

# A Challenge for Statistical Instructors: Teaching Bayesian Inference Without Discarding the “Official” Significance Tests

BRUNO LECOUTRE, MARIE-PAULE LECOUTRE  
and JEAN-MARIE GROUIN  
C.N.R.S. and Université de Rouen, France

## Abstract

The use of frequentist Null Hypothesis Significance Testing (NHST) is so an integral part of scientists' behavior that its uses cannot be discontinued by flinging it out of the window. Faced with this situation, our teaching strategy involves a smooth transition towards the Bayesian paradigm. Its general outlines are as follows. (1) To present natural Bayesian interpretations of NHST outcomes to draw attention to their shortcomings. (2) To create as a result of this the need for a change of emphasis in the presentation and interpretation of results. (3) Finally to equip students with a real possibility of thinking sensibly about statistical inference problems and behaving in a more reasonable manner. Our conclusion is that teaching the Bayesian approach in the context of experimental data analysis appears both desirable and feasible.

*Keywords:* EXPERIMENTAL DATA ANALYSIS; TEACHING BAYESIAN PROCEDURES; BAYESIAN INTERPRETATION OF  $P$ -VALUES; EFFECT SIZES;.

## 1 INTRODUCTION

Many recent papers have stressed on the necessity of changes in reporting experimental results. This has been recently made official by the American Psychological Association (Wilkinson *et al.*, 1999). A more and more widespread opinion is that inferential procedures that provide genuine information about the size of effects must be used in addition or in place of Null Hypothesis Significance Testing (NHST). Such procedures have been developed both in the frequentist and Bayesian frameworks. However they are again rarely used, in spite of the fact that they are nowadays straightforward to implement. In consequence it must be urged to reform the teaching of statistical inference and to include these procedures, even in introductory courses for non-statistician students. For more than twenty years now, we and other colleagues have gradually introduced the Bayesian approach to experimental data analysis, in courses and seminars for audiences of various backgrounds, especially in psychology.

The present paper is divided into three sections. (1) We briefly recall the shortcomings of NHST. (2) We present and criticize the recommendations proposed by the American Psychological Association to overcome these shortcomings. (3) As an alternative, we suggest teaching Bayesian methods as a therapy against the misuses and abuses of NHST. The feasibility of this teaching is illustrated.

## 2 THE CURRENT CONTEXT OF THE EXPERIMENTAL RESEARCH

**The shortcomings of Null Hypothesis Significance Testing.** Experimental research is facing a paradoxical situation. On the one hand Null Hypothesis Significance Testing (NHST) is required in most scientific publications as an unavoidable norm and it often appears as a label of scientificness. But on the other hand NHST leads to innumerable misinterpretations and misuses (Lecoutre *et al.*, 2001).

**Mistaking statistical significance for scientific significance.** The more significant a result is, the more scientifically interesting it is, and/or the larger the true effect is. This has been one of the most often denounced error. From a survey of research articles published in three different psychology journals, Craig *et al.* (1976) concluded that “researchers and journal editors as a whole tend to (over)rely on significant differences as the definition of meaningful research.”

**Improper uses of nonsignificant results as “proof of the null hypothesis”.** About half of the articles published in the 1994 issue of the *Journal of Abnormal Psychology* contained (in most cases unjustified) conclusions such as “there is no difference between groups” or “there is no interaction effect” based on nonsignificant tests (Poitevineau, 1998).

**“Non frequentist” interpretations of  $p$ -values.**  $1 - p$  is most often interpreted (even by experienced users) as the probability that the alternative hypothesis is true or as evidence of the replicability of the result (Oakes, 1986; Freeman, 1993; Falk and Greenbaum, 1995).

## 3 TIME FOR CHANGE IN REPORTING EXPERIMENTAL RESULTS

### 3.1 New Guidelines in Psychology

While users’ uneasiness towards NHST is ever growing (Lecoutre, 2000), it seems to be nowadays a crucial period of time. Many recent papers and editorials have stressed on the necessity of changes in reporting experimental results, especially in presenting and interpreting effect sizes (see e.g., Loftus, 1993; Serlin and Lapsley, 1993; Rouanet, 1996; Schmidt, 1996; Thompson, 1996; Heldref Foundation, 1997; Murphy, 1997; Brandsttter, 1999; Lecoutre and Poitevineau, 2000).

The majority trend is to advocate the use of confidence intervals. This has been recently made official by the American Psychological Association (Wilkinson *et al.*, 1999). The following extracts are proposed guidelines for revising the statistical

section of the American Psychological Association Publication Manual (italics are ours).

