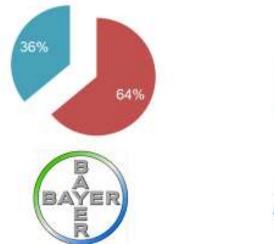


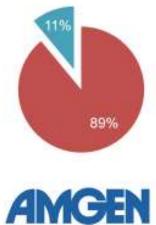
Agenda

- ➤ The reproducibility crisis, i.e. the replicability crisis
- ➤ The p-value crisis
- ➤ How to make a decision? What is the question?
- ➤ The Bayesian learning process in Pharmaceutical R&D
- ➤ The posterior predictive distribution
- ➤ The power, the Bayesian power and the Assurance
- The missing component: the elephant in the room?
- Take away message



The Replicability crisis: the beginning







Nature, 2014



STATISTICAL ERRORS

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

BY REGINA NUZZO



As the American Statistical Association officially reminded in March 2016....

REPRODUCIBILITY

Statisticians issue warning on P values

Statement aims to halt missteps in the quest for certainty.

BY MONYA BAKER

isuse of the P value — a commontest 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 biostalistician at the Dana Farber Cancer Institute in Boston, Masachusetts, 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

isuse 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



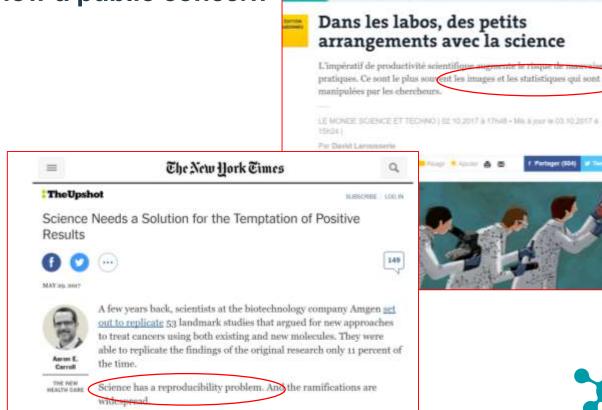
The American Statistical Association reminded in a press release some key points:

- 1. P-values can indicate how incompatible the data are with a specified statistical model.
- 2. 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.
- 3. 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.



Reproducibility now a public concern







Vidéos Archéologie Supplément partenaire : Les Prix EDF Pulse Affaire de logique

Emergency meeting: ASA Symposium in October 2017.



Scientific Method for the 21st Century: A World Beyond p < 0.05



The American Statistician - 2019









Maybe the key issue is the training in statistics

Why It Is Hard to Eliminate P-Values?

This brings us to the question of why eliminating P-value is so hard. The basic explanation is neither philosophical nor scientific, but sociologic; everyone uses them. It is the same reason we can use money. When everyone believes in something's value, we can use it for real things; money for food, and P-values for knowledge claims, publication, funding, and promotion. It does not matter if the P-value does not mean what people think it means; it becomes valuable because of what it buys.

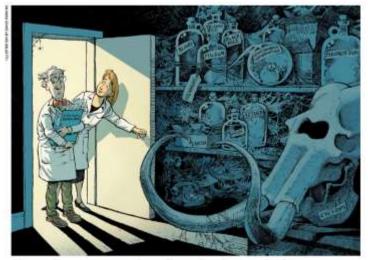
I CAN'T BELIEVE SCHOOLS
ARE STILL TEACHING KIDS
ABOUT THE NULL HYPOTHESIS.

I REMEMBER READING A BIG
STUDY THAT CONCLUSIVELY
DISPROVED IT YEARS AGO.

S. Goodman: The American Statistician - 2019



Nature March 2019





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.

When was the last time you heard a sentine speaker claim there was 'no difference' between two groups because the difference was statistically accomplished.

If your experience matches ones, there's a good chatace that this happened at the last talk you satemed. We hope that a least someone in the audience was perplexed if, as frequently happenes, a plot or table showed that there actually was a difference. How do statistics so often lead scientists to demy difference that those not obtained in statistics can plainly see? For several generations, researchers have been varied that a statistically non-significant arrest does not 'prove' the null hypothesis (the hypothesis had there in a difference between groups or no effect of a reatment on sever measured outcome?). Nor do entatistically significant results 'prove' some other hypothesis, Such misconoptism have famously warped the

literature with eventured claims and, less famously, led to claims of conflicts between studies where more exists.

We have some proposals to keep scientists from falling prey to these misconceptions.

