**BC Nature** (the Federation of BC Naturalists) represents more than 50 naturalists' clubs and over 6000 members province-wide.

BC Nature strongly recommends that the Review Panel reject the VFPA application for port expansion on Roberts Bank, as posing too great a threat to the wildlife and ecosystems that depend on the habitats of the Fraser delta and Salish Sea.

Attached is a report from Prof. PG Beninger which addresses our concerns.

**Review of Report: 'Biofilm Dynamics during 2018 Northward Migration' (File: 102738-10, January 2019), prepared for Vancouver Fraser Port Authority.**

**Prof. PG Beninger, Université de Nantes, France**

**On behalf of BC Nature**

**Expert Bio**

- Professor, Université de Nantes
- Over 100 publications in refereed primary scientific journals
-  7 book chapters
- Editor/author 'Mudflat Ecology' (Springer Nature Publishers)
- Co-Author, 3 refereed scientific publications concerning WESA and Roberts Bank
- Past Editorial Board, Journal of Shellfish Research, Journal of Experimental Marine Biology and Ecology, Marine Ecology Progress Series
- Current Contributing Editor, Aquatic Biology
- H-index = 39
- Web: http://www.peter-beninger.com/

**Preface**

Having worked on mudflats for several decades, among which the Roberts Bank system, it is my duty to bring to the Panel's attention some of the weaknesses of the most recent report, 'Biofilm Dynamics during 2018 Northward Migration' (File: 102738-10, January 2019), submitted by the Vancouver Fraser Port Authority (the proponent for RBT2.)  These weaknesses undermine the credibility of many of the proponent's conclusions. My review focuses on two broad categories of problems: (1) Statistical Approach, Treatment, and Interpretation, and (2) Other Serious Flaws – considered more briefly, yet of equal concern.

**1. Problems with statistical approach, treatment, interpretation**

On the surface, this report gives the impression of meticulous statistical methodology.  However, for reasons to be outlined below, it will become apparent that the statistical approach, treatment, and interpretation contain flaws which obviate the Proponent's major conclusions concerning environmental effects.

At the outset, it must be noted that the statistical usage in this report does indeed follow much common practice; however, this 'common practice' is a compendium of misuse and misinterpretation which has been decried by the most knowledgeable statisticians for decades, but which, under its own sheer weight, continues to afflict scientific publications today, including this one.  It is not because so many have been wrong, that we are right to imitate them – especially in a context in which the existence of an entire shorebird species is at stake.

*1.1  Conflation of descriptive and experimental approaches*

The use of frequentist statistical methods throughout the report, in various forms of 'Null Hypothesis Significance Testing' (NHST, itself an inadequate term) would situate this study in an experimental context.  However, the basic requirements of experimental protocol (pre-stated null and alternate hypotheses, randomization procedures, and controlled variables) are almost never stipulated.  As with most field work, it appears to be descriptive (also termed observational) in nature, yet it is experimental in the choice of statistical treatment.  This is akin to expecting a bicycle to also operate as a deep-sea submersible.  It is possible to compare data from different sites in a non-experimental context, but it is not valid to use frequentist statistical techniques in a non-experimental context.  In the present study, descriptive methods are indeed used to some extent (e.g. box plots), but these are usually then followed by some form of NHST.  The resulting P-values are then presented as though they were somehow meaningful.  As previously expressed (Beninger et al 2012),

'Many observational studies apply classical hypothesis-testing statistics in their data treatment, and this generates seemingly meaningful numbers, but in reality such studies should rely much more on descriptive statistics and comparisons of effect sizes (Greenland, 1990; Rothman, 1990), regardless of peer and reviewer pressure to the contrary. A second basic reason for eschewing classical statistics in observational studies is the requirement for hypothesis formulation based on a plausible theoretical framework, without which it is impossible to engage in precise interpretations of P-values or confidence intervals (Poole, 2001). Observational studies precede the first experimental studies, since it is impossible to formulate hypotheses when we know nothing at all about the systems studied.' (Beninger et al 2012).

