

ON THE CONCEPT OF P -VALUE

Carlos Alberto de Bragança Pereira and Sergio Wechsler

Universidade de São Paulo
Instituto de Matemática e Estatística
Caixa Postal 20.570 - Ag. Jd. Paulistano
01452-990, São Paulo, Brazil

Summary

Simple examples illustrate how misleading a p -value constructed with no regard to the alternative hypothesis can be. A p -value which regards the alternative hypothesis, called here P -value, is precisely defined. It is shown that the use of the P -value avoids the kind of inconsistencies illustrated by the examples. Although P -values could be considered useless by Bayesians, the use of prior distributions (to obtain weighted likelihoods) is a way by which classical statisticians could regard alternative composite hypotheses when performing significance tests.

Key words: Bayes factor, likelihood principle, null and alternative hypotheses, p -value, significance tests, weighted likelihood ratio.

1. Introduction

Significance testing is a widespread statistical procedure in scientific studies. It consists of the computation of a p -value and then of the judgement of the plausibility of a null hypothesis **H**. Unfortunately, some usual constructions of p -values completely disregard the alternative hypothesis, **A**. This paper

aims to show, by simple examples, how wrong the conclusions based on such misleading p -values can be. A p -value which regards \mathbf{A} , called here P -value, will be precisely defined in Section 4. The P -value is a well-defined (even if not a useful) quantity under the Bayesian view. It is also well-defined under Fisherian and Neyman-Pearson-Wald schools when both \mathbf{H} and \mathbf{A} are simple hypotheses.

Discussion on the concept of p -values have an extremely large reference list – see Good (1983). Most recently, Casella & Berger (1987), Berger & Sellke (1987), Berger & Delampady (1987), and Berger & Mortera (1991) contributed to this list. The subject, however, remains controversial as pointed out by Pratt (1987) and Good (1987). Obviously, after observing data \mathbf{x} and having defined a subset \mathbf{C} of the parameter space, the numbers $\Pr \{\mathbf{C}|\mathbf{x}\}$ and $\Pr \{\mathbf{x}|\mathbf{C}\}$ have completely different meanings. Therefore, comparison of values of posterior probabilities to p -values (as done by Casella & Berger, (1987)) seems to be a very confusing matter. Another confusing aspect of significance tests appears when relations between sample sizes and p -values are made. The recipe of Lindley & Scott (1986) is contrary to the recipe of Peto et al. (1976), as pointed out by Royall (1986). This is briefly discussed in Section 6.

The object of this paper is to illustrate practical aspects of the philosophical point of view of I.J. Good on p -values (pp. 22–55 of Good (1983)). The authors believe that Good's work completely illuminates the matter. However, p -values are still being carefully presented for users of Statistics (see Miettinen (1985), Ch. 9 and Peto et al. (1976)).

References supporting the use of significance tests are Cox (1977), Cox & Hinkley (1974), and Kempthorne & Folks (1971). That one should disregard \mathbf{A} when constructing significance tests are not suggested in these texts. In fact, as Spjøtvoll remarks in his discussion on Cox (1977), explicit regard to \mathbf{A} should be emphasized whenever p -values are being defined.

In many cases the p -value (or even the P -value) turns out to be a tail area. This fact is so common that p -values are often called tail areas. That this

nomenclature is infortunate is shown in Section 4 where we obtain a P -value which is the sum of three areas, one of which is central while the remaining are tail. A P -value that is the sum of three areas was already presented by Good (1983) and called “triple tail” since none of those areas were central.

2. Definition of p -value

The intuitive notion of p -value is captured by the following definition. Consider an experiment producing data \mathbf{x} , an observation of \mathbf{X} , to test a simple¹ hypothesis \mathbf{H} versus \mathbf{A} .

Definition 2.1. *The p -value is the probability, under \mathbf{H} , of the event composed by all sample points that favor \mathbf{A} (against \mathbf{H}) at least as much as \mathbf{x} does.*

However, many textbooks and papers present definitions which make no explicit reference to \mathbf{A} (Cox (1977)) or definitions which completely ignore \mathbf{A} (Freedman, Pisani & Purves (1978) or Pratt & Gibbons (1981)). Typical examples of such definitions are presented next, respectively. Definition 2.2 is in the direction of Mosteller & Rourke (1973) and Definition 2.3 is given in Berger & Sellke (1987).

Definition 2.2. *The p -value is the probability, under \mathbf{H} , of the event composed by all sample points at least as **extreme** as \mathbf{x} is.*

Let now $\mathbf{T} = T(\mathbf{X})$ be a statistic for which large values cast doubt on \mathbf{H} and at point \mathbf{x} it takes the value $\mathbf{t} = T(\mathbf{x})$.

