



Using Bayesian methods to improve reproducibility in preclinical research

Bruno Boulanger Clément Laloux October 2022

NCS 2022



- The reproducibility crisis (the replicability crisis ?)
- > What is the question (in discovery & preclinical research) ?
- ► The design
 - The power and the Assurance
 - The ignored components
 - · Lessons learned from bioassay development and validation
- Conclusions







STATISTICAL ERRORS

P values, the 'gold standard' of statistical validity, are not as reliable as many scientists assume.

BY REGINA NUZZO

ASA, 2016

Statisticians issue warning on P values

Statement aims to halt missteps in the quest for certainty.

BY MONYA BAKER

In the produced set of the P value — a common test for judging the strength of scientific evidence — is contributing to the number of research findings that cannot be reproduced, the American Statistical Association (ASA) warned on 8 March. The group has taken the unusual step of issuing principles to guide use of the P value, which it says cannot determine whether a hypothesis is true or whether results are important.

This is the first time that the 177-year-old ASA has made explicit recommendations on such a foundational matter, says executive director Ron Wasserstein. The society's members had become increasingly concerned that the *P* value was being misapplied, in ways that cast doubt on statistics generally, he adds.

cannot indicate the importance of a finding; for instance, a drug can have a statistically significant effect on patients' blood glucose levels without having a therapeutic effect.

Giovanni Parmigiani, a biostatistician at the Dana Farber Cancer Institute in Boston, Massachusetts, says that misunderstandings about what information a *P* value provides often crop up in textbooks and practice manuals. A course correction is long overdue, he adds. "Surely if this happened twenty years ago, biomedical research could be in a better place now."

FRUSTRATION ABOUNDS

Criticism of the *P* value is nothing new. In 2011, researchers trying to raise awareness about false positives gamed an analysis to reach a statistically significant finding: that listening to music by the Beatles makes undergraduates younger



WHAT FACTORS CONTRIBUTE TO IRREPRODUCIBLE RESEARCH?

Many top-rated factors relate to intense competition and time pressure.

	 Always/often contribute Sometimes contribute
	Selective reporting
	Pressure to publish
Statistics & DoE	Low statistical power or poor analysis
Statistics & DoE	Not replicated enough in original lab
	Insufficient oversight/mentoring
	Methods, code unavailable
Statistics & DoE	Poor experimental design
	Raw data not available from original lab
	Fraud
	Insufficient peer review
Statistics & DoE	Problems with reproduction efforts
	Technical expertise required for reproduction
Statistics & DoE	Variability of standard reagents
	Bad luck
	onature 0 20 40 60 80 100%

> PHARMALEX

Nature, 2016

WHAT FACTORS COULD BOOST REPRODUCIBILITY?

Respondents were positive about most proposed improvements but emphasized training in particular.

Very likely Likely

		Better understanding of statistics	Statistics & DoE	
		Better mentoring/supervision		
		More robust design	Statistics & DoE	_
		Better teaching		
		More within-lab validation	Statistics & DoE	
		Incentives for better practice		
		Incentives for formal reproduction	Statistics & DoE	
		More external-lab validation	Statistics & DoE	
		More time for mentoring		
		Journals enforcing standards		
		More time checking notebooks		
60 80 100%	20 40	©nature 0		

► PHARMALEX





Retire statistical significance

Valentin Amrhein, Sander Greenland, Blake McShane and more than 800 signatories call for an end to hyped claims and the dismissal of possibly crucial effects.

hen was the last time you heard a seminar speaker claim there was 'no difference' between two groups because the difference was 'statistically non-significant'?

If your experience matches ours, there's a good chance that this happened at the last talk you attended. We hope that at least someone in the audinerow sare prejugation of the there is no affreence between groups or no effect of a treatment on some measured someone in the audinerow sare prejugation of the treatment on some measured to come). Nor do attainstally significant results 'prove' some other hypothesis. Such that there actually was difference.

How do statistics so often lead scientists to deny differences that those not educated in statistics can plainly see? For several generations, researchers have been warned that a

tions, researchers have been warned that a statistically non-significant result does not 'porve' the null hypothesis (the hypothesis

PERVASIVE PROBLEM

Let's be clear about what must stop: we should never conclude there is 'no difference' or 'no association' just because a P value is larger than a threshold such as 0.05 \blacktriangleright

21 MARCH 2019 | VOL 567 | NATURE | 305

The National Academies Press, 2019





The question: is my product effective ?

How to make a decision ?



What is the probability of obtaining the observed data, if the product is not effective?



What is the probability that the product is effective, given the observed data?





