

Statistical mistakes and misunderstandings: More common than you might think

David Schlesinger, Ph.D.



Lars Leksell Gamma Knife Center University of Virginia







still The truth is hard to come by

One study's result is not necessarily the truth



RESEARCH ARTICLE

Estimating the reproducibility of psychological science

Attempted replication of 100 studies published in 3 psychology journals

97% of original but only 36% of replicated studies had statistically significant results

Only 47% of original effect sizes were within the 95% confidence interval of replication effect size

....including the study about non-reproducibility



On the Reproducibility of Psychological Science Valen E. Johnson, Richard D. Payne, Tianying Wang, Alex Asher, and Soutrik Mandal Department of Statistics, Texas A&M University, College Station, TX

Reanalyzes the previous paper's data using a model that accounts for publication bias toward significant p-values and estimates distribution of effect sizes.

Model suggests that 90% of the studies tested negligible effects Publication bias towards p<=0.05 is a main cause of lack of reproducibility.

V.E.Johnson, et al., JASA 112, 2017



nature.com

Number of retraction notices has increased by factor of 10

30-40% of retractions are for technical reasons

A lack of statistical fluency may be part of the problem



Many medical physicists receive little training in practical statistics as applied to clinical outcomes studies.

However...these studies are at the heart of our profession.

How to recognize when the statistics don't quite add up?

Session educational objectives

- 1. Learn about the presence of statistical problems in published studies
- 2. Identify common signs and symptoms of potential problems in various types of statistical tests
- 3. Learn methods for correctly implementing statistical analyses of the type commonly found in medical physics publications and in routine clinical activities

How we will spend the morning....

David Schlesinger University of Virginia

Jenghwa Chang Hofstra/Northwell Health

William Sensakovic Florida Hospital/UCF/FSU/Adventist

Mike Altman Washington University Study irreproducibility, philosophy of inferential statistics, how to use and misuse p-values

Normal and non-normal distributions: Why understanding distributions are important when designing experiments and analyzing data

Linear and logistic regressions, what they try to explain and how to interpret the results

Statistical significance, statistical power, and clinical significance. How to explain your results in context.

A little quiz (not a <u>SAMS question!</u>)

You have a treatment you expect might alter performance on certain task:

Test with control and experimental groups of 20 subjects each Compare the means of each group's performance Result is significant per independent means t-test (t=2.7, df=18, p=0.01) Which statements are true? (several or none may be correct) Remember: (r=2.7, df=18, ρ =0.01)

1) You have absolutely disproved the null hypothesis

- 2) You have found the probability of the null hypothesis being true
- 3) You have absolutely proved your experimental hypothesis
- 4) You can deduce the probability of the experimental hypothesis being true

 ${\rm 5})$ If you decide to reject the null hypothesis, you know the probability that you are making the wrong decision

6) You have a reliable experimental finding – if you repeat the experiment a large number of times you would obtain a significant result 99% of the time.

. Gigerenzer et al., pub in The Sage handbook of quantitative methodology for social sciences, 2



Cellphone use causes cancer - maybe?

Study that found cell phones cause cancer in rats is riddled with red flags Besetches straight release partial results without formal review, cause alarm.

ars TECHNICA A #

Study was released before complete peer-review on a pre-publication website Control rats showed less than expected natural rate of tumor incidence and died early Incidence of tumor development correlates with age, so the early control death may have magnified the statistical findings

Publication Year	Study	Туре	# participants	Outcome
2010	Interphone Study Group	Case-control study	~5000 cases; ~5000 matched controls; 13 countries	No overall risk*
2001 (updated 2007, 2011)	Danish cohort study	Cohort study	358,000	No association
2013 (updated 2014)	Million Women Study	Prospetive cohort study	791,710	Yes (acoustic neuroma), then no association
2014	CERENAT	Multicenter case control	447 cases, 892 matched controls	No association with regular use; yes association with heaviest use
2011	Swedish pooled analysis	Pooled analysis of 2 case control studies	1251 cases, 2438 controls	Increased risk of glioma

Cellphone use causes cancer?

We still don't know!

What question are we trying to ask?

Does cellphone use cause cancer?

Or more specifically since we are scientists running a study:

What is the probability that cellphone use causes cancer given the data from our study?

Two ways to consider the problem

Bayesian methods	Frequentist methods - Null-hypothesis significance testing (NHST)
Directly calculates conditional probability of a hypothesis	Determines how extreme the observed data is Never gives the probability of a hypothesis
Requires an estimate of the prior probability of a hypothesis	Does not require an estimate of prior probability
Data can be used as it comes in Assumes data is fixed and hypotheses vary	Requires the exact specification of the experiment in advance
	Assumes hypothesis is fixed and data varies

Bayes' Theorem

$$P(H|D) = \frac{P(D|H) \cdot P(H)}{P(D)}$$











Develop a probability distributions for our prior beliefs (from prior studies, population incidence, etc.) – Prior probability

Collect some sample data from a population This can also be modelled as a distribution likelihood.

Calculate the probability of the hypothes given the data using Bayes' theorem – posterior probability



What is the probability that cellphones cause cancer given the data from our study?