Definition 2.3. The p -value at point \mathbf{t} is the probability

$$p(\mathbf{t}) = \Pr \{ \mathbf{T} \geq \mathbf{t} | \mathbf{H} \} .$$

¹Also for cases where there exist similar regions under \mathbf{H} . The general case is treated in Section 5.

Definition 2.2 presupposes an ordering of the sample points which gives meaning to “extreme”. If the ordering regards **A**, Definitions 2.1 and 2.2 are essentially the same. However, if the ordering disregards **A**, Definitions 2.2 and 2.3 are essentially the same and the p -value can completely ruin the statistical analysis as shown by simple examples such as the ones presented in the following section. It must be remarked that although Berger & Wolpert (1984) present Definition 2.3, in their comments they prescribe that, at least *informally*, **A** must be regarded. A *formal* working definition of the kind of Definition 2.1 is introduced in Section 4.

3. Examples

This section presents three examples of misleading conclusions obtained by the use of Definition 2.3, the one that completely disregards **A**. All the examples concern simple **H** versus simple **A**. To stress the point of this paper simple examples were chosen. However, more standard examples could be constructed at the price of sacrificing simplicity. Section 5 discusses how to deal with composite hypotheses when constructing P -values.

Example A. Consider an urn containing exactly three marbles: one black, one white, and one green. Three marbles were randomly selected from this urn. Consider the following two hypotheses:

H : the selection was done *with* replacement.

vs

A : the selection was done *without* replacement.

Suppose the data consist of the vector $\mathbf{X} = (X_1, X_2)$, where X_1 = number of black marbles in the sample and X_2 = number of white marbles in the sample. Note that the null probability function of \mathbf{X} (multiplied by 27) is displayed in Table 3.1.

Suppose that the point $(1, 1)$ is observed. A person who disregards \mathbf{A} will construct a *p*-value using Definition 2.3 (or Definition 2.2 under an ordering which disregards \mathbf{A}). For any statistic \mathbf{T} which disregards \mathbf{A} , such a *p*-value at $\mathbf{t} = T(1, 1)$ will be $p(\mathbf{t}) = 1$. For instance, if the χ^2 statistic

$$\mathbf{T} = (X_1 - 1)^2 + (X_2 - 1)^2$$

is used, then the *p*-value will be $p(0) = 1$ since $T(1, 1) = 0$. By disregarding \mathbf{A} , no point is more supportive of \mathbf{H} than $\mathbf{x} = (1, 1)$, the vector of expected frequencies under \mathbf{H} . Also note that, under \mathbf{H} , all other sample points have smaller probability (are more “extreme”) than $\mathbf{x} = (1, 1)$ has. Hence, any ordering equivalent to the probability ordering under \mathbf{H} (e.g., the χ^2 statistic above) will produce the unity as the *p*-value. \mathbf{H} would be rejected only upon observation of a point other than $\mathbf{x} = (1, 1)$.

Yet the probability of $\mathbf{x} = (1, 1)$ under \mathbf{A} is one! The conclusion here is that one rejects \mathbf{H} (accepts \mathbf{A}) *only* when observing points which are *impossible* under \mathbf{A} !

Table 3.1

Null probability function (times 27) of \mathbf{X} in Example A

3	1			
2	3	3		
1	3	6	3	
0	1	3	3	1
x_2	0	1	2	3
x_1				

Example B. Let \mathbf{X} be a normal variable with mean zero and unknown variance σ^2 producing data \mathbf{x} . A minimal sufficient statistic here is $\mathbf{T} = T(\mathbf{X}) = \mathbf{X}^2$ with $\mathbf{t} = \mathbf{x}^2$. Using Definition 2.3 (or Definition 2.2) to test $\mathbf{H} : \sigma = 2$, one will compute the *p*-value as the tail area

$$p(\mathbf{t}) = \Pr \{ \mathbf{T} \geq \mathbf{t} \mid \sigma = 2 \} = 2\Phi \left(-\frac{\sqrt{2}}{2} \right) ,$$

where Φ is the standard normal distribution function.