PENVASINE PROBLEM

Let's be clear about what must stop, we should never conclude there is 'no difference' or 'no secciation' just because a P-value is larger than a threshold such as 0.05. *



National Academies May 2019

The National Academies of MEDICINE

ENGINEERING THE NATIONAL ACADEMIES PRESS

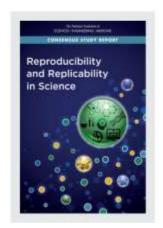
This PDF is available at http://nap.edu/25303

SHARE









Reproducibility and Replicability n Science (2019)

DETAILS

218 pages | 6 x 9 | PAPERBACK ISBN 978-0-309-48616-3 | DOI 10.17226/25303



To « p » or not to « p »: what is the question?





The objective: is my product effective?

How to make a decision?

Α

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

В

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



Two different ways to make a decision based on



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")
- Classical statistics perspective (Frequentist)



Pr(product effective | observed data)

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

The Bayesian perspective allows to directly address the question of interest.



Nature 2017

POINTS OF SIGNIFICANCE

Interpreting P values

A *P* value measures a sample's compatibility with a hypothesis, not the truth of the hypothesis.

Although P values are convenient and popular summaries of experimental results, we can be led astray if we consider them as our only metric¹. Even in the ideal case of a rigorously designed randomized study fit to a predetermined model, P values still need to be supplemented with other information to avoid misinterpretation.

A *P* value is a probability statement about the observed sample in the context of a hypothesis, not about the hypotheses being tested.

Nothing has changed in 20 years

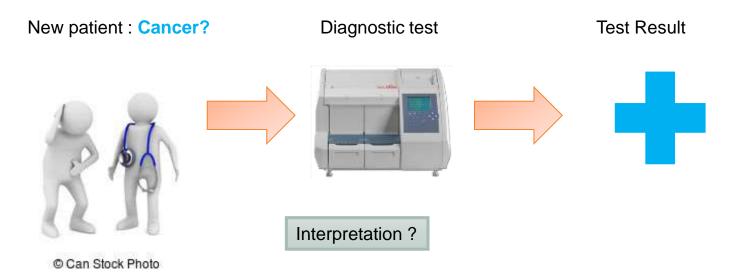
1999

Toward Evidence-Based Medical Statistics. 1: The P Value Fallacy

Steven N. Goodman, MD, PhD



What is the question ?



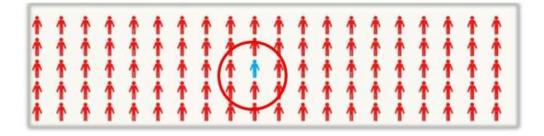
What is the probability that the patient has Cancer given the observed positive results?



Context

A disease **D** with a low prevalence

1 % of the population is diseased = D+



Major consequences if the disease is not detected



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:

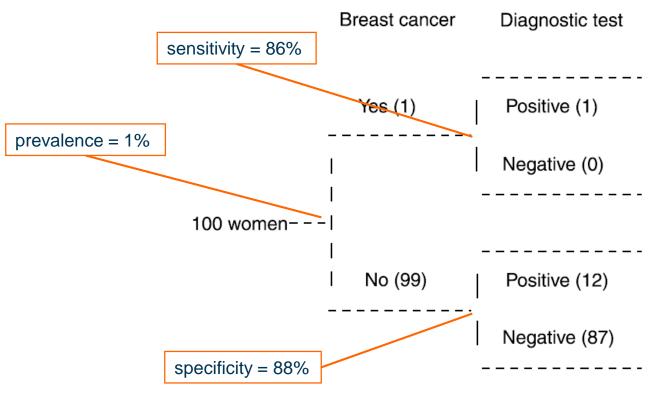
Given that the diagnostic test result is positive, what is the probability you truly have cancer?

Pr(cancer | positive result) = ?



Example

$$Pr(\text{cancer} \mid \text{positive result}) = \frac{1}{12+1} = 0.077$$



How can that be so low?
The small proportion of errors for the large majority of women who do not have breast cancer swamps the large proportion of correct diagnoses for the few women who have it.

The probability of interest dependents on the underlying prevalence of the disease.



