inference in the 21st Century: A World Beyond p < 0.05 The American Statistician ; Volume 73 2019 - Issue sup1:

11:12
"Out of the four studies, half reported no significant difference therefore......"



No effect Observed effect and confidence



Statistical Inference in the 21st Century: A World Beyond p < 0.05 The American Statistician ; Volume 73 2019 - Issue sup1:

P-Hacking; HARKing"0.05 cliff"

The 6 variables "determining" the chance of statistical significance

 Significance level = a [arbitrary, set at P= 0.05 min]
Desired Power of experiment = a [arbitrary, set at 0.80 - 0.90]
Alternative Hypothesis (1 vs 2 tailed)

- Size of the effect of biological interest
- Variation (i.e. Standard Deviation)
- Sample size (N)

Note: this is a closed system, i.e. fix any five and the sixth can be derived

= about False Positive

= about False <u>Negative</u>

(= SIGNAL) (= NOISE)





How many animals do we need ? (to have enough power)

How many animals are needed?







Power calculation, using Power and Precision software (they are many others)

+ GPOWER

Predicted effect = 10 units "Noise" = SD = 10 units

How many animals do we need ? (to have enough power)





1000 hypotheses tested. 100 are true. How many significant results expected with p< 0.05, 80% power. 75

43.33%



What we find

POLL OPEN

	Effect	No effect
No Effect	Type 1 error $(\alpha = 0.05)$	
Effect	(1-β= 0.80)	Type 2 error (β)
	No Effect Effect	EffectNoType 1 errorEffect $(\alpha = 0.05)$ Effect $(1-\beta=0.80)$

16.67%

1000 hypotheses tested. 100 are true. How many significant results expected with p< 0.05, 80% power.

900

100

1000

What we find

		Effect	No effect
ţ	No Effect	Type 1 error $(\alpha = 0.05)$	(1-a = 0.95)
Γru		$900 \ge 0.05 = 45$	900 x 0.95 = 855
he -	Effect	$1 - \beta = 0.80$	Type 2 error $(\beta = 0.20)$
H		$100 \ge 0.8 = 80$	$100 \ge 0.2 = 20$
		45 + 80 = 125 (25% overestimate)	855 + 20 = 875 (3% underestimate)





1000 hypotheses tested. 100 are true. How many significant results expected with p< 0.05, **20% power.**

900

10

100

What we find

			Effect	No effect
)	ţ	No Effect	Type 1 error $(\alpha = 0.05)$	$(1-\alpha = 0.95)$
	Tru		$900 \ge 0.05 = 45$	900 x 0.95 = 855
	່ອ	Effect	$1-\beta = 0.80-0.20$	Type 2 error ($\beta = 0.80$)
0	È		$100 \ge 0.8 = 20$	$100 \ge 0.2 = 80$
0		()	45 + 20 = 65 35% underestimate)	855 + 20 = 935 (4% overestimate)
	Bl	JT most (69	%) of significant results ar	e wrong (45 out of 65).





Sources of Noise

Age, Sex, Weight Stress Subclinical disease Temperature etc Animal House/Barn/Cage Position in rack

Time (hours/months) Experimenter/carer (competence, unintentional bias)

The experimental

-unit?



Look at the data !

. . .

This the topic of the next section...

Experimental Designs & Analyses

Example of designs:

Factorial

Randomised Blocks

- •Latin Square
- Cross-Over
- Repeated Measure
- •Covariance

Two common mistakes in biomedical-literature:

- <u>multiple</u> t-tests
- no blocking



N= 27 rats, assigned to 3 levels of treatment

ANOVA table

Source	DF	SS	MS	F	P
Treatment	2	2861	1430.5	13.13	0.00014
Error	24	2614	108.9		
Total	26	5475			



 \succ

=overall mean

Source DF SS MS F Ρ 2 2861 1430.5 13.13 0.00014 Treatment Error 24 2614 108.9 Total 26 5475



=overall mear	1
---------------	---

Source	DF	SS	MS	F	P
Treatment	2	2861	1430.5	13.13	0.00014
Error	24	2614	108.9		
Total	26	5475			



Why Squares?

1. To remove negative signs (or they add up to zero)

62.96%

- 2. To give more weight to outliers
- 3. Because that is best way to partition variation

11.11%

4. Because it works

7.41%

5. Dont know

3.7%



<= total amout of varation around the mean











=overall mear	1
---------------	---

Source	DF	SS	MS	F	P
Treatment	2	2861	1430.5	13.13	0.00014
Error	24	2614	108.9		
Total	26	5475			





Source	DF	SS	MS	F	P
Treatment	2	2861	1430.5	13.13	0.00014
Error	24	2614	108.9		
Total	26	5475			

rat size shown

Blocking (i.e. controlling) for size

Big rats "block"Medium ratsSmall ratsRat Rat Rat Rat...Rat Rat Rat...Rat Rat Rat...