Now let the alternative hypothesis be $\mathbf{A} : \sigma = 1$. Clearly, the tail area under \mathbf{A} is smaller than the one under \mathbf{H} . Indeed,

$$p_{\mathbf{A}}(\mathbf{t}) = \Pr(\mathbf{T} \geq \mathbf{t} \mid \mathbf{A}) = 2\Phi(-\sqrt{\mathbf{t}}) < 2\Phi\left(-\frac{\sqrt{\mathbf{t}}}{2}\right) = p(\mathbf{t}) .$$

Hence, small values of $p(\mathbf{t})$ correspond to even smaller values of $p_{\mathbf{A}}(\mathbf{t})$ and favor \mathbf{H} , not \mathbf{A} . For instance, by following the decision procedure prescribed by Burdette & Gehan (1970), if $\mathbf{t} = \mathbf{x}^2 = 16$, one may wrongly state that there is moderate evidence against \mathbf{H} since

$$p(16) = .0454 < .05 .$$

However, since $p_{\mathbf{A}}(16) \cong 0$, the evidence against \mathbf{A} is much stronger.

Example C. Let f be the density of a random variable \mathbf{X} from which an observation \mathbf{x} is obtained. Consider the significance test of “ $\mathbf{H} : f$ is normal with zero mean and variance 4” versus “ $\mathbf{A} : f$ is the standard Cauchy density”. Suppose again that $\mathbf{t} = \mathbf{x}^2 = 16$. Then, as in Example B, the p -value is $p(16) = .0454$. Note that even with Example B having a different hypothesis \mathbf{A} , the p -value has exactly the same value. Of course, this is because the alternative hypotheses were disregarded.

It must be pointed out that contrary to Example B, here the tail area under \mathbf{A} is $p_{\mathbf{A}}(16) = .156$, which is greater than the p -value. However, as we shall see in Section 5, the value $\mathbf{t} = 16$ again highly supports \mathbf{H} , *not* \mathbf{A} .

In both Examples B and C the sample size is $n = 1$. This is not a restriction of our point and it was used just for simplicity. The same inconsistencies occur in more realistic situations where n is large. Yet, simple hypotheses do not restrict the criticism against p -values. One may consider composite hypotheses as well (Section 5).

4. P -value, a p -value that regards the alternative hypothesis

The examples presented in Section 3 show clearly that Definition 2.3 should be avoided. It is clear that Definition 2.1 is the one to be considered. In this section we present a formalized and workable version of Definition 2.1 which is based on likelihood ratios, following the suggestions of Good (1983). In the next section, by using weighted likelihood ratios, we extend the definition for the case of composite hypotheses.

In the spirit of Definition 2.1, one needs to characterize the set of all sample points that “favour” \mathbf{A} against \mathbf{H} at least as much as \mathbf{x} , the observed data, does. This is done in the following definition by considering the likelihood ratio ordering. Let $f_{\mathbf{H}}$ and $f_{\mathbf{A}}$ be the probability density functions under \mathbf{H} and \mathbf{A} , respectively. The likelihood ratio statistic is denoted by $\mathbf{R} = R(\mathbf{X}) = f_{\mathbf{A}}(\mathbf{X})/f_{\mathbf{H}}(\mathbf{X})$ which takes the value $r = R(\mathbf{x})$ at point \mathbf{x} .

Definition 4.1. *Suppose \mathbf{H} and \mathbf{A} are simple hypotheses and \mathbf{R} is the likelihood ratio statistic. The P -value at point r (or at point \mathbf{x}) is*

$$P(r) = \Pr \{ \mathbf{R} \geq r \mid \mathbf{H} \} .$$

The following result shows the consistency of P -values. On the other hand, p -values lack this kind of consistency.

Lemma 4.1. *For any positive r ,*

$$P_{\mathbf{A}}(r) = \Pr \{ \mathbf{R} \geq r \mid \mathbf{A} \} \geq P(r) .$$

Proof. For convenience we present the proof for the discrete case. Suppose first that $r \leq 1$. We then have

$$\begin{aligned} P_{\mathbf{A}}(r) &= \Pr \{ \mathbf{R} \geq r \mid \mathbf{A} \} = 1 - \Pr \{ \mathbf{R} < r \mid \mathbf{A} \} \\ &= \sum_{y:R(y)<r} f_{\mathbf{H}}(y) + \sum_{y:R(y)\geq r} f_{\mathbf{H}}(y) - \sum_{y:R(y)<r} f_{\mathbf{A}}(y) \\ &= P(r) + \sum_{y:R(y)<r} [f_{\mathbf{H}}(y) - f_{\mathbf{A}}(y)] . \end{aligned}$$

But $R(y) < r$ implies $f_{\mathbf{A}}(y) < rf_{\mathbf{H}}(y)$ and $r \leq 1$ implies $f_{\mathbf{A}}(y) < f_{\mathbf{H}}(y)$. Therefore,

$$\sum_{y:R(y)<r} [f_{\mathbf{H}}(y) - f_{\mathbf{A}}(y)] \text{ is a sum of non-negative terms.}$$

On the other hand, if $r > 1$, Lemma 4.1 follows by noting that $r^{-1} < 1$ and using the result above with the roles of \mathbf{H} and \mathbf{A} interchanged.