*1.2. The 'truth machine' approach*

An erroneous mindset has evolved within the scientific community, which assumes that the purpose of statistics is to function as an impartial, electronic 'truth machine': we load the data, press the 'return key', and wait for the computer to generate a 'magic' number, the P-value, which will determine for us whether we have proved a point or not.  Such a procedure is found throughout the Proponent's report (e.g. p 25, 27 – in the latter case, the significance fallacy is particularly egregious: 'No difference was documented between replicate samples…p=0.87').  Again, as previously expressed (Beninger et al 2012),

'A mechanistic approach to data treatment has often replaced intelligent data interpretation, and this has been lamented by many statisticians in many fields.

Scientists often feel that they must treat data in a certain stereotypical fashion in order to be taken seriously by their peers (Stewart-Oaten, 1995). There are two very important points to make on this subject: (1) all statistical treatments rely on reasoned judgment, whether the scientist uses it or not, so it is impossible to think of statistics as a simple, blind, 'scale of justice'; (2) a failure to use reasoned judgment in statistical treatment of data is a fundamental abdication of responsibility which calls into question any subsequent conclusions.' (Beninger et al 2012).

## 1.3. Even within the inappropriate experimental framework, problems abound

Beyond the inappropriate experimental framework, many of the 'usual suspects' of experimental statistical misuse are found in the proponent's report, as explained below.

### 1.3.1 Sampling and spatial distribution

The choice of sites for the sampling stations is obviously the crucial underpinning of any field study. In the report, we are told that the chosen sites were known areas where WESA accumulated and fed; this is eminently reasonable. However, within these large areas, the 'replicate biofilm and invertebrate sampling plots' (p. 5) were said to be 'randomly located' (p 5). There are two problems with this: (1) no randomization procedure was given, so we may presume that it was probably 'haphazard' rather than random, and (2) much more importantly, we have here a context where organism spatial distribution is known to be aggregated, at multiple spatial scales, such that even true random sampling can fail to capture the patchiness of the biofilm, meiofauna, and macrofauna (Beninger and Boldina 2018). There is a complete lack of spatial distribution study in this report, without which it is impossible to adequately sample any of these compartments. As so aptly put by Fortin and Dale (2009), '. . .knowledge of the characteristics of spatial structure is almost always the first step to understanding ecological complexity'. Without such baseline knowledge, how could we possibly evaluate the impact of RB2 on the essential features of organism spatial distribution? How could we determine if the patches of this or that organism increased, decreased, or remained the same in size? With respect to biofilm, did this resource become more or less abundant following the (future) construction of RBT2?

> '…readers may well exclaim that they will never have the time or the financial resources to do the often long and painstaking work involved in spatial analyses, especially in situations where the objects of interest are not readily visible, e.g. shipboard sampling of epibenthic organisms, infauna sampling on mudflats, or meiofaunal sampling anywhere. We can only reply that we empathize

completely, but unless they find the necessary time and resources, confidence in their study results will suffer from the failure to have done so. The next-best approach is simply to sample such a large number of randomly-selected points that it is virtually impossible not to capture the spatial pattern. However, this is usually both infeasible and unverifiable.' (Beninger and Boldina 2018)

### 1.3.2. Ambiguity of P-values

As summarized previously (Beninger et al 2012),

'A conventional expression in the medical statistical world is that misinterpretation of P-values (for α, obviously) has killed more people than any other type of scientific misconduct. In marine ecology, it has doubtless been responsible for many errors in ecological interpretation, which have certainly unduly influenced environmental policy, and brought about negative economic consequences. In fact, there is only one thing we may conclude from a P-value: the probability of obtaining the result (or test statistic, generically termed Evidence, E) if $H_0$ were, in fact, true: $P(E|H_0)$.' (Beninger et al 2012).

