

STATISTICS IN ANAESTHESIA: A SPECIAL SERIES

Theory and practical use of Bayesian methods in interpreting clinical trial data: a narrative review

David Ferreira^{1,2,*}, Mael Barthoulot³, Julien Pottecher⁴, Klaus D. Torp⁵, Pierre Diemunsch⁴ and Nicolas Meyer^{2,6}

¹Anesthesiology and Intensive Care Department, CHU de Besançon, Besançon, France, ²Université de Strasbourg, iCUBE, UMR7357, Illkirch Cedex, France, ³Institut Pasteur de Lille, Unité d'Epidémiologie et de Santé Publique, INSERM-U1167, Lille, France, ⁴Anesthesiology and Intensive Care Department, IHU-Strasbourg, CHU de Strasbourg, Strasbourg, France, ⁵Department of Anesthesiology, Mayo Clinic, Jacksonville, FL, USA and ⁶Public Health Department, GMRC, CHU de Strasbourg, Strasbourg, France

*Corresponding author. E-mail: dferreira@chu-besancon.fr

Summary

The critical reading of scientific articles is necessary for the daily practice of evidence-based medicine. Rigorous comprehension of statistical methods is essential, as reflected by the extensive use of statistics in the biomedical literature. In contrast to the customary frequentist approach, which never uses or gives the probability of a hypothesis, Bayesian theory uses probabilities for both hypotheses and data. This statistical approach is increasingly used for analyses of clinical trial data and for applied machine learning. The aim of this review is to compare general Bayesian concepts with frequentist methods to facilitate a better understanding of Bayesian theory for readers who are not familiar with this approach. The review is intended to be used in combination with a checklist we have devised for reading reports analysed by Bayesian methods. We compare and contrast the different approaches of Bayesian vs frequentist statistical methods by considering data from a clinical trial that lends itself to this comparative approach.

Keywords: Bayesian methods; clinician; frequentist; randomised controlled trial; statistical analysis; theory

Editor's key points

- Bayes' theorem was devised by the 18th century English theologian and mathematician Thomas Bayes.
- The Bayesian view of probability is related to degree of belief, reflecting the plausibility of an event given incomplete knowledge.
- The Bayesian approach uses probabilities for both hypotheses and data.
- This statistical approach is increasingly used in clinical trials, genetic analyses and for applied machine learning.

To keep up to date, physicians frequently consider the results of clinical trials. To appropriately interpret trial results, clinicians must understand the statistical methods reported. Currently, two kinds of methods coexist in inferential statistics: frequentist and Bayesian methods. Frequentist statistics are predominant in the field of biomedical research, based on the null hypothesis test (NHT). The probability of an event is defined as the long-term frequency of occurrence of this event, in a series of repeated trials or in a set of 'identically' conducted experiments. The name *frequentist statistics* is derived from this definition of probability. This probability is empirical and deemed objective because it relies on past observed data only.

Received: 15 August 2019; Accepted: 6 April 2020

© 2020 British Journal of Anaesthesia. Published by Elsevier Ltd. All rights reserved.
For Permissions, please email: permissions@elsevier.com

However, recent advances in the implementation of older Bayesian statistical approaches^{1,2} have led to a renewed interest in Bayesian methods, where probability is a measure of the degree of confidence or knowledge (or belief) in the occurrence of an event. This definition is consistent with the meaning of probability in everyday language. Bayesian probability is subjective and relates to statement on the credibility of an event. In this approach, the parameter of interest, with unknown values, is used in a probability distribution (the set of values that the parameter can take, with their probabilities of occurrence). Bayesian probability distributions express our degree of knowledge of the parameter with *prior* probabilities (knowledge of the parameter before further study), *posterior* probabilities (conditional on study data), and predictive probabilities (relating to data yet to be observed).