Now we apply the P -value to the examples of Section 3. The conclusions will now be consistent and contrary to those obtained when p -values were used.

Example A: (continuation) The statistic \mathbf{R} takes the value $27/6$ at point $(1, 1)$ and 0 at any other point (x_1, x_2) . If point $(1, 1)$ has been observed, the P -value is

$$P(27/6) = \Pr \{ \mathbf{R} \geq 27/6 \mid \mathbf{H} \} = 6/27 < 1 = p(1, 1) .$$

Note that for any other sample point, the P -value is the unity, supporting \mathbf{H} and rejecting, indeed, \mathbf{A} .

Example B: (continuation) Recall that $\mathbf{H} : \sigma = 2$ and $\mathbf{A} : \sigma = 1$, and note that $\{ \mathbf{x} : R(\mathbf{x}) \geq R(4) \} = \{ \mathbf{x} : T(\mathbf{x}) \leq 16 \} = \{ \mathbf{x} : -4 \leq \mathbf{x} \leq 4 \}$. Hence the P -value at point $\mathbf{x} = 4$ (or $\mathbf{x} = -4$) is

$$P(2e^{-6}) = 1 - 2\Phi(-2) = .9546 ,$$

the central area, *not* the tail area p -value. In fact, here $P = 1 - p$.

Example C: (continuation) Recall that “ $\mathbf{H} : \text{Normal}(0, 2)$ ” and “ $\mathbf{A} : \text{Cauchy}(0, 1)$ ”. These two densities are illustrated in Figure 1. Figure 2 introduces the possible values of $R(\mathbf{x})$. From both figures one understands why $\mathbf{t} = \mathbf{x}^2 = 16$ is less “extreme” than $\mathbf{t} = .04$, although the latter has a much higher density than the former.

If the observed sample point is $\mathbf{x} = 4$ (or $\mathbf{t} = 16$), then the set $\mathcal{T} = \{ \mathbf{x} : R(\mathbf{x}) \geq R(4) \}$ is the union of the following three intervals: $(-\infty, -4]$, $[-1, 388, 1.388]$, and $[4, \infty)$.

Figure 1

Standard Cauchy density and Normal density with mean and standard deviation equal to zero and 2 respectively.

Figure 2

Likelihood ratio, $R(x)$, between standard Cauchy and $N(0, 2)$.

The P -value (area under the Normal $(0, 2)$ at set \mathcal{T}) at $\mathbf{x} = 4$ is then $P = .5553 > .0454 = p(16)$.

Now, if $\mathbf{x} = .2$, a point that is less “extreme” in the light of Definition 2.3, the P -value is $P = .0904 < .5553$. This shows that, under Definition 4.1, $\mathbf{x} = 4$ is less extreme than $\mathbf{x} = .2$. Also note that the former P -value is the tail area plus a center slice of the central area, while the latter is the central area plus a small sub-tail of the tail area.

This section shows that the idea of significance level being necessarily a tail area is wrong. P -values can be tail areas, central areas, triple tail areas (as in Good, 1983), triple areas obtained as a sum of central and tail areas, or even a finite sum of areas not all necessarily tail or central.

5. Composite hypotheses

An extension of Definition 4.1 for the case of composite hypotheses is presented in this section. When in presence of composite hypotheses, the statistician usually faces the problem of nuisance parameter elimination. If λ is an indicator function defined as $\lambda = 1$ if \mathbf{H} is true and $\lambda = 0$ if \mathbf{A} is true, then λ clearly becomes the parameter of interest in the significance testing problem. In fact, in the composite hypotheses case, one must use a method of nuisance parameter elimination in order to define $P(r) = \Pr\{\mathbf{R} \geq r \mid \lambda = 1\}$ and $P_{\mathbf{A}}(r) = \Pr\{\mathbf{R} \geq r \mid \lambda = 0\}$. The aim here is to interpret these two quantities as a summary of infinitely many numbers. Those numbers are all the possible P -values (and $P_{\mathbf{A}}$ -values) obtained by consideration of simple sub-hypotheses of \mathbf{H} against simple sub-hypotheses of \mathbf{A} .

The partial likelihood approach is a celebrated way of treating composite hypotheses. For example, the Fisher exact test is an important example of the use of such an approach. However, as shown in Irony & Pereira (1986), it

is possible to improve Fisher's solution by considering a Bayesian tool for a correct P -value evaluation. This Bayesian tool is now used to extend Definition 4.1 to the composite hypotheses case.