Agresti, A. (2007). *An Introduction to Categorical Data Analysis*. Wiley, 2nd ed. Colguboup, D. (2014). An investigation of the false discovery rate and the misir

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

The clinical trial analogy

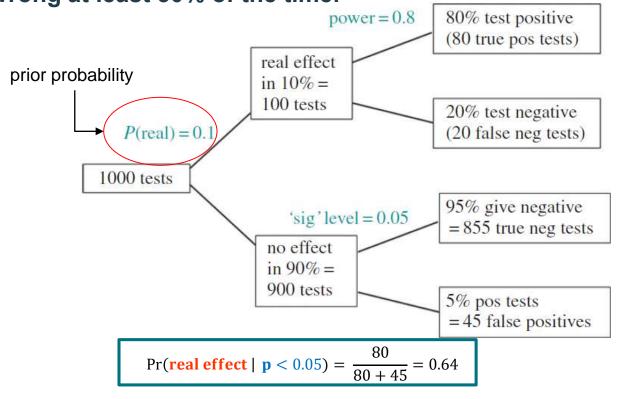


Pr(drug effective | data) = ?

depends largely on **prior probability** that there is a real effect



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

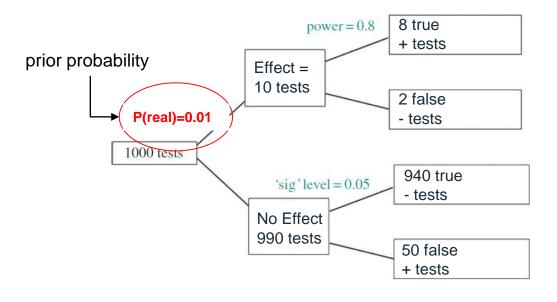




Colquhoun, D. (2014). An investigation of the false discovery rate and the misinterpretation of *p*-values. *R. Soc. Open sci.* **1**(3): 140216.



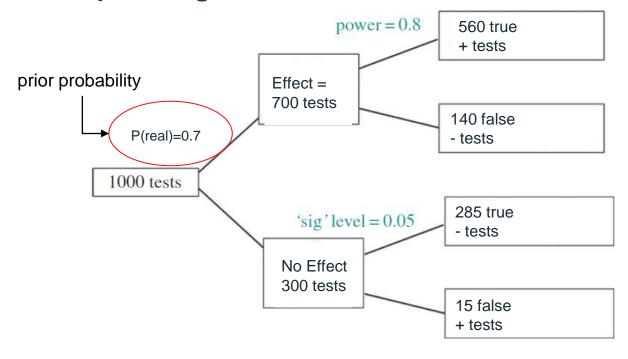
"If you use p = 0.05"....when you are in early discovery



$$Pr(real\ effect \mid p < 0.05) = \frac{8}{8+50} = 0.14 !!!!$$



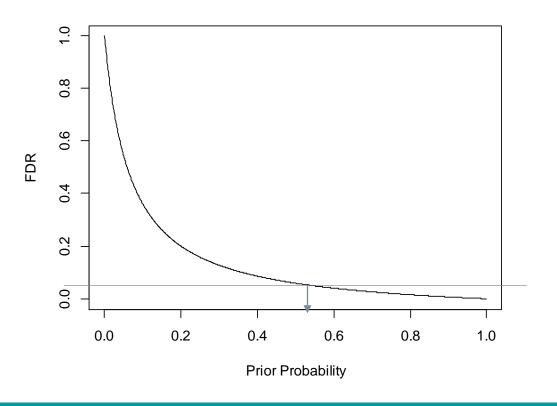
.... if the prior is good



$$Pr(real\ effect \mid p < 0.05) = \frac{560}{560 + 15} = 0.97$$



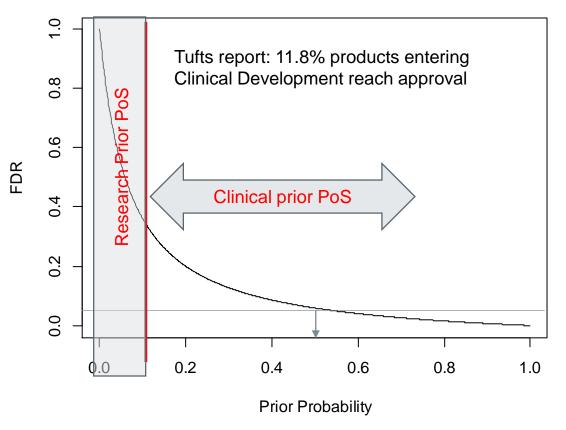
False "Discovery" Rate for p<0.05, power=0.8 as function of Prior Probability