Many studies have now demonstrated the feasibility and relevance of these statistical methods in clinical trials.^{3–7} Consequently, an increasing number of therapeutic trials with results analysed by Bayesian methods are being published in major journals.^{8–19} Recommendations for Bayesian analyses have been developed,^{20–22} but these were primarily developed for researchers. They are not aimed at readers unfamiliar with Bayesian methods. For this reason, we have also developed a tool destined for practitioner to aid them in the understanding of clinical trials with analyses described by Bayesian terminology.²³ To further facilitate understanding of the Bayesian approach, here, we compare Bayesian and frequentist approaches using the IMMERSION clinical trial as an example.²⁴

Methods

Clinical trial

We apply the Bayesian approach to a clinical trial following the introduction, methods, results, and discussion (IMRAD) structure.²⁵ IMMERSION was an open prospective randomised controlled study with parallel groups. Non-pregnant women were assigned to either water immersion (2 h in a bathtub [bath group]) or to bed rest for the same duration at neutral temperature (bed group). Diuresis (primary endpoint) was assessed by measuring voiding volume. A Bayesian statistical analysis was performed. The mean difference, its 95% credible interval (CrI) and the associated posteriors probabilities were computed for the main outcome. The main analytical objective was to determine the difference in diuresis levels, hereafter called θ , between the intervention group (partial immersion: 'bath') and control group ('bed rest'). The required sample size was 20 subjects per group for an expected mean diuresis difference of 100 ml, with a standard deviation (sd) of 100 ml, a Type I error rate of 5%, and a Type II error rate of 20% in an equitail test.²⁴

Part 1: theory and general concepts of frequentist and Bayesian methods

Study objective defined by probability

Considering a coin toss, the frequentist approach dictates that, if one tosses a coin 10 times and, for example, gets six tails, then the probability of tails is, according to this experiment, 60%. In contrast, in Bayesian theory, for a coin assumed to be fair, the prior probability of heads or tails is 50%. If, for example, one gets tails on six of 10 tosses, then the posterior probability of tails, according to the results of this experiment, will be a combination of the prior probability of 50% and of the observed 6/10 tail tosses.

Therefore, this view distinguishes the notion of frequency observed during the experiment (6/10) from that of the estimated value of the unknown probability of getting heads, in a population of total tosses. The final posterior estimated value of the unknown probability of heads can thus be close to the prior 50% probability hypothesis or close to the observed 60% probability depending on the relative weight of the prior and of the data.

Hypothesis

To compare two means, the frequentist methodology asks the clinician to *define two hypotheses*:

- (i) H0: there is no difference between the two groups.
- (ii) H1: there is a difference between the two groups.

The magnitude of the difference is only specified for calculating the number of required subjects, but is not the formal subject of the test.

In the IMMERSION trial, the null hypothesis is H0 (i.e. there is no difference between the mean diuresis of the IMMERSION group and the mean diuresis of the control group). Alternative hypothesis is H1 (i.e. there is a difference between the two groups in terms of diuresis).

The Bayesian methodology asks the clinician to *define the prior knowledge of plausible means of each group*, which the statistician translates into a probability distribution (corresponding to a description of all possible values of the estimated mean and their probabilities of occurrence), called a *prior* distribution to summarise this knowledge.

In most cases, partial knowledge of the difference of the means (θ) is available before the data are collected. This prior information may come from previous similar experiments, expert opinions of the phenomenon, or basic physiological knowledge.

In the IMMERSION trial, authors assume that the hourly diuresis for a healthy woman is in most cases (i.e. with a very high probability) between 0 and 1 L. It is never negative or >4 L. The starting hypothesis is that the diuresis of the experimental group is 50 ml greater than that of the control group. Therefore, we can start from a hypothesis of a mean diuresis level of 100 ml h⁻¹ in the control group and 150 ml h⁻¹ in the experimental group with identical sds of approximately 40 ml.

Theoretical framework

In the NHT approach, only H0 is formally tested and the whole process of the test is carried out under the consideration that this hypothesis is true. However, the process attempts to show that it is false. This hypothesis test establishes a decision rule that will lead us to consider non-rejection or rejection of the null hypothesis H0. As this decision is based on the results of a sample that is only a part of the population, we cannot make a decision with certainty, and therefore, a risk of error must be considered. The decision to accept or reject H0 is made so as to minimise the risk of a wrong decision in a hypothetical series of test repetitions.

By contrast, Bayesian theory does not use test in this way. The Bayes theorem allows to determine probabilities, such as the probability of a given value of a parameter of interest θ (in IMMERSION, probability that the difference in mean diuresis of the immersion group is greater than that of the control group) knowing the data observed in the experiment (termed *y*) (i.e. here, measurements of diuresis values in each participant in both groups). If we call θ the difference in mean diuresis, then according to the Bayes theorem:

Table 1 Interpretation of theory, according to frequentist and Bayesian methods.