**Hypothesis tests.** “It is hard to imagine a situation in which a dichotomous accept-reject decision is better than reporting an actual p value or, better still, a confidence interval. *Never use the unfortunate expression ‘accept the null hypothesis.’ Always provide some effect-size estimate when reporting a p value.*”

**Interval estimates.** “*Interval estimates should be given for any effect sizes involving principal outcomes.* Provide intervals for correlations and other coefficients of association or variation whenever possible.”

**Effect sizes.** “*Always present effect sizes for primary outcomes.* If the units of measurement are meaningful on a practical level (e.g., number of cigarettes smoked per day), then we usually prefer an unstandardized measure (regression coefficient or mean difference) to a standardized measure.”

**Power and sample size.** “Provide information on sample size and the process that led to sample size decisions. *Document the effect sizes, sampling and measurement assumptions, as well as analytic procedures used in power calculations.* Because power computations are most meaningful when done before data are collected and examined, it is important to show how effect-size estimates have been derived from previous research and theory in order to dispel suspicions that they might have been taken from data used in the study or, even worse, constructed to justify a particular sample size.”

### 3.2 Further Difficulties

If the research community accepts the above recommendations, confidence intervals could quickly become a compulsory norm in experimental publications. However, for many reasons due to their *frequentist* conception, confidence intervals can hardly be viewed as the ultimate method. Indeed the appealing feature of confidence intervals is the result of a fundamental misunderstanding. It is so strange to treat the data as random even after observation that the “right” frequentist interpretation of confidence intervals does not make sense for most users. Ironically the incorrect natural interpretation of confidence intervals in terms of probability about parameters is encouraged by the duplicity of most statistical instructors who tolerate and even use it: “We can be 95% confident that the population mean is between 114.06 and 119.94” (Kirk, 1982, page 43). It is undoubtedly this Bayesian interpretation of confidence intervals in terms of “a fixed interval having a 95% chance of including the true value of interest” which is their appealing feature. After many attempts to teach the “correct” interpretation of frequentist procedures, we completely agree with Freeman (1993) that in these attempts “we are fighting a losing battle”.

Furthermore the proposed guidelines are both partially technically redundant and conceptually incoherent. They should result in replacing the ritual of NHST by another set of rituals, without supplying a real statistical thinking. As the authors of these guidelines state, it is probably true that “*statistical methods should guide and discipline our thinking but should not determine it.*” However it is no less true that “*we need statistical thinking, not rituals*” (Gigerenzer, 1998).

## 4 OUR PROPOSALS

The use of NHST is so an integral part of scientists' behavior that its misuses and abuses should not be discontinued by flinging it out of the window. Our conclusion is that the sole effective therapy against its "damages" is a smooth transition towards the Bayesian paradigm.

### 4.1 Standard Bayesian Methods for Experimental Data Analysis

For many years we have worked with colleagues in France within this perspective in mind in order to develop standard "noninformative" Bayesian methods for the most familiar situations encountered in experimental data analysis (see e.g., Rouanet and Lecoutre, 1983; Lecoutre *et al.*, 1995; Lecoutre, 1996; Lecoutre and Charron, 2000; Lecoutre and Poitevineau, 2000; Rouanet *et al.*, 2000). Based on more useful working definitions than frequentist procedures, these methods are fully justified, at least as objective, and they can be used and taught as easily as the  $t$ ,  $F$  or *chi-square* tests. They are nowadays available and they are concrete proposals for bypassing the shortcomings of NHST and improving the current statistical methodology and practice.

Our statistical teaching and consulting experience, especially in psychology, revealed us that these methods were far more intuitive and much closer to the thinking of scientists than frequentist procedures (see also Kadane, 1995). They have been applied many times to real data and well accepted by psychology journals (see e.g., Hoc and Leplat, 1983; Ciancia *et al.*, 1988; Lecoutre, 1992; Hoc, 1996; Clment and Richard, 1997; and many experimental articles published in French).

### 4.2 Our Teaching Strategy

Our strategy in front of the misuses of NHST is to introduce Bayesian methods as follows.

- (1) To present natural *Bayesian interpretations* of NHST outcomes to call attention about their shortcomings.
- (2) To create as a result of this the need for *a change of emphasis in the presentation and interpretation* of results.
- (3) Finally to equip students with a real possibility of *thinking sensibly about statistical inference* problems and behaving in a more reasonable manner.