Over the two sets of parametric points that characterize \mathbf{H} and \mathbf{A} , consider respectively two probability measures $\Pi_{\mathbf{H}}$ and $\Pi_{\mathbf{A}}$. These two measures define two weighting systems that can be interpreted as conditional priors on parametric points in \mathbf{H} and \mathbf{A} , respectively. Now, for every possible observation \mathbf{x} , define the weighted likelihood values $f_{\mathbf{H}}(\mathbf{x})$ and $f_{\mathbf{A}}(\mathbf{x})$ which are the weighted averages of the likelihood function under \mathbf{H} and \mathbf{A} , respectively; i.e., $f_{\mathbf{A}}(\mathbf{x}) = \int f(\mathbf{x}|\theta) d\Pi_{\mathbf{A}}(\theta)$ and $f_{\mathbf{H}}(\mathbf{x}) = \int f(\mathbf{x}|\theta) d\Pi_{\mathbf{H}}(\theta)$. Also, define the weighted likelihood ratio statistic $\underline{\mathbf{R}} = \underline{R}(\mathbf{X}) = f_{\mathbf{A}}(\mathbf{X})/f_{\mathbf{H}}(\mathbf{X})$ which takes the value $\mathbf{r} = \underline{R}(\mathbf{x})$ at point \mathbf{x} .

Definition 5.1. *Suppose \mathbf{H} and \mathbf{A} are composite hypotheses and $\underline{\mathbf{R}}$ is the weighted likelihood ratio statistic. The P -value at point \mathbf{r} (or at point \mathbf{x}) is*

$$P(\mathbf{r}) = \Pr \{ \underline{\mathbf{R}} \geq \mathbf{r} | \mathbf{H} \} = \int I_{\mathbf{r}}(x) f_{\mathbf{H}}(x) dx ,$$

where $I_{\mathbf{r}}(x)$ is the indicator function of the set $\{x : \underline{R}(x) \geq \mathbf{r}\}$.

It is interesting to note that, by considering the ordering defined by $\underline{\mathbf{R}}$, it would be possible to compute the probability of $\{ \underline{\mathbf{R}} \geq \mathbf{r} \}$ (a p -value) for all elements of \mathbf{H} . The above P -value is the weighted average (using $\Pi_{\mathbf{H}}$) of these probabilities. Analogously, we could characterize the $P_{\mathbf{A}}$ -values using \mathbf{A} and $\Pi_{\mathbf{A}}$. Consequently, the property included in Lemma 4.1 still holds if we consider $f_{\mathbf{H}}$ and $f_{\mathbf{A}}$ as the sample models under \mathbf{H} and \mathbf{A} , respectively. Now, the situation of Example B (and C) is considered under a more realistic situation of composite hypotheses.

Example D: In Example B, let \mathbf{A} be composite. It is not difficult to see that if " $\mathbf{A} : \sigma < 2$ ", then the P -value is a central area. Analogously, if " $\mathbf{A} : \sigma > 2$ ", then the P -value is a tail area. The interesting question is what form the P -

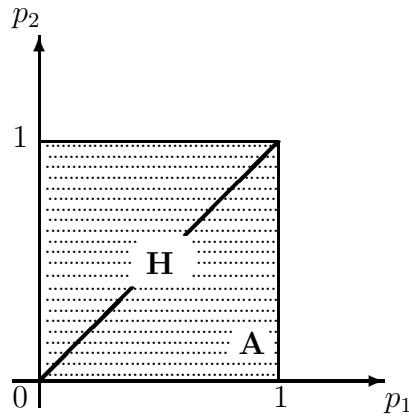


Figure 3

*Hypotheses **H** and **A** in Example D*

value has when “**A** : $\sigma \neq 2$ ”. For simplicity consider that $\Pi_{\mathbf{H}}$ is degenerated at point 2 and $\Pi_{\mathbf{A}}$ is characterized by a one-degree-of-freedom χ^2 (prior) density for $1/\sigma^2$. Simple integration shows that $f_{\mathbf{A}}$ is the standard Cauchy density and therefore the computation of the P -value is reduced to the computation of the P -value of Example C, where a triple area can be obtained.

To end this section, our version of the Fisher significance test for comparing proportions is presented. Note that, although with different dimensions, both hypotheses are composite. Basu (1979) and Pereira & Lindley (1987) present examples questioning the partial likelihood method used in the Fisher significance test. The same examples could be used to support the P -value of Definition 5.1.