Null hypothesis test (NHT) approach	Bayesian approach
The probabilities of error are known as <i>alpha</i> (α) and <i>beta</i> (β). The risk α , or Type I error rate, is the probability of rejecting the null hypothesis (H0) as false when it is true.	The expression (1) can be simplified: Posterior \propto Prior \times Likelihood (‘ \propto ’ symbol reads ‘is proportional to’)
The risk β , or Type II error rate, is the probability of not rejecting the null hypothesis (H0) when the alternate hypothesis (H1) is true.	The <i>prior</i> distribution summarises the information available on the parameter of interest before the collection of our data. It corresponds to all possible values that the diuresis in the two groups can plausibly have before the study is carried out (see example in hypothesis section). It may be based on data from previous trials (other cases are described in this paper). It is obtained by the combination of the <i>prior</i> distribution (what we know about the parameter before the experiment) and the <i>likelihood</i> (what the data tell us about the parameter according to its prior probability). The data are formally turned into accumulating statistical knowledge through the use of Bayesian theorem into the posterior distribution.
The complementary probability of the Type II error (1– β) defines the power of the test and represents the probability of rejecting the null hypothesis (H0) when the hypothesis (H1) is true.	
In other words, the conclusion provided by the NHT does not measure the probability that H0 is true or false, but the probability of a given result in a repetitive process. ²⁶	
The P-value is the probability, under H0, of obtaining a statistic as extreme as the value observed in the sample. Given a threshold of significance, we compare P and α to decide whether to reject or not reject H0.	The posterior distribution describes all we know about the parameter after the experiment. It thus provides us with the parameter estimate and its credible interval.
<ul style="list-style-type: none"> • If $P < \alpha$, we reject the null hypothesis H0 (in favour of H1). • If $P > \alpha$, we do not reject H0 (in favour of H0). 	
We can then interpret the P-value as the smallest threshold of significance for which the null hypothesis is accepted.	
Two results are possible:	
(i) H0 is not rejected, and we admit there is no difference between the mean diuresis levels of the two groups.	
(ii) We reject H0 and accept H1, and we admit there is a difference between the mean diuresis levels of the two groups.	

$$\Pr(\theta|y) = \Pr(y|\theta) \Pr(\theta) / \Pr(y) \tag{1}$$

Each term of the theorem has a usual denomination. The explanation of these terms is presented in Table 1. The term $\Pr(\theta)$ is the *prior* probability of θ . The term $\Pr(y|\theta)$ is the likelihood function of θ . The term $\Pr(\theta|y)$ is the *posterior* probability.

Concept of estimation

In the frequentist approach, the parameter of interest θ is unknown, but is considered constant. The parameter θ is estimated, considering that everything we know about θ comes from the data. Therefore, this estimation relies solely on the likelihood and probability of the data under H0. H0 has no probability in itself.

In the IMMERSION study, the Bayesian approach reports a primary outcome measure comparing the mean diuresis between the two study groups. Previous knowledge of an expected mean difference in diuresis of 100 ml between the groups, under the assumption of an SD of 100 ml diuresis in each group, was only used in the sample size determination.

The principle of Bayesian estimation is to consider the parameter θ as unknown. What is unknown is uncertain and is thus given a probability distribution (see Bayesian definition of probability). We then estimate the probability that θ is within a

certain range of values. To estimate θ consists in adjusting the *prior* knowledge on θ using the information provided by the data for the experiment, through the likelihood. We then examine the conditional distribution of θ knowing the data y (i.e. the *posterior* distribution).

In the IMMERSION study, the chosen *prior* distribution for mean diuresis levels (a normally distributed outcome) in each group was a normal, or Gaussian, distribution, with two parameters: the mean and the SD. The prior was $N(m=2.68; SD=1)$ in the bath group and $N(m=1.75; SD=1.56)$ in the bed group, based on the results of Katz and colleagues,²⁷ expressed in millilitres standardised on the mean study participant weight (59.9 kg).

Bayesian and frequentist estimators can be numerically very close, especially when there are many data points. When there are few data points, the difference may be great and depends on the choice of the *prior* distribution, and we could compensate for this lack of information with a *prior* knowledge. The more observations we have, the more the relative importance of the *prior* information decreases.

Use of data

For frequentists, in analysing a clinical trial, the previous data are not used explicitly. They are used for computing the sample size, but they do not appear in the final analysis of the

data. The clinician must define an expected mean difference and SD to calculate the sample size of the study. In the context of a meta-analysis, previous data are formally included in the computation, but only as 'pure data' and not as an accumulating knowledge.