Currently two different ways to make a decision based on



В

Pr(observed data | product is not effective)

- Better known as the **p-value** concept
- Used in the null hypothesis test (or decision)
 - This is the likelihood of the data assuming an hypothetical explanation (eg the "null hypothesis")

This essentially depends on the very question of interest

Pr(product effective | observed data)

- Bayesian perspective
- It is the probability of efficacy given the data





A problem of decision making

The accuracy of a diagnostic test is assessed as follows:

- Sensitivity: Pr(positive result | cancer)
- Specificity: Pr(negative result | no cancer)

In practice:

Given that the diagnostic test result is positive, what is the probability you truly have cancer?

Pr(cancer | positive result) = ?





"If you use p = 0.05 to suggest that you have made a discovery, you will be wrong at least 30% of the time."



Colquhoun, D. (2014). An investigation of the false discovery rate and the misinterpretation of p-values. R. Soc. Open sci. 1(3): 140216.



► PHARMALEX



► PHARMALEX

© PharmaLex

11

A lesson from bioassay and diagnostic world



$$PPV = P(D^{+} | T^{+}) = \frac{P(T^{+} | D^{+}) \times P(D^{+})}{P(T^{+})}$$

$$=\frac{SE \times PR}{SE \times PR + (1 - SP) \times (1 - PR)}$$



Bayesian inference is the mechanism used to update the state of knowledge



The process to arrive at a posterior distribution makes use of Bayes' formula.

> PHARMALEX



Decision rules based on Posterior Probability



Direct answer to the question





Design I: Power and Assurance





► PHARMALEX

© PharmaLex

Power vs assurance

independent samples t-test ($H_0: \mu_1 = \mu_2 \text{ vs } H_1: \mu_1 \neq \mu_2$)

bayesian approach (assurance)

- ► In order to reflect the uncertainty, a large number of effect sizes, i.e. $(\mu_1 \mu_2)/\sigma_{\text{pooled}}$, are generated using the prior distributions.
- A power curve is obtained for each effect size
- the expected (weighted by prior beliefs) power curve is calculated

assumptions:



▶ PHARMALEX

An example: Power vs Assurance



► PHARMALEX

Design II: The missing components



• You know this: Meta-analysis showing study-to-study differences





► PHARMALEX

Different scenarios may happen



Groups vary independently (ρ =0)





Groups vary with some dependencies (ρ ~0.5)





If you do one trial you may get one of those outcomes....



> PHARMALEX

© PharmaLex

Impact of study-to-study variability (and lab-to-lab)

- Everyone know there are such variabilities but this is ignored in design, power calculation, evaluation,
- It is even consciously avoided !
 - To have a "Better precision" !
- It is related to the "replicability" issue, achieving a robust conclusion regardless of the study
- If ignored and existing:
 - then there is a major risk of type I error-inflation!
 - the estimates are biased
 - It violates fundamental DoE practices

Bayesian method

→ Novick S., and Zhang T. Mean Comparisons and Power Calculations to Ensure Replicability in Preclinical Drug Discovery, Stat. in Medicine, 2020.





PHARMALEX

© PharmaLex

Study "formats": example in pre-clinical pharmacology





► PHARMALEX

Improving precision of measurements

2

- Assume that:
 - θ is the parameter of interest
 - you can perform R studies of r animals

• The variance of
$$\theta$$
 is: $V(\theta) = \frac{\sigma_{Study}^2}{R} + \frac{\sigma_r^2}{R \times r}$

Currently most consider that:

$$V(\theta) = \frac{\sigma_r^2}{1 \times r}$$

But in reality, it is:

$$V(\theta) = \frac{\sigma_{Study}^2}{1} + \frac{\sigma_r^2}{1 \times r}$$

• How to design trials / allocate animals to have best precision of θ ?







Received: 18 July 2020 Revised: 18 November 2020 Accepted: 21 November 2020

DOI: 10.1002/sim.8848

RESEARCH ARTICLE

Statistics WILEY

Mean comparisons and power calculations to ensure reproducibility in preclinical drug discovery

Steven Novick[®] | Tianhui Zhang

- Use predictive distribution
- Use of informative priors is justified in preclinical research



FIGURE 3 Panels (A) and (B) show the marginal posterior distributions for μ_E , μ_F , and $\mu_F - \mu_E$ from the single-study analysis. Panels (C) and (D) show the predictive distributions for $\mu_E + \gamma_E$, $\mu_F + \gamma_F$, and $(\mu_F + \gamma_F) - (\mu_E + + \gamma_E)$

Conclusions

- ► What's the question ?
- In discovery the prior probability of success is low
- Broad use of Bayesian statistics in discovery and preclinical research will help to tackle the replicability crisis
- …combined with better design of experiments as well
- Assurance instead of Power