False "Discovery" Rate for p<0.05, power=0.8 as function of Prior

Probability

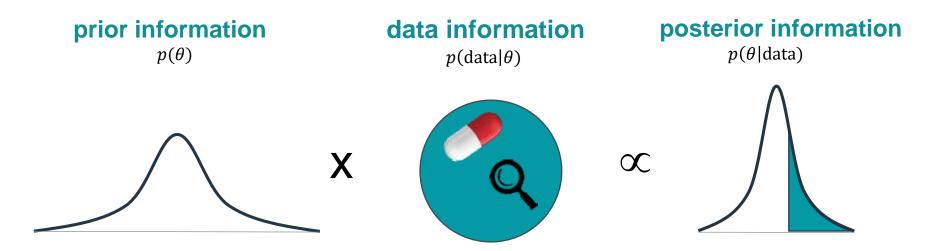




THE VALUE OF
BAYESIAN APPROACH
IN DRUG
DEVELOPMENT

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

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



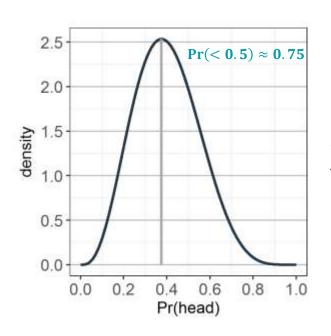


The process to arrive at a posterior distribution makes use of Bayes' formula.



The coin flipping experiment

prior information

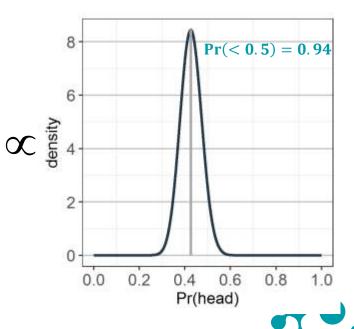


data information

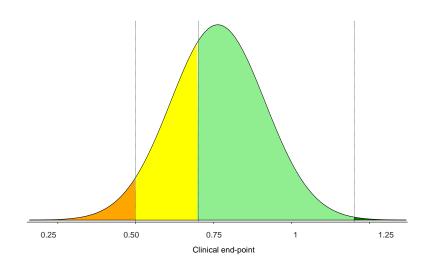


43 heads in 100 flips

posterior information

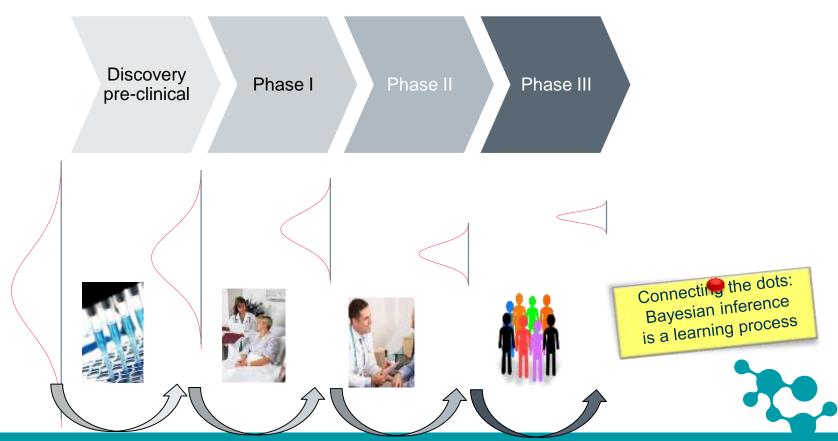


Decision rules based on Posterior Probability





Drug development is a learning process



We now have computing power to apply Bayesian statistics



Regulatory point of view

2010 & 2016 Guidance for medical device clinical trials

Guidance

for Industry and FDA Staff
Guidance for the Use of
Bayesian Statistics in
Medical Device Clinical Trials

Document issued on: February 5, 2010

Leveraging Existing Clinical Data for Extrapolation to Pediatric Uses of Medical Devices

Guidance for Industry and Food and Drug Administration Staff

Document issued on June 21, 2016.

This document will be in effect as of September 19, 2016.

The draft of this document was issued on May 6, 2015.

For questions regarding this document, contact Jacqueline Francis (CDRH) at (301) 796-6405 (Jacqueline Francis/@fda.hhs.gov). CDRHPediatricExtrapolation@fda.hhs.gov. or the Office of Communication, Outreach, and Development (CBER) at 800-835-4709 or 240-402-8010.