'In the Fisherian–Pearson amalgam, P is the probability that a given effect could arise if the null hypothesis, or any hypothesis not envisaged in the alternative hypothesis, were true (Type I error), whereas α is the probability value above which we reject our alternative hypothesis, also called the significance level. Although most statistics textbooks do vaguely mention that the significance levels of 0.05, 0.01 etc. are subjective, workers in marine ecology often recognize the 0.05 level as a sort of mechanical Occam's razor: 0.05, reject null hypothesis, ≥0.05, accept null hypothesis. Statisticians have decried this reasoning for decades (Cohen, 1994; Fisher, 1959; Gelman and Stern, 2006; Gigerenzer, 2004; Gigerenzer et al., 2004; Hubbard and Bayarri, 2003; Stephens et al., 2005; and the numerous references cited in these works). Quite apart from the fact that we cannot accept a null hypothesis (this is tantamount to 'proving' a null hypothesis, which is simply not possible), it is argued that the onus is on the researcher to establish an α in line with his/her 'evidence and ideas' (Fisher, 1959), or at least in line with the 'risk posed to science or society of false positive or false negative results' (Mapstone, 1995; Stephens et al., 2005). While this may be possible in some research fields (especially medicine), it either does not make much sense, or is impossible to accomplish, in most marine ecological studies. We must therefore be open to the possibility of significance levels different from 0.05, and to the justification for such levels, especially with respect to the minimization of Type 2 errors (i.e., the probability of accepting $H_0$ when it is in

fact false; see below). An α= 0.1 might be a sufficient Type 1 error level for concluding a difference in coloration of reef fish, for example, while 0.05 may be required to conclude that there is a difference in heavy metal concentration in species of fish consumed by humans.  (Beninger et al 2012).

The 'Fisherian–Pearson amalgam' is the approach to NHST implicitly followed in the Proponent's study; in itself, it is flawed because it conflates two distinct approaches, such that it satisfies the conditions of neither.  Even within this flawed framework, the preceding excerpt makes an important point with respect to the report: it is incorrect to 'accept the null hypothesis' at any P-value; we can only say that it is likely that the null hypothesis is true.  A large P-value indicates that there is a strong probability of the observed result being obtained either because the null hypothesis is true, *or due to some other, uncontrolled/unforeseen factor*.  It is, therefore, a leap of logic to interpret high P-values as indicating that two or more test statistics, and more importantly the statistical descriptors (e.g. means, slopes), are not 'significantly different', and to conclude from there, that there is no *appreciable* difference (see 'effect size' below).  One of many examples is found on p 27: 'Modelling found no difference in chlorophyll *a* abundance among sampling dates over the duration of the study ($F_{5,64}$= 2.2, p= 0.07)' – in this case, compounding the significance fallacy is a flagrant example of autocorrelation.  PoThe overall, unstated null hypothesis of this study is that the construction of RBT2 will have little effect on the migrating WESA population.  Of particular pertinence in this regard is the following (Beninger et al 2012):

'…lack of statistical significance could be due to the null hypothesis being true, but it could also be due to a myriad of other causes which we have not identified. So even a very general null hypothesis such as 'there is no effect' may not be meaningful — there may not have been an effect because of an interfering variable for which we have not controlled because we could not imagine its existence, and without which the effect would be manifest at the pre-determined α level. There is no reason to suppose that there is NO effect when P exceeds the chosen significance level, particularly when the data do point in the direction of the effect. The automatic acceptance of the null hypothesis when P exceeds the significance level is a widespread form of 'corrupt science' in marine ecology, which, paradoxically, its practitioners often consider one of the highest and most rigorous forms of science (Carver, 1978).'

*1.3.3. Statistical power*

In view of the importance of the null hypothesis in the present study, the most important error to be minimized is Type 2: accepting the null hypothesis when in fact the alternate hypothesis (considerable effect) is true.  The corresponding error probability term is 'β'.

' Reducing β is called increasing the power of the statistical test, defined as 1−β (note that this still does not allow us to calculate the value of β). Without increasing α, there are only 2 avenues available for increasing power: reducing the variability of the data (e.g. re-doing the experiment with more efficient instrumentation or methodology, when possible), or increasing sample size (see e.g. Green, 1989 for a discussion of the determination of the necessary N to achieve a desired power level for the detection of a given response magnitude in pollution impact studies). Both options are usually associated with increased material costs. However, since β is inversely proportional to √N, relatively large sample size increments translate to much more modest gains in power (reductions in β). Keeping in mind the potential gravity of Type 2 errors in marine ecology, and the difficulty of increasing power by either reducing data variability or increasing sample size, it is therefore clear that in many cases the optimal compromise would be to increase the level of α, e.g. doubling it to 0.1 (and therefore increasing the risk of a Type 1 error), as this will automatically increase statistical power (Peterman, 1990), without incurring any additional material costs. Obviously, such decisions can only be made if we have some knowledge of the relative consequences of Type 1 and 2 errors for each particular study (informed judgment once again).'