Blocking (i.e. controlling) for size

Big rats "block" N=9 Rat Rat Rat Treatment 1 Treatment 2 Treatment 3 Medium rats Rat Rat Rat

Small rats Rat Rat Rat

Blocking (i.e. controlling) for size

Big rats "block" N=9 Rat Rat Rat Treatment 1 Treatment 2 Treatment 3 Medium ratsSRat Rat RatF

Small rats Rat Rat Rat

Etc...



	Source	DF	SS	MS	F	P
	Treatment	2	2861	1430.5	13.13	0.00014
	Error	24	2614	108.9		
	Total	26	5475			
	Source	DF	SS	MS	F	P
WITH BLOCKING	Block(=size)	2	2279	1139.7	74.94	0.0000
	Treatment	2	2861	1430.5	94.06	0.0000
	Error	22	334	15.2		
	Total	26	5475			

ТЗ







	Source	DF	SS	MS	F	P
	Treatment	2	2861	1430.5	13.13	0.00014
	Error	24	2614	108.9		
	Total	26	5475			
	Source	DF	SS	MS	F	P
WITH BLOCKING	Block(=size)	2	2279	1139.7	74.94	0.0000
	Treatment	2	2861	1430.5	94.06	0.0000
	Error	22	334	15.2		
	Total	26	5475			

df SS MS F pSex(overall diff. between the sexes ?)Drug(overall diff. between the drug treatments ?)Sex*DrugInteraction ?





df SS MS F pSex(overall diff. between the sexes)Drug(overall diff. between the drug treatments)Sex*DrugInteraction





POLL OPEN

91.3%



2. NO

- 3. Not sure 0%
- Panic: 4. 0%
- Please explain the principle again 4.35% 5.





POLL OPEN

93.75%

- $1. \quad \underline{YFS}$
- 2. **NO**3.13%
- 3. Not sure 0%
- 4. *Panic:*
- 5. Please explain the principle again 3.13%





84.62%

1. YES

- 2. NO 7.69%
- 3. Not sure 7.69%
- 4. *Panic:* 0%
- 5. Please explain the principle again 0%
Multiple comparisons



Variety of tests. Most common are:

1. Comparison of <u>selected</u> pairs of mean: Bonferroni tests; (equivalent of multiple t tests with correction for multiplicity) but harsh (low power); not recommended for 5 groups or more.

2. One group (e.g control) against all the others= Dunnett's test

3. Compairing means of <u>preselected</u> groups A <u>&</u> B vs C<u>&D&E</u> = Contrast

4. All pairs of means = Tukey's or Student-Newman's test (Roughly the same: Tukey said to be more conservative ie more false negatives (Type II error), and SN more false positives (Type I error)

Statistics, Experimental Design, and Animal Experimentation.

- Searching for Truth
- Power (Calculations)
- Exp. Design & Analyses: controlling variation
- Quick Recap
- Refinement vs Reduction

Thinking (List) Exercises

Document 1

Final look at 1. Test of two drugs on maze ability









In Short: Take home messages

- Remember the law of diminishing returns (power curve)
- Remember the (squared) effect of variation on numbers
- Use biggish experiments (factorial) rather many small ones (+ remember: additional bonus of "interactions")
- Remember that power is affected by variation <u>and effect size</u> (Refinement vs Reduction)
- Identify and reduce sources of unwanted variation
- Include them in your experimental design (rather than just worry about it afterwards)
- Know about experimental design (ignorance is no defence)
- Talk to someone

Power

Precision

Design

• Do it <u>before</u> you start

Statistics, Experimental Design, and Animal Experimentation.

- Searching for Truth
- Power (Calculations)
- Exp. Design & Analyses: controlling variation
- Quick Recap
- Refinement vs Reduction

Quick Post course Quiz

9 questions



A highly statistically significant effect (e.g. p <0.001) means that the biological effect has to be:







When comparing groups, twice as much noise in the data means that we require how many animals to show the same statistical significance?

1. Same number

0%

- 2. 50% more animals 0%
- 3. Twice as many 2.56%
- 4. Four times as many

97.44%

5. Not sure

0%



What is the effect of Variation on Sample size (needed to obtain statistical significance).



0%



What is the relationship between Sample Size and the Power of an experiment



Statistically speaking, which of these results show a drug * sex interaction (multiple answers allowed)



POLL OPEN



The analysis table below corresponds to which graph?





What is the relationship between SE, SD and N?

¹ 1. SE = SD * √N

0%

- ^{2.} 2. SE = SD * N 0%
- ^{3.} 3. SE = SD / N 0%
- ^{₄.} 4. SE = SD /√N

100%

- ^{5.} 5. SE = √N / SD 0%
- ^{6.} 6. Not sure
 - 0%

What is a type 1 error about?

- 1. Chance of obtaining a false positive
- Chance of obtaining a false negative
 2.7%
- 3. Fundamental error at the data collection stage
- 4. Fundamental error at the analysis stage
- 5. Fundamental error at the experimental design stage 0%
- Not sure (but it sounds worse than a type 2 error !)



Thank you