Example E: Let $\mathbf{X} = (X_1, X_2)$ where X_1 and X_2 are independent binomial with probability of success p_i and sample size n_i , $i = 1, 2$. After observing data $\mathbf{x} = (x_1, x_2)$ suppose one wants to evaluate the P -value for “**H** : $p_1 = p_2$ ” versus “**A** : $p_1 \neq p_2$ ”. These two hypotheses are illustrated in Figure 3.

If $\Pi_{\mathbf{H}}$ is uniform over the line **H** and $\Pi_{\mathbf{A}}$ is uniform over the set **A**, then,

considering $\mathbf{x} = x_1 + x_2$ and $N = n_1 + n_2$, we have

$$f_{\mathbf{H}} = \frac{\binom{x}{x_1} \binom{N-x}{n_1-x_1}}{\binom{N}{n_1}} \frac{1}{(N+1)},$$

and

$$f_{\mathbf{A}} = \frac{1}{(n_1+1)(n_2+1)}.$$

For the case of $n_1 = n_2 = 5$, Table 5.1 displays all possible values of $\mathbf{R}(x_1, x_2)$ divided by 77. To obtain $f_{\mathbf{H}}$ it is enough to divide the inverted values in the table by 2772. Suppose that $\mathbf{x} = (4, 1)$ is observed. Then $\mathbf{r} = 77/25$ and $P = .0577$. Using now Fisher significance procedure we have $p = .2063$. Note that p is four times the value of P and that Irony & Pereira (1986) show by simulation that the frequency of more extreme (ordering of Table 5.1) points than \mathbf{x} , under \mathbf{H} , is indeed much smaller than p . On the other hand, this frequency is very close to P .

Table 5.1

Values of $\mathbf{R}(\mathbf{x})$ (divided by 77) in Example E

x_2	x_1	0	1	2	3	4	5
0		1/252	1/256	1/56	1/21	1/6	1
1		1/126	1/140	1/105	1/60	1/25	1/6
2		1/56	1/105	1/120	1/100	1/60	1/21
3		1/21	1/60	1/100	1/120	1/105	1/56
4		1/6	1/25	1/60	1/105	1/140	1/126
5		1	1/6	1/21	1/56	1/126	1/252

Note that Definition 5.1 does make use of a prior probability over the parametric space. Then, by considering the full Bayesian probability space, P and $P_{\mathbf{A}}$ are in fact well-defined conditional probabilities. Therefore the computation of these quantities is a non-problem for a Bayesian. However, its

use in significance testing may be questioned by Bayesians since it is believed to violate the Likelihood Principle (Berger & Wolpert, 1984, p.105).

6. *P*-values: a Bayesian look

In this section we discuss the inappropriateness of the use of *P*-values, as a “measure of improbability” of **H**, for Bayesians.

Rejecting **H** whenever the Posterior Odds, *B*, (the ratio of posterior probabilities of **A** and **H**) is large, is equivalent to rejecting **H** whenever the *P*-value is small. In other words, for every $c > 0$, there exists an $\alpha > 0$ such that the event $\{B > c\}$ is equivalent to the event $\{P \leq \alpha\}$. The constant α depends on the prior probability of **H**, on the constant *c*, and on the sampling model. In fact,

$$\alpha = 1 - F_R\left(\frac{\pi c}{1 - \pi}\right) = P\left(\frac{\pi c}{1 - \pi}\right),$$

where F_R is the conditional distribution function of the statistic *R* given **H**, and π is the prior probability of **H**. To obtain this result, note that

$$B > c \Leftrightarrow R(\mathbf{x}) > \frac{c\pi}{1 - \pi} \Leftrightarrow \Pr\{R(\mathbf{X}) \geq R(\mathbf{x})|\mathbf{H}\} \leq \Pr\{R(\mathbf{X}) > \frac{c\pi}{1 - \pi} \mid \mathbf{H}\}.$$

Therefore, $B > c$ is equivalent to $P(\mathbf{R}) \leq \Pr\{R(\mathbf{X}) > \frac{c\pi}{1 - \pi} \mid \mathbf{H}\} = 1 - F_R(\frac{c\pi}{1 - \pi})$.

In this sense, comparison of the actual *P*-value to α can be viewed as a mere computational option in the implementation of a Bayes test. Note that the *P*-values and consequently α may change under a different sampling model, even when the Posterior Odds remains unchanged. Therefore, the use of the *P*-value as a “measure of improbability (Jeffreys, 1961) of **H** on the actual data” is clearly a violation of the Likelihood Principle. Furthermore, using *P*-values in hypothesis testing without any regard to the corresponding Bayes test (which produces the value of the corrected significance level $P(\frac{\pi c}{1 - \pi})$) again violates the Likelihood Principle and is consequently unacceptable to Bayesians. One should also keep in mind that *p*-values (not *P*-values) are unacceptable even

as computational tools, for the above kind of equivalence with Bayes tests no longer holds. Recall the examples of Section 3.