FDA CID initiative



Bayesian Applications



- · Safety monitoring
 - Large CV risk studies that leverage control patient data from other sources via Bayesian adaptive designs
- Oncology
 - Early phase dose-finding trial designs, e.g., CRM
 - Bayesian adaptive trials that use intermediate or accelerated approval endpoints for decision-making
- Rare diseases
 - Incorporate prior information from early phase trials
 - Use information about disease progression in analytical model
 - Compute shrinkage estimators of effects in rare subsets of disease
 - Incorporate prior information from adult trials to improve efficiency of pediatric trials

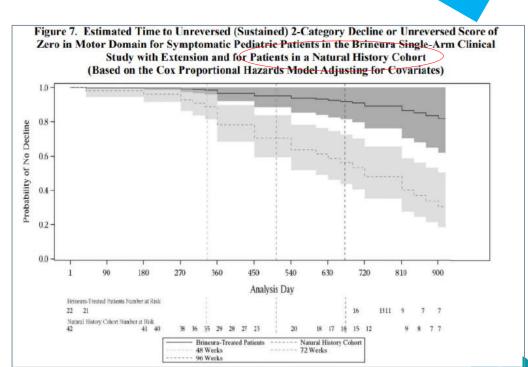






Historical control





Part II: Power, Bayesian power and Assurance

Power vs assurance

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

frequentist approach (power)

➤ A power calculation takes a particular value of the effect within the range of possible values given by H₁ and poses the question: if this particular value happens to obtain, what is the probability of coming to the correct conclusion that there is a difference?

assumptions:

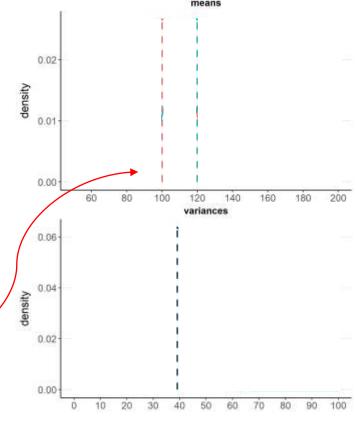
$$\mu_1 = 100;$$

$$\mu_2 = 120;$$

—— very strong priors!

$$\sigma_1^2 = \sigma_2^2 = 39$$

assumptions:



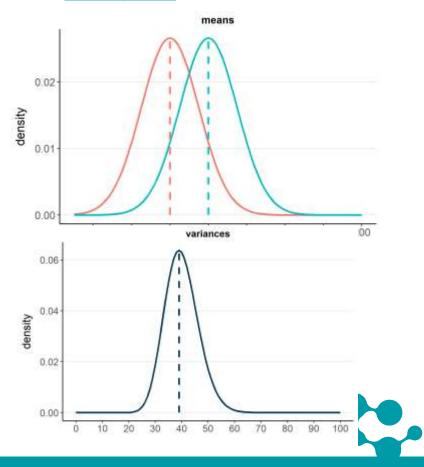
Power vs assurance independent samples t-test ($H_0: \mu_1 = \mu_2 \text{ vs } H_1: \mu_1 \neq \mu_2$)

bayesian approach (assurance)

- In order to reflect the uncertainty, a large number of effect sizes, i.e. $(\mu_1 \mu_2)/\sigma_{pooled}$, 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

- Note1: Given those priors, using the Frequentist power approach, the Probability of Success of the trial is 50%
- Note2: about 50% of Phase III trials are failing because of lack of efficacy (S. Wang, FDA, 2008)

assumptions:



(Frequentist) Power

Let R denote the rejection of the null hypothesis, the power is, assuming parameter values of $\theta = \theta^*$

$$\pi(\boldsymbol{\theta}^*, n) \coloneqq \Pr(R|\boldsymbol{\theta}^*, \boldsymbol{n})$$

➤ It is a conditional probability. It is conditional on the parameters of the model, e.g. the "true effect size" in a frequentist test and the sample size.