Develop a probability distributions for our prior beliefs (from prior studies, population incidence, etc.) – Prior probability

Collect some sample data from a population. This can also be modelled as a distribution likelihood.

Calculate the probability of the hypothesis given the data using Bayes' theorem – posterior probability





What is the probability that cellphones cause cancer given the data from our study?

Develop a probability distributions for our prior beliefs (from prior studies, population incidence, etc.) – Prior probability

Collect some sample data from a population. This can also be modelled as a distribution likelihood.

Calculate the probability of the hypothesis given the data using Bayes' theorem – posterior probability







But...we mostly don't do this....







11

and we turn the proble	em around
All we are left with is	P(data hypothesis)
We create a reference hypothesis called the Null Hypothesis (H_0) $\qquad \qquad \qquad$	$\longrightarrow P(data H_0)$
We calculate how likely our data is assuming that this reference hypothesis is true	

Philosophically, we measure how unexpected our experimental data is The more unexpected, the less likely the null hypothesis is true.

The null hypothesis (H₀)

Serves as a reference hypothesis

Is usually the opposite of the hypothesis you hope to be true

Is frequently stated as "no difference", but doesn't have to be

Any statistically significant result assumes the null hypothesis is true

A. Field, et al., Discovering Statistics Using R, 2012

What is the probability I would get data this (or more) extreme, assuming cellphones do not cause cancer?

1. Identify null and experimental hypotheses

2. Determine the appropriate test statistic and its distribution

3. Specify the significance level you are going to use and get critical value

4. Calculate value of test statistic from your data

5. See if this is more extreme than critical value (i.e. calculate a p-value)

A. Field, et al., Discovering Statistics Using R, 20

What is the probability I would get data this (or more) extreme, assuming cellphones do not cause cancer?

Н_о,

Determine the appropriate test statisti and its distribution

3. Specify the significance level you are going to use and get critical value

4. Calculate value of test statistic from you data

5. See if this is more extreme than critical value (i.e. calculate a p-value)

A. Field, et al., Discovering Statistics Using R. 20









What is the probability I would g extreme, assuming cellphones of	et data this (or more) to not cause cancer?
dentify null and experimental hypotheses.	H _o ,
Determine the appropriate test statistic	X ^{2,} t, F, Z, etc.
nts distribution Specify the significance level you are ng to use and get critical value	p=0.05 (convention)
Calculate value of test statistic from your a	Accurant to the second se
See if this is more extreme than critical ue (i.e. calculate a p-value)	ž
	- · · · · · ·

What should I do with the p-value?

Option 1 (Fisherian method)

2.1 and 3.2 goi 4.0 dat

va

Report the p-value without any statement of "rejecting the null hypothesis" Option 2 (Neyman-Pearson method)

"Reject" the null hypothesis if p-value is below significance threshold Assumes you specified alternate hypothesis (Balance type I and type II errors)

In both cases, report the statistical test, test statistic, degrees of freedom, etc. Not just the $p\mbox{-value}!$

Ok...but what is a p-value?

Ok...but what is a p-value?

"a *p*-value is the probability under a specified statistical model that a statistical summary of the data (e.g., the sample mean difference between two compared groups) would be equal to or more extreme than its observed value."

Ok...but what is a p-value?

"a *p*-value is the probability under a specified statistical model that a statistical summary of the data (e.g., the sample mean difference between two compared groups) would be equal to or more extreme than its observed value."

Not clear?

Some facts about p-values and NHST

A p-value is the probability of obtaining data equal to or more extreme than the data you actually collected if you assume the null hypothesis is true.

A p-value is always between 0.0 and 1.0

Scientific conclusions should not be solely based on whether a p-value passes a threshold $% \left({{{\mathbf{x}}_{i}}^{2}}\right) = \left($

Using a threshold of p<0.05 is completely arbitrary

p-values depend on exact experimental setup (some of which can be implicit)

p-values depend on the sample size and spread of the data

10 Intuision Intuitive Biostatistics: A nonmathematical Guide to Statistical Thinking, 4th ed., Oxford University Press, 2018

What a p-value is NOT

A p-value is not the probability that the null hypothesis is true

1.0 minus the p-value is not the probability the alternative hypothesis is true

1.0 minus the p-value is not the probability the results will hold up under repeated experiments

A high p-value does not mean the null hypothesis is true

a p-value is not the probability of your results being "a random coincidence"

p-values are not a measure of the effect size or importance of a result

The ASA's statement on p-values: context, process, and purpose

Ronald L. Wasserstein & Nicole A. Lazar

1. P-values can indicate how incompatible the data are with a specified statistical model.

 P-values do not measure the probability that the studied hypothesis is true, or the probability that the data were produced by random chance alone.

 Scientific conclusions and business or policy decisions should not be based only on whether a p-value passes a specific threshold.

4. Proper inference requires full reporting and transparency

5. A p-value, or statistical significance, does not measure the size of an effect or the importance of a result.

6. By itself, a p-value does not provide a good measure of evidence regarding a model or hypothesis.

R. Wasserstein et al., The American Statistician

So why do we use NHST?