To add more confusion to the subject, Royall (1986) pointed out that it is listed in important literature (Peto et al. (1976) and Lindley & Scott (1984)) two opposite prescriptions for the important role of sample sizes in significance testing.

A given P -value in a large trial is usually stronger evidence that the treatments differ than the same value in a small trial of the same treatments would be. (Peto et al. (1976), p.593)

All significance tests are dubious because the interpretation to be placed on the phrase “significant at 5%” depends on the sample size: it is more indicative of the falsity of the null hypothesis with a small sample than with a large one. (Lindley & Scott (1984), p.1)

To explicate their statement Peto et al. (1976) used a prior probability for \mathbf{H} and computed $\frac{\Pr(\mathbf{H}|\text{significant})}{\Pr(\mathbf{A}|\text{significant})}$, where significant means that the event $\{\mathbf{T} \geq t\}$ has occurred. Hence, in their argument, they did not use the whole information given by the event $\{\mathbf{T} = \mathbf{t}\}$. Had they used the correct Posterior Odds, $\frac{\Pr(\mathbf{H}|\mathbf{T}=\mathbf{t})}{\Pr(\mathbf{A}|\mathbf{T}=\mathbf{t})}$, their conclusion would trivially be the opposite (the one of Lindley & Scott (1984)), as demonstrated by DeGroot (1986, p.380-1). In the context of hypothesis testing, Hodges & Lehmann (1954) also questioned the effect of the sample size in the conclusion of a test. As in DeGroot (1986), they recommend that the choice of a significance level must depend on the sample size.

7. Conclusion

The consideration of \mathbf{A} when defining significance levels is not new (although p -values have always been largely used) and we can refer here to

de Finetti (1972, p.163), Good (1983, p.140), Jeffreys (1961, p.383), Lindley (1978), and Neyman (1981). This paper just shows how dangerous it is not to consider \mathbf{A} in the definition of significance levels. We also present a working definition that does regard \mathbf{A} , even in the case of composite hypotheses. Using prior distributions (or equivalently, weighting systems), the definition of P -value for the case of composite hypotheses is an extension of the one introduced for the case where both hypotheses are simple. Such a P -value is therefore well-defined for Bayesians who can *compute* it with no discomfort (of course, whenever the sample distribution under \mathbf{H} is known). Classical statisticians may also consider the P -value (and even *use* it in the way Burdette and Gehan (1970) use p -values in significance testing) if they accept the idea of replacing, in the likelihood ratio, maximum of likelihoods with corresponding weighted average likelihoods, i.e., if they accept the use of averages for suprema.

Yet, the *use* of P -values other than for computational purposes is unacceptable for Bayesians, since a procedure based exclusively on them would violate the Likelihood Principle. Consequently, it would be incoherent for Bayesians. An inferential procedure that does not violate the Likelihood Principle ought to be based only on the *observed* sample point, not on others (more “extreme”) points that could be observed but were not (Basu, 1975). As pointed out in Section 6, however, the computation of a Bayes Factor can be replaced – probably disadvantageously – by the computation of the corresponding P -value, which is to be appreciated only in the light of the Bayes Factor scale determined by the loss function. Still, this fact does not qualify P -values as Bayesian quantities. Note that Bayes tests are based only on Bayes Factors and these do not change with proportional likelihoods. Since P -values do change, we may well have two Bayes Factors (related to two different models) with equal values corresponding to different P -values.

We close this article with the following quotation from Professor Dennis V. Lindley (1978, p.5):

One can only judge something in relation to the alternatives – a principle that is often not appreciated either in statistics or in politics. It was a great achievement of Neyman and Pearson to recognize this.

Acknowledgements

We wish to thank the referee for many valuable comments that greatly improved the paper. The financial support of CAPES and CNPq is also acknowledged.

(Received June 1993. Revised June 1994.)