Assurance

"Assurance is the unconditional probability that a trial will lead to a specific outcome"

$$\gamma(n) := \int \pi(\theta, n) f(\theta) d\theta$$
$$\gamma(n) := \Pr(R) = E_{\theta}[\pi(\theta)]$$

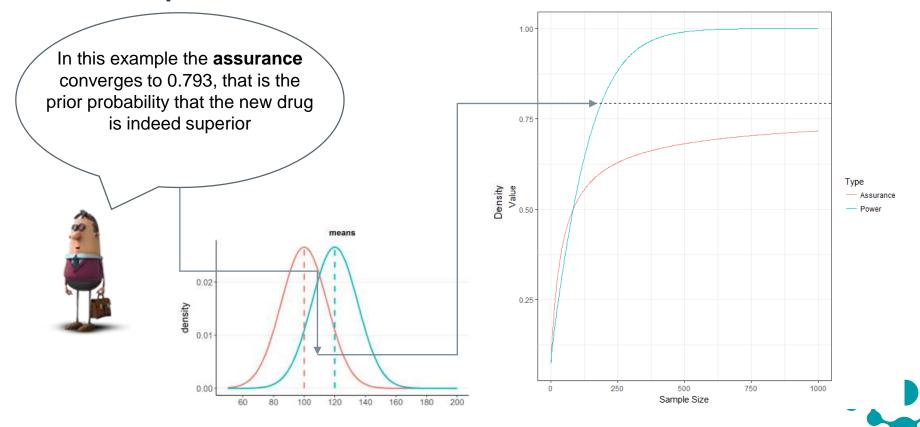
It is thus also a function of n (and eventually other nuisance parameters)

The assurance is the expected power over all possible values of theta (-> over its prior distribution...)





An example: Power vs Assurance



Difference Simulations/Predictions

Simulations

the "new observations" are drawn from distribution "centered" on estimated location and dispersion parameters (treated as "true values").

Predictions

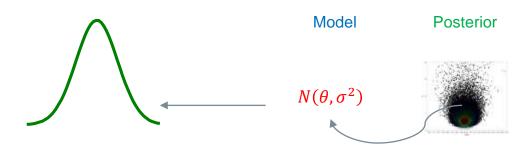
the uncertainty of parameter estimates (location and dispersion) is taken into account before drawing "new observations" from relevant distribution



Predictions

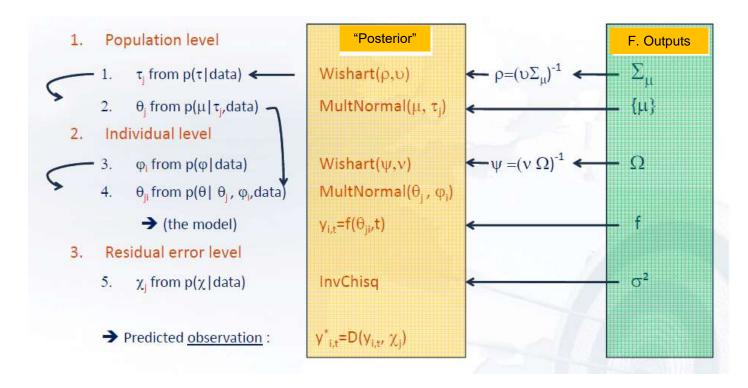
 \triangleright Given the model and the posterior distribution of its parameters, what are the plausible values for a future observation \tilde{y} ?

$$p(\tilde{y}|data) = \int p(\tilde{y}|\theta) p(\theta|data) d\theta$$





Note: It's easy to approximate the predictive distribution from Frequentist outputs