For Bayesians, in analysing a clinical trial, the previous data can be used explicitly, by introducing them in the prior distribution. Considered as a prior knowledge, they are thus formally included in the estimation process and mixed with the observed data. In the context of a meta-analysis, previous data are included in the estimation process through the prior, and the Bayesian meta-analysis can reach conclusion much sooner than its classical counterpart that considers the study as independent.²⁸

Results

Practical application of statistical analyses

Study objective defined by probability

For the frequentist, a single question using the NHT is addressed, namely, that there is no difference in diuresis reduction after partial immersion (bath) vs after bed rest. For the Bayesian, the main objective is determined by asking a single question in a probabilistic form (i.e. as a degree of belief). In IMMERSION, the main objective would be to determine what the probability is that diuresis is reduced by at least x ml after partial immersion (bath) vs after bed rest.

Process of data analyses

The classical test is performed as follows:

- (i) The null and alternative hypotheses are expressed, in particular by specifying the targeted effect size, based on previous parameter estimations, considered as known, or as hypothetical values of clinical interest. The hypotheses of the IMMERSION study are defined in the hypothesis section.
- (ii) Determination of the specific test or model to use in accordance with the nature of the variable (t-test for a quantitative outcome, etc.) and the conditions of application of the test: the main objective of the IMMERSION study is to compare a quantitative variable (diuresis) between the interventional group (bath) and the control group (bed rest) in a unilateral situation (higher diuresis level in the bath group). If the test assumptions are met (normality of distributions and homogeneity of variances), the test to be used is Student's t-test. If the conditions are not met, a non-parametric Mann–Whitney test will be used instead.
- (iii) Determination of the test values after defining the acceptance and rejection zones based on the Type I rate and the power of the test: in our example, we use the t-test. The t-test threshold value delineating the acceptance and rejection zones was defined with an α risk of 5% for $n=40$.
- (iv) Calculation of the test value from the sample data and conclusion of the test: the t-test value is calculated and compared with the t-test threshold value at a, say, 5% level. Similarly, a P-value can be computed and compared with the alpha level.

According to Gelman,²⁹ the process of Bayesian data analysis can be described by dividing it into three steps:

- (i) A probability distribution is determined for all variables of the studied problem. The sources used to construct the prior distribution can be derived from a meta-analysis, previous studies,^{4,9} expert opinion,^{4,30} or a biophysical theory. However, the use of expert opinion is debated because this introduces a more subjective nature into the study outcome analysis. The process of the expert expressing knowledge and formulating it mathematically in the prior distribution is called *elicitation*. The model should be consistent with what is known about the scientific problem involved and the nature of the parameters involved in the analysis (mean, proportions, etc.). The prior distribution specifies what is known about the difference of interest. Clinically relevant differences to be tested are specified before the study.
- (ii) From the observed data and the prior distribution, the posterior distribution is calculated by simulation (using software such as BUGS, JAGS, etc.). The main objective of the IMMERSION study is the same as in a classical framework.
- (iii) Evaluation of the goodness of fit of the model and the implications on the resulting posterior distribution. Using the suitable model, software, and technique, for the main outcome, the mean difference (in favour of the bath group) and its 95% CrI are computed, thanks to the posterior distribution, and the probability that the difference is positive and that the difference is $>0.835 \text{ ml kg}^{-1} \text{ h}^{-1}$, which corresponds to a diuresis difference of 50 ml h^{-1} , standardised on the mean study participant weight (59.9 kg).

Presentation of the results

A frequentist result usually provides the estimated parameter of interest with its 95% confidence interval (CI). When using a statistical test in addition to the test results and its 95% CI, the P-value is computed.

In the IMMERSION study, the difference between the bath and bed rest group hourly diuresis was $1.23 \text{ ml kg}^{-1} \text{ h}^{-1}$ (95% CI: 0.42; 2.05); $P=0.0039$. The P-value means that, under the null hypothesis (typically that the intervention has no effect), there is a 0.39% chance of observing a mean diuresis gap as large as that observed in the study.

The posterior distribution is the main result of a Bayesian analysis and it encompasses all the values that the parameter of interest can take *a posteriori*. From this distribution, one can deduce a mean or median with a 'range' for the parameter of interest, called a CrI. A Bayesian result is thus obtained by estimating the posterior value of the parameter of interest with its 95% CrI.