### 4.3 Bayesian Alternatives to the Guidelines

**Hypothesis tests: Bayesian interpretation of  $p$ -values.** A well-known feature of Bayesian inference is that it provides insightful interpretations of many frequentist procedures. For most usual situations of experimental data analysis, it offers the student a smooth transition from the traditional techniques to the Bayesian method. Moreover the Bayesian interpretation of  $p$ -values clearly points out the methodological shortcomings of NHST. In particular, it becomes apparent that a "nonsignificant" outcome is hardly worth anything.

**Interval estimates: Bayesian interpretation of the usual CI.** It becomes correct to say that “there is a 95% chance of the parameter being included between the fixed bounds of the interval” (conditionally on the data).

**Effect sizes: straight Bayesian answers.** Beyond the reinterpretations of the usual frequentist procedures, other Bayesian statements give straight answers to the question of effect sizes. For instance, it can be reported that “there is a 80% posterior probability of a large positive difference (e.g.  $\delta > +2$ ), a 20% probability of a small difference ( $-2 < \delta < +2$ ), and a 0.01% probability of a large negative difference ( $\delta < -2$ )”. Such a statement has no frequentist counterpart.

**Power and sample size: Bayesian data planning and monitoring.** Bayesian predictive probabilities are efficient tools for designing (“how many subjects?”) and monitoring (“when to stop?”) experiments. The predictive distribution of a test statistic can be used to include and extend the frequentist notion of power in a way that has been termed *predictive* power (Spiegelhalter *et al.*, 1986). More generally, predictive procedures give the researcher a very appealing method to evaluate the chances that the experiment will end up showing a conclusive result, or on the contrary a non-conclusive result. The prediction can be explicitly based on either the hypotheses used to design the experiment, expressed in terms of prior distribution, or on partial available data, or on both (see especially Berry, 1991; Lecoutre *et al.*, 1995; Dignam *et al.*, 1998; Lecoutre, 2001).

**Introducing “informative” priors.** When a standard Bayesian analysis suggests a given conclusion, another line of attack is to investigate the impact of *skeptical* or *handicap* prior distributions. In this way, the experiment will only stop if the partial data give sufficient evidence to counterbalance it (see e.g. Spiegelhalter, Freedman and Parmar, 1994). Here is a very appealing way to introduce “informative” priors. Contrasting the resulting posterior with the noninformative solution allows the students to understand the relative roles of sample sizes, data and external information.

## 5 TEACHING BAYESIAN METHODS FOR ANALYSIS OF VARIANCE

### 5.1 The Specific Analysis Approach

In the graduate statistics course in psychology, we especially developed Bayesian methods in the analysis of variance framework. A specificity of experimental data analysis is that experimental investigations frequently involve complex designs, especially repeated-measures designs. This complexity must be explicitly taken into account for teaching Bayesian procedures. A simple way to do it is to use the *specific analysis approach*. Roughly speaking, a specific analysis for a particular effect consists in handling only data that are relevant for it. Further justifications can be found in Rouanet and Lecoutre (1983) (see also Rouanet, 1996).

This conception brings a simple way for teaching the analysis of variance methods (including the multivariate procedures) with a Bayesian perspective. This point will be illustrated here from a typical numerical example.

## 5.2 A Typical Example: Reaction Time Experiment

As an elementary illustration, let us consider the following example. In a psychological experiment, the subject must react to a signal. The experimental design involves two crossed repeated factors: Factor  $A$  (signal frequency) with two levels, frequent ( $a1$ ) and rare ( $a2$ ), and Factor  $B$  (foreperiod duration), also with two levels, short ( $b1$ ) and long ( $b2$ ). The  $n = 12$  subjects are divided into three groups of four subjects each. The main research hypothesis is a null (or about null) interaction effect between factors  $A$  and  $B$  (*additive model*). Since  $A$  and  $B$  are both two-level factors, their interaction can be represented by a single contrast among the four conditions. We take the contrast  $[+1 - 1 - 1 + 1]$ . The value of this contrast can be computed for each subject. The twelve individual interaction effects, reported in table 1, constitute a set of derived data that may be called relevant data for interaction.