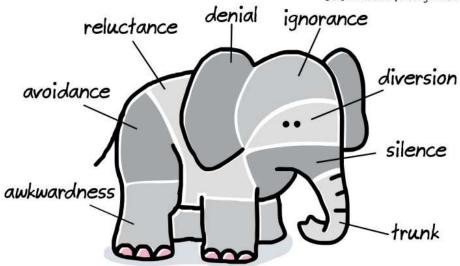


Part III: The missing component

The elephant in the room? The study-to-study variability

PARTS OF THE ELEPHANT IN THE ROOM

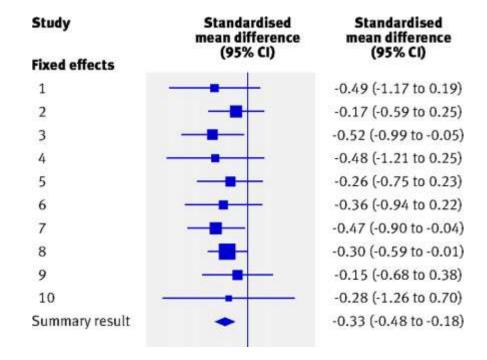
@ John Atkinson, Wrong Hands



@ John Atkinson, Wrong Hands . gocomics.com/wrong-hands . wronghands1.com

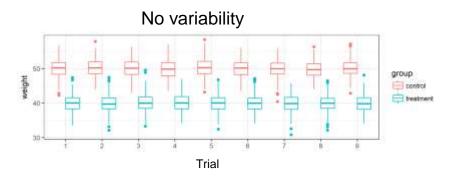


You know this: Meta-analysis showing study-to-study differences

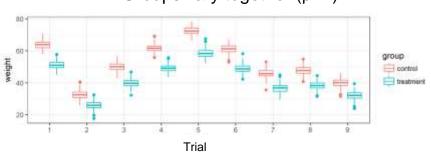




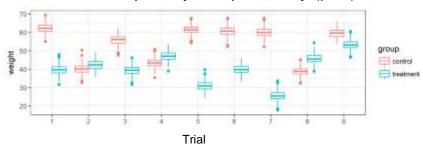
Different scenarios may happen



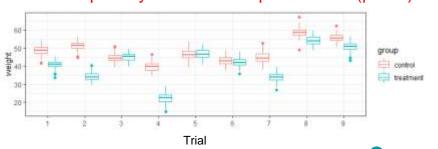
Groups vary together (ρ =1)



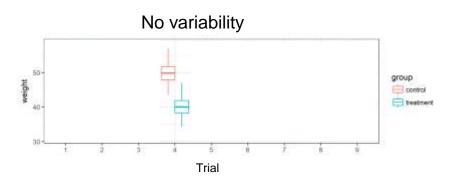
Groups vary independently (ρ =0)

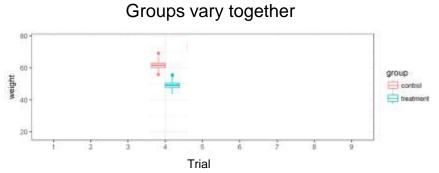


Groups vary with some dependencies ($\rho \sim 0.5$)

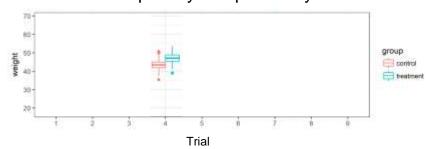


If you do one trial you may get one of those outcomes....

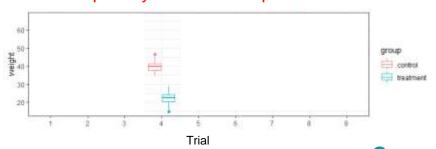




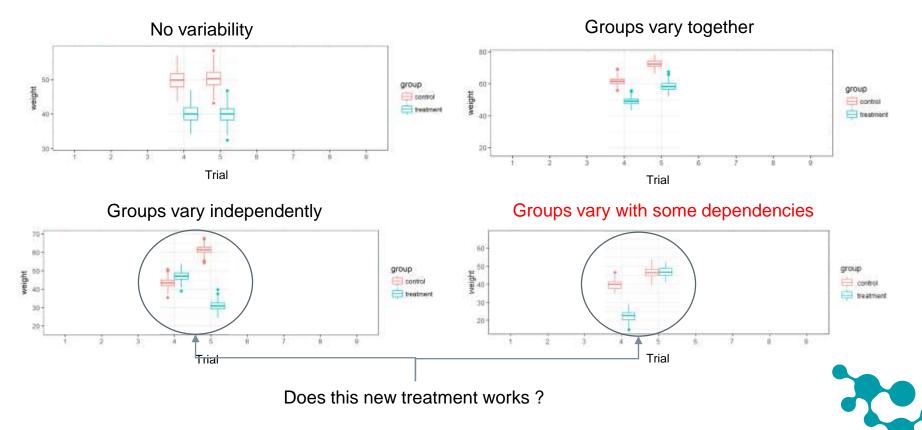
Groups vary independently



Groups vary with some dependencies

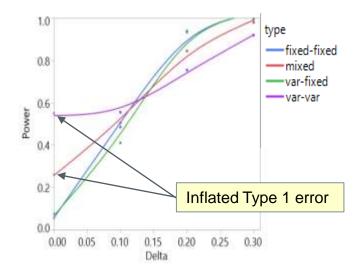


If you do two trials you may get one of those outcomes....