In the IMMERSION study, the posterior probability that the diuresis difference is at least $0.835 \text{ ml kg}^{-1} \text{ h}^{-1}$ in bath vs bed rest (considering an informative prior from Katz and colleagues,²⁷ $N [2.68, 1]$) was estimated to be 0.782, with a mean difference of $1.26 \text{ ml kg}^{-1} \text{ h}^{-1}$ (95% CrI: 0.20; 2.32). It was positive, so hourly diuresis was larger for the bath group than the bed rest group. The probability of this difference being positive was 0.99.

Interpretation of the results

The frequentist situation corresponds to a deductive inference, which begins with a hypothesis about the world and tests whether the observations are consistent with this hypothesis.³¹

We see that the *P*-value is not the probability that H_0 is right or false. It does not give any indication of the magnitude of the difference in the population, which is the criterion of clinical interest. It is in fact confounded with the difference magnitude and the sample size. Moreover, the fact that the threshold of significance is set at 0.05 is totally arbitrary^{12,13} and completely ignores the clinical–biological likelihood and previous knowledge.¹⁴

For Bayesians, different thresholds can be tested on the same *posterior* distribution without having to correct the tests for multiple comparisons. For example, the probabilities that the diuresis difference is greater than, say, 0.5, 0.7, or 1 ml kg⁻¹ h⁻¹ can be computed and compared without considering the number of comparisons made, because they all derive from the same distribution.

The estimation of the value of the unknown parameter is based on a random sample. This estimate will result from the combination of the known information on the parameter before the experiment and the data resulting from the experiment. The goal is not to estimate the mean difference in a given sample (this is of no interest), but to make a general estimate of this difference for the population of interest, based on the observation provided by a single sample. The use of the concept of subjective probability makes it possible to really calculate the probability, in the population, that the parameter of interest is within a given range, based on a single sample. The effect magnitude can thus be ‘isolated’ from the sample size effect.

Interpretation of intervals

Common statistical analyses rely on both descriptive and inferential analyses. The descriptive analysis is done giving the point estimation (mean, proportion, quartiles, etc.), whilst the inferential analysis relies on formal testing of the parameter. These two approaches are separated in the classical statistical context, whereas in Bayesian statistics they are both computed on a unique object (i.e. the *posterior* distribution). In both paradigms, the point estimate is given with an interval: CI in the frequentist method and CrI in the Bayesian world.

The 95% CI means that, if an experiment was repeated an infinite number of times under the same conditions, 95% of the estimated intervals would contain the true value of the parameter, whose value remains unknown.^{32,33} It is a description of what value a parameter can take under repeated sampling. Contrary to what intuition may make one believe, it does not provide the probability that the value is within a given interval.

The 95% CrI indicates that there is a probability of 0.95 (‘95% chance’) that the true value of the parameter is within the interval. The 95% *posterior* CrI depends in part on the *prior* distribution (credible values taken on diuresis levels before the study are performed) and in part on the observed data (diuresis values during the study). When the initial information is vague (lowly or uninformative *prior*), the 95% CrI, in its usual version, will have approximately the same bounds as the classical 95% CI.

Initial *prior* information that is relatively precise will reduce the dispersion of the *posterior* distribution, and thus reduce the

95% CrI, making it more precise. Moreover, one can calculate the probability that the parameter of interest value is greater or less than a threshold (set before the study), or whether it is in a predefined given interval.

Sensitivity analyses

Sensitivity analysis is defined as ‘a method to determine the robustness of an assessment by examining the extent to which results are affected by changes in methods, models, values of unmeasured variables, or assumptions’ with the aim of identifying ‘results that are most dependent on questionable or unsupported assumptions’.^{34,35}

There are different kinds of scenarios:³⁵

- (i) Modification of cut-offs or definition of outcomes
- (ii) Methods of inclusion of outliers
- (iii) Use of missing data or not
- (iv) Intention-to-treat or per-protocol analysis, etc.

Statistical methods can be modified, for example:

- (i) Parametric and non-parametric methods
- (ii) Use of different methods of adjustment (baseline characteristics and the kind of method of adjusting)

All the elements listed previously concern both approaches, but the Bayesian approach also requires that assumptions about prior distributions be included in the sensitivity analysis. Indeed, sensitivity analyses are frequently used to test the impact of several *prior* distributions on the estimation of the parameter of interest and on its *posterior* distribution. They are essential because they allow for checking the result stability under varying initial assumptions.