References

- Basu, D. (1975). Statistical information and likelihood (with discussion). *Sankhyā A* **37**, 1–71.
- Basu, D. (1979). Discussion of Joseph Berkson's paper "In dispraise of the exact test". *Journal of Statistical Planning and Inference*, **3**, 189–92.
- Berger, J.O. and Delampady, M. (1978). Testing precise hypotheses. *Statistical Science*, **2**, 317–52.
- Berger, J.O. and Mortera, J. (1991). Interpreting the stars in precise hypothesis testing. *International Statistical Review*, **59**, 337–53.
- Berger, J.O. and Sellke, T. (1987). Testing a point null hypothesis: the irreconcilability of P values and evidence. *Journal of the American Statistical Association*, **82**, 112–22.
- Berger, J.O. and Wolpert, R.L. (1984). *The likelihood principle*. Hayward, California: Institute of Mathematical Statistics Monography Series.
- Burdette, J.O. and Gehan, E.A. (1970). *Planning and analysis of clinical studies*. Springfield, Illinois: Charles C. Thomas.

- Casella, G. and Berger, R.L. (1987). Reconciling Bayesian and frequentist evidence in the one-sided testing problem. *Journal of the American Statistical Association*, **82**, 106–11.
- Cox, D.R. (1977). The role of significance testing (with discussion). *Scandinavian Journal of Statistics*, **4**, 49–70.
- Cox, D.R. and Hinkley, D.V. (1974). *Theoretical statistics*. London: Chapman & Hall.
- De Finetti, B. (1972). *Probability, induction, and statistics*. New York: John Wiley.
- DeGroot, M.H. (1986). *Probability and statistics*. 2.ed. London: Addison-Wesley.
- Dickey, J.M. and Lientz, B.P. (1970). The weighted likelihood ratio, sharp hypotheses about chances, the order of a Markov chain. *Ann. Math. Stat.*, **41**, 214–66.
- Freedman, D., Pisani, R. and Purves, R. (1978). *Statistics*. New York: Norton.
- Gibbons, J.D. and Pratt, J.W. (1975). *P*-values: interpretation and methodology. *American Statistician*, **29**, 20–25.
- Good, I.J. (1983). *Good thinking: the foundations of probability and its applications*. Minneapolis: University of Minnesota Press.
- Good, I.J. (1987). Comment on Casella & Berger (1987) and Berger & Sellke (1987). *Journal of the American Statistical Association*, **82**, 125–8.
- Hodges Jr., J.L. and Lehmann, E.L. (1954). Testing the approximate validity of statistical hypotheses. *Journal of the Royal Statistical Society B*, **16**, 261–8.
- Irony, T.Z. and Pereira, C.A. de B. (1986). Exact tests for equality of two proportions: Fisher v. Bayes. *Journal of Statistical Computation and Simulation*, **25**, 93–114.
- Jeffreys, H. (1961). *Theory of probability*. 3.ed. London: Oxford University Press.
- Kempthorne, O. and Folks, L. (1971). *Probability, statistics, and data analysis*. Ames: The Iowa State University Press.
- Lindley, D.V. (1978). The Bayesian approach (with discussion). *Scandinavian Journal of Statistics*, **5**, 1–26.
- Lindley, D.V. and Scott, W.F. (1984). *New Cambridge Statistical Tables*. London: Cambridge University Press.
- Miettinen, O.S. (1985). *Theoretical epidemiology principles of occurrence research in medicine*. New York: John Wiley.
- Mosteller, F.R. and Rourke, R.E. (1973). *Sturdy statistics: nonparametrics and order statistics*. Mass: Addison-Wesley.

- Morrison, D.E. and Henkel, R.E. (1970). *The significance test controversy*. Chicago: Aldine Publishing Company.
- Neyman, J. (1981). Egon Pearson (August 11, 1895 – June 12, 1980). An appreciation. *The Annals of Statistics*, **9**, 1–12.
- Pereira, C.A. de B. and Lindley, D.V. (1987). Examples questioning the use of partial likelihood. *The Statistician*, **36**, 15–20.
- Peto, R., Pike, M.C., Armitage, P., Breslow, N.E., Cox, D.R., Howard, S.V., Mantel, N., McPherson, K., Peto, J. and Smith, P.G. (1976). Design and analysis of randomized clinical trials requiring prolonged observation of each patient, I. Introduction and design. *British Journal of Cancer*, **34**, 585–612.
- Pratt, J.W. (1987). Comment on Casella & Berger (1987) and Berger & Sellke (1987). *Journal of the American Statistical Association*, **82**, 123–5.
- Pratt, J.W. and Gibbons, J.D. (1981). *Concepts of nonparametric theory*. New York: Springer-Verlag New York Inc.
- Royall, R.M. (1986). *The effect of sample size on the meaning of significance tests*. Technical Report # 587. Department of Biostatistics, John Hopkins University. Washington, DC.
- Spjøtvoll (1977). Discussion of Cox's paper "The role of significance testing". *Scandinavian Journal of Statistics*, **4**, 49–70.