**Table 1.** *Reaction time experiment: basic data and relevant data for interaction (times in ms)*

	subject	$a1b1$	$a2b1$	$a1b2$	$a2b2$	Individual interaction effects
	1	387	435	416	473	$387-435-416+473 = +9$
group	2	321	336	343	368	+10
$g1$	3	333	362	358	390	+3
	4	344	430	352	393	-45
	5	368	432	432	504	+8
group	6	357	367	394	411	+7
$g2$	7	336	346	340	421	+71
	8	387	454	438	496	-9
	9	345	408	417	479	-1
group	10	358	389	372	407	+4
$g3$	11	317	375	341	392	-7
	12	386	510	464	513	-75
	mean	353.3	403.7	388.9	437.3	$d = -2.08$ ms $s = 33.28$ ms $t = \frac{d}{s/\sqrt{n}} = -0.22$ [9 df]

Here the relevant data constitute a simple one-way layout, so that the  $A.B$  effect amounts to the overall mean  $\delta$  of the relevant data. Thus it only involves the elementary Bayesian inference about a normal mean, with only two parameters, the true interaction effect  $\delta$  and the within-group standard deviation  $\sigma$  (with the usual assumption of a common variance for the three groups).

Assuming the usual noninformative prior, the posterior distribution of  $\delta$  is a generalized  $t$  distribution. It is centered on the observed mean effect  $d = -2.08$  ms and has scale factor  $e = s/\sqrt{n} = 9.61$ , where  $s = 33.28$  ms is the within-group standard deviation of the relevant data. Remark that  $e$  is the denominator of the usual  $t$  test statistic, that is  $e = \frac{d}{t}$ . In consequence, the posterior distribution of  $\delta$  can be directly derived from  $t = -0.22$ , or again from the usual ANOVA  $F$  ratio  $F = t^2$ . This ensures the technical and conceptual link between Bayesian and

frequentist procedures, in particular the Bayesian interpretation of  $p$ -values and confidence intervals.

The posterior distribution can be interactively investigated by means of visual software. The credibility limits for a given probability, or conversely the probability of a given interval can be interactively computed. Both for  $p$ -values and usual confidence intervals, Bayesian interpretations can be made explicit.

Here the  $t$  test for the  $A.B$  interaction turns out to be “perfectly nonsignificant” ( $p = 0.83$ ): the hypothesis of a null interaction ( $\delta = 0$ ) is “not rejected” by the data. But this is only a negative result, which does not really bring out any evidence in the data that might be positively “in favor of a small interaction effect”. This is clearly illustrated by the Bayesian interpretation of the  $t$  test outcome: there is a 42% ( $100(\frac{2}{5})\%$ ) posterior probability of a positive interaction effect ( $\delta > 0$ ) and a 58% complementary probability of a negative interaction effect ( $\delta < 0$ ). But, since the observed effect can be assessed small, the situation can be perceived as *favorable* to the acceptance of the additive model. Is this impressionistic judgment really justified? Bayesian procedures yield direct statements about the smallness of the interaction effect. We find here that with posterior probability 0.95,  $\delta$  is less than 22.1 milliseconds in absolute value. If one is willing to regard this limit as a tolerable deviation, this statement constitutes a positive statement in favor of the additive model. On the other hand, if the value 22.1 is not deemed to be acceptable, the data can be declared “inconclusive”.

Other analyses can be carried out. In particular prior distributions favorable to the additive model can be investigated. For instance, let us consider a prior such that  $\delta | \sigma^2 \sim N(0, \frac{1}{12}\sigma^2)$  with  $\sigma^2$  still having the usual noninformative prior. We find that with posterior probability 0.95,  $|\delta|$  is less than 14.5 milliseconds. The difference from the value 22.1 found earlier reflects the increase of evidence in favor of the additive model. Alternatively a lump of probability on the hypothesis of a null interaction can be incorporated in the prior.

### 5.3 Teaching Bayesian Methods for Complex Experimental Designs

Three decisive advantages of the specific analysis approach can be stressed. (1) All the traditional analysis of variance procedures can be derived as a direct extension of the basic procedures used in descriptive statistics (means, standard deviations) and inferential statistics (Student’s  $t$  tests). (2) Complex designs involving several factors can easily be handled. In particular, the exact validity assumptions for each inference can be made explicit and comprehensible. (3) Bayesian procedures for assessing the magnitude of effects become straightforward to implement. Furthermore the possibility of teaching Bayesian methods in the context of realistic complex experimental designs is an essential requirement for motivating students.

Statistical computer programs based on the specific inference approach have been developed (Lecoutre, 1996). They incorporate both current practices (significance tests, confidence intervals) and Bayesian procedures. These procedures are applicable to general experimental designs (in particular, repeated measures designs), balanced or not balanced, with univariate or multivariate data, and co-variables.