If the results are essentially identical when different *prior* distributions are used, then we can consider the data to be of sufficient weight and that the conclusion has been reached. Otherwise, ideally, the number of individuals included in the analysis should be increased until a stable conclusion is reached. However, this is not possible, most of the time, in a fixed design. This calls for tempering the results on the one hand, and planning a new study to increase the amount of data available on the other.

Thus, in addition to the ‘principal’ *prior* distribution, complementary analyses using an uninformative *prior* distribution, an ‘enthusiastic’ *prior* distribution (in favour of the hypothesis tested), and a ‘pessimistic’ *prior* distribution (not in favour of the hypothesis tested) can be performed. It is also possible to perform a sensitivity analysis on the model or the estimation methods used. If different priors yield different results and conclusion, the authors should acknowledge this.

Discussion

This comparison shows that, beyond the differences in their use, the difference in the interpretation of these two methods is also very important. All clinicians naturally have a moderate sense of Bayesian reasoning, often without knowing it. In the search for a diagnosis, the question is, ‘What is the probability that my patient has this or a different pathology?’ The clinician intuitively uses a pretest probability for each disease (in the form of a list of diagnoses consistent with early observations, each diagnosis being more or less probable, depending on the frequency of the disease in the general population and the frequency of each symptom for each disease). The pretest probability of each disease under consideration is increased or

decreased according to the results of the diagnostic test (resulting in a post-test probability). The same reasoning is used when making decisions about treatments. The initial confidence about whether or not a treatment will work is influenced by the clinician's knowledge of pathophysiology and pharmacology, by reading reports of high-quality RCTs,³⁶ and should be confirmed by practice.³⁷

The *prior* distribution is called 'informative' if it conveys a lot of information on the parameter (i.e. is precise on its prior values). 'Non-informative' or better 'lowly informative' distribution can also be used: all possible values of the parameter of interest θ are, in the eyes of the expert, equally (or almost equally) likely, that is, he (she) will not bet more on one value than on another. Using this kind of *prior*, it is the data (the likelihood) that will have the greatest impact on the *posterior* distribution.^{4,38,39} A 'typical' non-informative *prior* distribution is a normal distribution with a zero mean and very large variance.³⁸ In Phase III trials, non-informative (lowly) informative priors are generally preferred by regulators.

To follow a hypothetical deductive scientific approach, **the *prior* distribution must be defined before data collection.** The *prior* distribution is susceptible to biases if it is constructed after data collection, especially if it is based on expert opinion. However, the *prior* could be revised during the study if other information becomes available, provided that the revisions occur independently of the results of the study.

In Bayesian methodology, the *prior* knowledge is explicitly formalised, enabling authors to use informative knowledge as soon as possible under the control of the reader's common sense. In the IMMERSION study, assuming that a mean diuresis level is between -10^6 and $+10^6$ L h⁻¹ would be difficult. However, these are values considered possible by the alternative hypothesis in the frequentist NHT. With Bayesian methods, the combination of the plausible values for the parameter of interest before the experiment (*prior*) and the information from the experiment allows for obtaining our result: the *posterior* distribution (all the values that *posterior* parameters can take and their probabilities). When the sample size is large ($n \rightarrow \infty$), the influence of the *prior* distribution fades and the likelihood of the data makes up the bulk of the *posterior* distribution. When the state of knowledge about the problem in question permits only a very vague *prior*, the likelihood of the data will also have a predominant impact on the *posterior* distribution. Therefore, there will be disagreement only when the data are insufficient to yield the same *posterior* distribution for different prior assumptions. If the data are of sufficient quantity, all the different *prior* values will approximately approach the same *posterior* distribution, and the debate will be settled. Furthermore, the *posterior* distribution is often difficult to obtain analytically, which requires simulation through specific statistical software. This is one of the reasons why Bayesian methods have sometimes been deemed to be difficult to use in the past.

In frequentist methodology, this *prior* knowledge is implicitly used for the parameter of interest when calculating the sample size. This calculation can never be reassessed without prior agreement in the protocol. In practice, it undergoes the fluctuations inherent in the daily constraints of research (patient lost to follow-up, protocol violations, etc.), therefore making the conclusions of the study obsolete.⁴⁰ The final usable result of the statistical tests is the *P*-value, which determines a threshold of arbitrary significance to reject the null hypothesis. The *P*-value does not measure the probability

that the hypothesis studied is true or that the data are produced by chance alone.