It is objective in that everyone will agree on the p-value for given data/statistical test/experimental design

Doesn't require prior probabilities

Requires careful description of the experiment and choice of $\operatorname{p-value}$ thresholds ahead of time

Computationally simple (and widely available)

Used for a long time (over 100 years)

The key is to understand the limits of each method

Example - Two approaches, two different answers

What is the probability P(fault | test+) that a positive test means machine really has the fault?

Using Bayes' theorem

Using NHST (H₀ = machine does NOT have a fault)

P(fault | test+) = 0.50

P(test+ | no fault) = 0.01 so p=0.01 (in this case)

Some of our QA tests show this difference in action!





So now...back to our quiz

Which statements are true? (several or none may be correct) Remember: (t=2.7, df=18, p=0.01)

- 6) You have a reliable experimental finding if you repeat the experiment a large number of times you would obtain a significant result 99% of the time. Gigerenzer et al., pub in The Sage handbook of quantitative methodology for social sciences, 2004.

	in the Three G	roups of Figure 1	viarice as frue)	
		Germany 2000		United Kingdon 1986
Statement (abbreviated)	Psychology students	Professors and lec- turers: not teaching statistics	Professors and lecturers: teaching statistics	Professors and lecturers
1. H ₀ is absolutely disproved	34	15	10	1
2. Probability of H ₀ is found	32	26	17	36
3. H, is absolutely proved	20	13	10	6
4. Probability of H, is found	59	33	33	66
5. Probability of wrong decision	68	67	73	86
6. Probability of replication	41	49	37	60

Reporting statistical tests with p-values has become a de-facto requirement



But there are problems.....

p-values just below p=0.05 are over-represented



B. Ginsel, et al., BMC Res Notes 8, 2015.

The distribution of probability values in medical abstracts: an observational study Battan Greet^(*), Achrone Aggewal^{*}, Wei Xuar² and Un Hend^{**}

Conclusions:

p-values immediately below 0.05 appear to be over-represented in the literature relative to their expected frequency Shows evidence of systematic error including publication bias, selective reporting, methodological errors, or fraud.

Authors love to invent ways of getting around the "rules"

- "a trend towards significance (p=0.06)"
- "not absolutely significant but very probably so (p>0.05)"
- "not significant in the narrow sense of the word (p=0.29)"
- "teetering on the brink of significance (P=0.06)"
- "tantalizingly close to significance (p=0.104)"
- "possibly statistically significant (p=0.10)"
- "a nonsignificant trend toward significance (p=0.1)"



The search for a p-value you like

P-hacking Data dredging Cherry-picking etc....



The search for a p-value you like

P-hacking Data dredging Cherry-picking etc....

1....



The search for a p-value you like

P-hacking Data dredging Cherry-picking etc....

nttps://xkcd.com/88



The search for a p-value you like

P-hacking Data dredging Cherry-picking etc....



The search for a p-value you like

Cherry-picking





Many papers are	e publishec	l with statistical	shortcomings
-----------------	-------------	--------------------	--------------

Category	New Med	New England J Med (n=31)		Nature Med (n=22)	
	N	%	Ν	%	
Use of a wrong or suboptimal statistical test	5	16.1	6	27.3	
No sample size / power calculation	13	41.9	22	100.0	
Failure to prove test assumptions are not violated	16	51.6	13	59.1	
Failure to define all statistical tests clearly and correctly	20	64.5	20	90.9	
Failure to state which values of p indicate statistical significance	14	45.2	15	68.2	

_

(Several) Concluding thoughts....

Characteristic	Article Year			
	1990 (n=133)	2000 (n=122)	2010 (n=106)	p-value (X ² test)
SAS	5.3%	25.4%	49.1%	<0.001
STATA	0.0%	5.7%	32.1%	<0.001
SPSS	2.3%	4.9%	13.2%	0.002
t-test	21.1%	25.4%	26.4%	0.577
Chi-square	40.6%	41.8%	41.8%	0.471
Multiple regression	24.1%	42.6%	48.1%	<0.001
Survival analysis	14.3%	22.1%	43.4%	<0.001

Excerpts from tables in Arnold LD, et al., PLOS ONE 8(10), 2013.



There is no substitute for scientific reasoning



Let's say we have 2 positioning devices and we are comparing position error

Device 1: -0.029 (0.067)mm Device 2: -0.097(0.436)mm

Welch t-test: t=1.097, df=51.325, p=0.2778

Levene's Test: F value=46.374, p<0.001

Some (I think) great resources

G. Gigerenzer et al., "The Null Ritual: What You Always Wanted to Know about Significance Testing but Were Afraid to Ask", in The Sage handbook of quantitative methodology for the social sciences, Sage, 2004.

H. Motulsky, Intuitive Biostatistics: A nonmathematical Guide to Statistical Thinking, 4th ed., Oxford University Press, 2018.

A. Field, et al., Discovering Statistics Using R, Sage, 2012. (There is a new SPSS edition available as well).

Special thanks to:

Jeff Sloan, Ph.D., Mayo Clinic

Michael Altman, Ph.D., Washington University