From an interactive use of the computer programs, a limited set of theoretic

notions is needed to introduce basic procedures, i.e. inferences about one degree of freedom effects in complex designs. An introductory course about descriptive statistics, and elementary inference techniques for the comparison of two means, is generally a sufficient background. Then the attention can be concentrated about the interpretations and the practical meaning of procedures. As a consequence, the principles of advanced techniques can be more easily understood, independently of their mathematical difficulty.

## 6 Conclusion

*“I stopped teaching frequentist methods when I decided that they could not be learned”* (Berry, 1997).

Nowadays Bayesian routine procedures for the familiar situations of experimental data analysis are easy to implement and to teach. They offer promising new ways in statistical methodology (Rouanet *et al.*, 2000). Their results can be taught to non-statistician students in intuitively appealing and readily interpretable form. Using the noninformative Bayesian interpretations of significance tests and confidence intervals in the natural language of probabilities about unknown effects comes quite naturally to students. In return the common misuses and abuses of NHST appear to be more clearly understood. Resorting to computers solves the technical problems involved in the use of Bayesian distributions. This gives the students an attractive and intuitive way of understanding the impact of prior distributions.

In summary, the Bayesian approach appears both desirable and feasible for teaching statistical inference in the context of experimental data analysis. On the one hand, it fulfills the requirements of scientists: objective procedures (including traditional  $p$ -values); procedures about effect sizes (beyond  $p$ -values); procedures for designing and monitoring experiments. On the other hand, it allows to overcome real difficulties, in particular the common misuses of null hypothesis significance tests, and the incorrect interpretations of frequentist procedures.

### REFERENCES

- Berry, D.A. (1991). Experimental design for drug development: a Bayesian approach. *J. Biopharmaceutical Statist.* **1**, 81–101.
- Berry, D.A. (1997). Teaching elementary Bayesian statistics with real applications in science. *Amer. Statist.* **51**, 241–246.
- Brandsttter, E. (1999). Confidence intervals as an alternative to significance testing. *Methods in Psychological Research Online* **4**, 33–46 (Internet: www.mpr-online.de).
- Ciancia, F., Maitte, M., Honoré, J., Lecoutre, B. and Coquery, J.-M. (1988). Orientation of attention and sensory gating: An evoked potential and RT study in cat. *Experimental Neurology* **100** 274–287.
- Clment, E. and Richard, J.-F. (1997). Knowledge of domain effects in problem representation: the case of Tower of Hanoi isomorphs. *Thinking and Reasoning* **3**, 133–157.
- Craig, J.R., Eison, C.L. and Metze, L.P. (1976). Significance tests and their interpretation: An example utilizing published research and  $\eta^2$ . *Bull. Psychonomic*