Many clinicians instinctively interpret the *P*-value by using an inductive reasoning that is characteristic of the Bayesian method.³¹ In a clinical case, the frequentist statistician asks, 'What is the probability of having a temperature $>39.5^\circ$ with a diagnosis of influenza?', whereas the Bayesian statistician asks, 'What is the probability of having the flu, knowing that the temperature is $>39.5^\circ$?', the typical reasoning used in a diagnostic procedure.

In summary, we have compared Bayesian methods with frequentist methods using the IMRAD structure to provide a better understanding of the general theories and concepts of these two methods. This document is intended to serve as an adjunct to our reading grid,²³ but can also be used individually to understand Bayesian methodology.

Authors' contributions

Article development: all authors

Article drafting: all authors

Final review of article: all authors

Declarations of interest

The authors declare that they have no conflicts of interest.

Funding

Centre Hospitalier Universitaire de Strasbourg.

Acknowledgements

The authors thank Laura Smales (BioMedEditing, Toronto, Canada) for English editing.

References

1. Gilks WR, Thomas A, Spiegelhalter DJ. A language and program for complex Bayesian modelling. *J R Stat Soc Ser Stat* 1994; **43**: 169–77
2. Lunn DJ, Thomas A, Best N, Spiegelhalter D. WinBUGS - A Bayesian modelling framework: Concepts, structure, and extensibility. *Stat Comput* 2000; **10**: 325–37
3. Berry DA. *Statistics: a Bayesian perspective*. Belmont, CA: Duxbury Press; 1996
4. Spiegelhalter DJ, Abrams KR, Myles JP. *Bayesian approaches to clinical trials and health-care evaluation*. London, U.K.: Wiley; 2004
5. Goodman SN. Introduction to Bayesian methods, I: measuring the strength of evidence. *Clin Trial*. 2005; **2**: 282–90. discussion 301–4, 364–78
6. Louis TA. Introduction to Bayesian methods, II: fundamental concepts. *Clin Trial*. 2005; **2**: 291–4. discussion 301–4, 364–78
7. Berry DA. Introduction to Bayesian methods, III: use and interpretation of Bayesian tools in design and analysis. *Clin Trial*. 2005; **2**: 295–300. discussion 301–4, 364–78
8. Roberts KA, Dixon-Woods M, Fitzpatrick R, Abrams KR, Jones DR. Factors affecting uptake of childhood immunisation: a Bayesian synthesis of qualitative and quantitative evidence. *Lancet* 2002; **360**: 1596–9