Impact of study-to-study variability

- Assumptions are made that there is no study-to-study variability.
- Everyone know there is such variability but this is ignored in design, power calculation, evaluation,
- This variance component is fundamental
- It is related to the "replicability" issue, achieving a conclusion regardless of the study
- If ignored and existing:
 - then there is a major risk of type I error-inflation!
 - the estimates are biased since confounded with study effect
 - It violates fundamental DoE practices: maximize
 D-optimality, ie sources of variability in studies





Study "formats": example in pre-clinical pharmacology

Classic (Common practice):

Model:
$$Y_{ij} = \mu + t_i + \varepsilon_{ij}$$

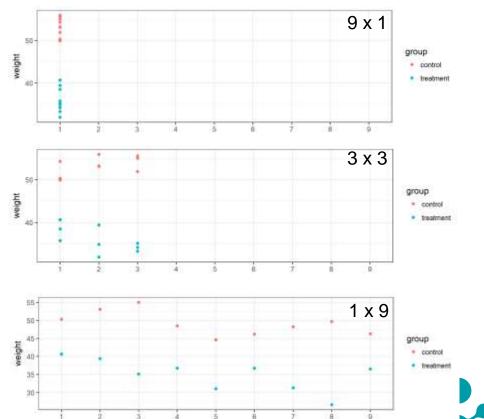
Optimal designs:

Intermediate design

Model:
$$Y_{ijk} = \mu + t_i + r_j + \varepsilon_{ijk}$$

- $r_j \sim N(0, \sigma_{study}^2)$, random effect due to the jth study, Same for both groups
- ➤ Extreme design

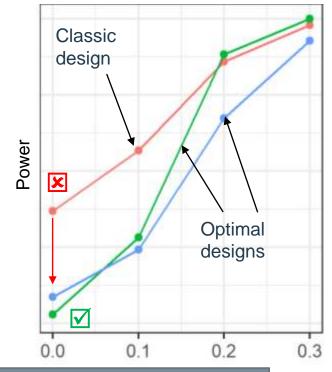
Model: $Y_{ij} = \mu + t_i + \varepsilon_{ij}$



Performance comparison of the tested designs ("formats")

- Current approach -> "classic" design (all in one study)
- Convention (USP <1032>)
 - "Convert bias into lack of precision of the estimate"
 - Control precision with sample size
- Concept of study "format":

N(total) = R (runs) x r (replicates)





Optimal designs allow to control for Type I error in all cases

Improving precision of measurements by adding noise sources

- Assume that:
 - θ is the parameter of interest
 - you can perform R studies of r patients
- The variance of θ 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 patients to have best precision of θ ?

$$\sigma_{Study}^2 \ll \sigma_r^2$$

1 trial, 10 patients

$$\frac{\sigma_{Study}^2}{R} + \frac{\sigma_r^2}{R \times r} = \frac{3}{1} + \frac{10}{1 \times 10} = 4$$

2 trials, 5 patients/trial

$$\frac{\sigma_{Study}^2}{R} + \frac{\sigma_r^2}{R \times r} = \frac{3}{2} + \frac{10}{2 \times 5} = 2.5$$



$$\sigma_{\text{Study}}^2 \gg \sigma_{\text{r}}^2$$

1 trial, 10 patients

$$\frac{\sigma_{Study}^2}{R} + \frac{\sigma_r^2}{R \times r} = \frac{10}{1} + \frac{3}{1 \times 10} = 10.3$$

2 trials, 5 patients/trial

$$\frac{\sigma_{Study}^2}{R} + \frac{\sigma_r^2}{R \times r} = \frac{10}{2} + \frac{3}{2 \times 5} = 5.3$$



Conclusions: To reduce partially the issue of Replicability:

- 1. What is the question?
- 2. Consider the study-to-study variance component in designing and sizing trials
 - Available via literature and control groups used in many trials
- 3. Consider the uncertainty of parameters estimates
- 4. Use prior distributions to compute the Assurance instead of the Power
 - Focus on probability of success of the trial beyond the power
- 5. Use Bayesian thinking and practices all the way through
 - This is a easy way to carry on the uncertainty
 - This this available now
 - This is the answer to most of your questions: Pr(drug is effective | data)



20-22 SEPTEMBER 2020

BETHESDA NORTH MARRIOTT HOTEL & CONFERENCE CENTER, ROCKVILLE, MARYLAND, USA

September 20 Short course on

Bayesian Complex Innovative Designs: a path for regulatory acceptance

Telba Irony FDA Scott Berry Berry Consultants

John Scott FDA Roger Lewis UCLA and Berry Consultants

Dionne Price FDA Karen Price Lilly

The Adolphe Quetelet Society



Contacts

Bruno Boulanger

Chief Scientific Officer

bruno.boulanger@pharmalex.com

+32 476 813936

Phone +32 10 461010

5 Rue Edouard Belin Mont-saint-Guibert

Belgium

Mobile



Thank you!

Follow us on social media