- Soc.* **7**, 280–282.
- Dignam, J.J., Bryant, J., Wieand, H.S., Fisher, B. and Wolmark, N. (1998). Early stopping of a clinical trial when there is evidence of no treatment benefit: protocol B-14 of the National Surgical Adjuvant Breast and Bowel Project. *Controlled Clin. Trials* **19**, 575–588.
- Falk, R. and Greenbaum, C.W. (1995). Significance tests die hard. The amazing persistence of a probabilistic misconception. *Theory & Psychology* **5**, 75–98.
- Freeman, P.R. (1993). The role of  $p$ -values in analysing trial results. *Statist. in Medicine* **12**, 1443–1452.
- Gigerenzer, G. (1998). We need statistical thinking, not statistical rituals. *Behavioral and Brain Sciences* **21**, 199–200.
- Heldref Foundation (1997). Guidelines for contributors. *J. Experimental Education* **65**, 95–96.
- Hoc, J.-M. (1996). Operator expertise and verbal reports on temporal data. *Ergonomics* **39**, 811–825.
- Hoc, J.-M. and Leplat, J. (1983). Evaluation of different modalities of verbalization in a sorting task. *Internat. J. Man-Machine Studies* **18**, 283–306.
- Kadane, J.B. (1995). Prime time for Bayes. *Controlled Clin. trials* **16**, 313–318.
- Kirk, R. E. (1982). *Experimental Design* (2<sup>nd</sup> edition). Belmont: Brook-Cole.
- Lecoutre, B. (1996). *Traitement Statistique des Données Expérimentales: Des Pratiques Traditionnelles aux Pratiques Bayésiennes* (with Bayesian Windows programs by B. Lecoutre and J. Poitevineau, freely available at Internet address <http://www.univ-rouen.fr/LMRS/Persopage/Lecoutre/pac.htm> or mail to [bruno.lecoutre@univ-rouen.fr](mailto:bruno.lecoutre@univ-rouen.fr)). Montreuil, FR: CISIA.
- Lecoutre, B. (2001). Bayesian predictive procedure for designing and monitoring experiments. ISBA and Eurostat: *ISBA 2000, Proceedings*, to appear.
- Lecoutre, B. and Charron, C. (2000). Bayesian procedures for prediction analysis of implication hypotheses in  $2 \times 2$  contingency tables. *J. Educational & Behavioral Statist.* **25**, 185–201.
- Lecoutre, B., Derzko, G. and Grouin, J.-M. (1995). Bayesian predictive approach for inference about proportions. *Statist. in Medicine* **14**, 1057–1063.
- Lecoutre, B., Lecoutre, M.-P., & Poitevineau, J. (2001). Uses, abuses and misuses of significance tests in the scientific community: won't the Bayesian choice be unavoidable? *Internat. Statist. Rev.* (to appear).
- Lecoutre, B., Poitevineau, J. (2000). Aller au del des tests de signification traditionnels: vers de nouvelles normes de publication. *L'Anne Psychologique* **100**, 683–713.
- Lecoutre, M.-P. (1992). Cognitive models and problem spaces in “purely random” situations. *Educational Studies in Math.* **23**, 557–568.
- Lecoutre, M.-P. (2000). And... What about the researcher's point of view. In H. Rouanet, J.-M. Bernard, M.-C. Bert, B. Lecoutre, M.-P. Lecoutre and B. Le Roux, *New Ways in Statistical Methodology: From Significance Tests to Bayesian Inference* (2<sup>nd</sup> edition), Bern, SW: Peter Lang, 65–95.
- Loftus, G.R. (1993). Editorial comment. *Memory and Cognition* **21**, 1–3
- Murphy, K.R. (1997). Editorial. *J. Applied Psychology* **82**, 3–5.

- Oakes, M. (1986). *Statistical Inference: A Commentary for the Social and Behavioural Sciences*. New York: Wiley.
- Poitevineau, J. (1998). *Mthodologie de l'Analyse des Donnes Exprimentales: Etude de la Pratique des Tests Statistiques chez les Chercheurs en Psychologie, Approches Normative, Prescriptive et Descriptive*. Ph.D. Thesis, Université de Rouen (France).
- Rouanet, H. (1996). Bayesian procedures for assessing importance of effects. *Psychological Bul.* **119**, 149–158.
- Rouanet, H. and Lecoutre, B. (1983). Specific inference in ANOVA: From significance tests to Bayesian procedures. *British J. Math. Statist. Philosophy* **36**, 252–268.
- Rouanet, H., Bernard, J.-M., Bert, M.-C., Lecoutre, B. , Lecoutre, M.-P. and Le Roux, B. (2000). *New Ways in Statistical Methodology: From Significance Tests to Bayesian Inference* (2<sup>nd</sup> edition). Bern, SW: Peter Lang.
- Schmidt, F.L. (1996). Statistical significance testing and cumulative knowledge in psychology: implications for training of researchers, *Psychological Methods* **1**, 115–129.
- Serlin, R.C. and Lapsley, D.K. (1993). Rational appraisal of psychological research and the good-enough principle. In *A Handbook for Data Analysis in the Behavioral Sciences. Vol 1: Methodological Issues* (G. Keren and C. Lewis, eds.). Hillsdale, NJ: Erlbaum, 199–228.
- Spiegelhalter, D.J., Freedman, L.S. and Parmar, M.K.B. (1994). Bayesian approaches to randomized trials. *J. Roy. Statist. Soc. A* **157**, 357–416.
- Thompson, B. (1996). AERA editorial policies regarding statistical significance testing: three suggested reforms. *Educational Researcher* **25**, 26–30.
- Wilkinson, L. and Task Force on Statistical Inference, APA Board of Scientific Affairs (1999). Statistical Methods in Psychology Journals: Guidelines and Explanations. *Amer. Psychologist* **54**, 594–604.

Bruno Lecoutre  
 ERIS, Laboratoire de Mathématiques Raphaël Salem,  
 UMR 6085 C.N.R.S. et Université de Rouen, Mathématiques, Site Colbert,  
 76821 Mont-Saint-Aignan Cedex, France.  
 E-mail: bruno.lecoutre@univ-rouen.fr.  
 Internet: <http://www.univ-rouen.fr/LMRS/Persopage/Lecoutre/Eris.htm>