9. Brophy JM, Joseph L. Placing trials in context using Bayesian analysis. GUSTO revisited by Reverend Bayes. *JAMA* 1995; **273**: 871–5
10. Muss HB, Berry DA, Cirincione CT, et al. Adjuvant chemotherapy in older women with early-stage breast cancer. *N Engl J Med* 2009; **360**: 2055–65
11. Baeten D, Baraliakos X, Braun J, et al. Anti-interleukin-17A monoclonal antibody secukinumab, in-treatment of ankylosing spondylitis: a randomised, double-blind, placebo-controlled trial. *Lancet* 2013; **382**: 1705–13
12. Shah PL, Slebos DJ, Cardoso PFG, et al. Bronchoscopic lung-volume reduction with Exhale Airway Stents for Emphysema (EASE trial): randomised, sham-controlled, multicentre trial. *Lancet* 2011; **378**: 997–1005
13. Wilber DJ, Pappone C, Neuzil P, et al. Comparison of antiarrhythmic drug therapy and radiofrequency catheter ablation in patients with paroxysmal atrial fibrillation: a randomized controlled trial. *JAMA* 2010; **303**: 333–40
14. Singh SM, Austin PC, Chong A, Alter DA. Coronary angiography following acute myocardial infarction in Ontario, Canada. *Arch Intern Med* 2007; **167**: 808–13
15. Tinetti ME, Baker DI, King M, et al. Effect of dissemination of evidence in reducing injuries from falls. *N Engl J Med* 2008; **359**: 252–61
16. Gausche M, Lewis RJ, Stratton SJ, et al. Effect of out-of-hospital pediatric endotracheal intubation on survival and neurological outcome: a controlled clinical trial. *JAMA* 2000; **283**: 783–90
17. Morris RK, Malin GL, Quinlan-Jones E, et al. Percutaneous vesicoamniotic shunting versus conservative management for fetal Lower Urinary Tract Obstruction (PLUTO): a randomised trial. *Lancet* 2013; **382**: 1496–506
18. Cannon CP, Shah S, Dansky HM, et al. Safety of anacetrapib in patients with or at high risk for coronary heart disease. *N Engl J Med* 2010; **363**: 2406–15
19. Holmes DR, Teirstein P, Satler L, et al. Sirolimus-eluting stents vs vascular brachytherapy for in-stent restenosis within bare-metal stents: the SISR randomized trial. *JAMA* 2006; **295**: 1264–73
20. Spiegelhalter DJ, Myles JP, Jones DR, Abrams KR. Bayesian methods in health technology assessment: a review. *Health Technol Assess* 2000; **4**: 1–130
21. The BaSiS Group. Bayesian standards in science (BaSiS) [Internet]. Draft: Sept. 13, 2001. Available at: <http://lib.stat.cmu.edu/bayesworkshop/2001/BaSiS.html>.
22. Sung L, Hayden J, Greenberg ML, Koren G, Feldman BM, Tomlinson GA. Seven items were identified for inclusion when reporting a Bayesian analysis of a clinical study. *J Clin Epidemiol* 2005; **58**: 261–8
23. Ferreira D, Barthoulot M, Pottecher J, Torp KD, Diemunsch P, Meyer N. A check-list for Bayesian clinical trials in anesthesiology. *Br J Anaesth* 2019
24. Effect of Immersion, Performed under the conditions of obstetrical dilatation bath, on Diuresis and Hemodynamic Variables in Young Women — Full Text View — ClinicalTrials.gov [Internet]. [cited 2018 Jul 24]. Available from: <https://clinicaltrials.gov/ct2/show/NCT02409953>.
25. International Committee of Medical Journal Editors. Uniform requirements for manuscripts submitted to biomedical journals. *N Engl J Med* 1997; **336**: 309–16
26. Le bayésianisme aujourd’hui. Fondements et pratiques Égré P, Drouet I, editors. *Econ Hist Methodol Philos* 2017: 453–8
27. Katz VL, Ryder RM, Cefalo RC, Carmichael SC, Goolsby R. A comparison of bed rest and immersion for treating the edema of pregnancy. *Obstet Gynecol* 1990; **75**: 147–51
28. Halsey LG. The reign of the p-value is over: what alternative analyses could we employ to fill the power vacuum? *Biol Lett* 2019; **15**: 20190174
29. Gelman A. *Bayesian data analysis*. 3rd Edn. Boca Raton, FL: CRC Press; 2014
30. Parmar MK, Griffiths GO, Spiegelhalter DJ, Souhami RL, Altman DG, van der Scheuren E. Monitoring of large randomised clinical trials: a new approach with Bayesian methods. *Lancet* 2001; **358**: 375–81
31. Wijesundera DN, Austin PC, Hux JE, Beattie WS, Laupacis A. Bayesian statistical inference enhances the interpretation of contemporary randomized controlled trials. *J Clin Epidemiol* 2009; **62**: 13–21. e5
32. Greenland S, Senn SJ, Rothman KJ, et al. Statistical tests, P values, confidence intervals, and power: a guide to misinterpretations. *Eur J Epidemiol* 2016; **31**: 337–50
33. Neyman J. Outline of a theory of statistical estimation based on the classical theory of probability. *Phil Trans R Soc Lond Ser Math Phys Sci* 1937; **236**: 333–80
34. Schneeweiss S. Sensitivity analysis and external adjustment for unmeasured confounders in epidemiologic database studies of therapeutics. *Pharmacoepidemiol Drug Saf* 2006; **15**: 291–303
35. Thabane L, Mbuagbaw L, Zhang S, et al. A tutorial on sensitivity analyses in clinical trials: the what, why, when and how. *BMC Med Res Methodol* 2013; **13**: 92
36. Patel A, Gronseth G, Glantz M. Carotid artery stenosis and randomized controlled trials—should we abandon the gold standard (P7.170). *Neurology* 2014; **82**(10 Supplement). P7.170
37. Wheatley K, Clayton D. Be skeptical about unexpected large apparent treatment effects: the case of an MRC AML12 randomization. *Control Clin Trial*. 2003; **24**: 66–70
38. Diamond GA, Forrester JS. Clinical trials and statistical verdicts: probable grounds for appeal. *Ann Intern Med* 1983; **98**: 385–94
39. Berry DA. Bayesian clinical trials. *Nat Rev Drug Discov* 2006; **5**: 27–36
40. Baker M. 1,500 scientists lift the lid on reproducibility. *Nature* 2016; **533**: 452–4