'Despite the lack of preoccupation with this error source, the consequences of insufficient attention to β may be very important; a striking example in the medical field was given by Streiner (1990), and in the field of marine environmental research, it has been argued that the consequences of a Type 2 error are usually even more serious, and certainly more pernicious, than those of a Type 1 error (Fairweather, 1991; Mapstone, 1995; Peterman, 1990)…' (Beninger et al 2012)

**Despite the primordial importance of reducing β in this study, the term is never mentioned, and no procedure to do so is even described.**

*1.3.4. Statistical significance vs Biological / Ecological significance*

At numerous points in the Proponent's report, reference is made to 'statistically significant' differences (e.g. pp 26, 28, 32, 34, 37 to point out just a few).  The term 'statistically significant' is such a semantic minefield, that it has been decried by statisticians for decades, and recently over 800 statisticians from around the world signed a paper  in Nature advocating that this term be used no longer (Amrhein et al 2019).  Setting aside the gross semantic inadequacy of that expression, which routinely leads to misinterpretation on the part of readers (even when that is not the intention of the authors), it must be emphasized that:

> 'A statistically significant difference may be obtained from data entirely bereft of biological significance. It is absolutely fundamental that we distinguish between statistical significance and biological significance (Johnson, 1999; Mapstone, 1995; Stefano et al., 2005; Yoccoz, 1991). This is one of the most important, and most-neglected, concepts in marine biology/ecology.' …'In other words, P-values merely specify whether or not an effect is likely to exist (Stefano et al., 2005); they give no information at all on the magnitude of the effect (effect size)' (Beninger et al 2012)

Conversely,

> 'In the absence of a statistically-significant difference for a given effect between groups, we may conclude that the groups are equal with respect to this effect. This is perhaps the most egregious misinterpretation of statistical non-significance. Effect equality may only be statistically ascertained, within specified limits of probability, using tests of equivalence and noninferiority… Moreover, such tests rely on reasoned judgment (e.g. in deciding what a significant effect size is), so they are not the objective statistical razors which many biologists so earnestly — and naïvely — seek.' (Beninger et al 2012)

Compared to the numerous references to 'statistical significance' in this report, there is scarcely a handful of references to effect size (e.g. p 83, point 5), and we note the near-total absence of pertinent ecological context within which such effects may be evaluated.  This is likely due to lack of information about the ecological context, because the large data base to date is, when compared with the amount of heretofore unknown knowledge, sketchy at best.  Yet this is precisely the point: the lack of knowledge about the mudflat ecosystem so hobbles our attempts to predict outcomes such as the potential effects of RBT2, that we are unable to present the arguments which really matter: how great an effect, and is that much serious, slight, or somewhere in-between?

### 1.4. Correlation-autocorrelation

At several points in the report, the authors proceed with statistical tests and correlations which rely on independence of data points, when in fact the data points are autocorrelated. Examples include the obvious temporal autocorrelation on p 81, as well as many of the 'significant positive correlations' on p 38. Of course the SFA, MUFA, and PUFA components will show high correlations with the Total Fatty Acid fraction (Table 3.7), since they are autocorrelated! Proceeding with such tests, when autocorrelation is obvious at the outset, diminishes the credibility of the report.

### 1.5. Conclusion: statistical approach, treatment, and interpretation

The inadequacy of the NHST approach, in the way in which it is applied in the Proponent's report, has been repeatedly pointed out by the world's leading statisticians over the past several decades (Beninger et al 2012):

'Beyond the problems associated with misunderstanding, misinterpretation, and misuse of $\alpha$ and P-values, the problems with NHST have been most succinctly and elegantly portrayed (within a literature singularly graced with eloquence), by Stephens et al. (2007). To sum up, NHST is, in most cases, a very inappropriate tool used in very inappropriate ways, to achieve a misinterpreted result. The driving force behind its use is the belief that it is a totally objective, mechanical procedure which will reveal objective truth precisely because we use it in this fashion. Not only is this obviously not the case, but there is no alternative, totally objective, mechanical procedure which will reveal objective truth in any approach, as has been eloquently underscored by several statistical luminaries, notably Jacob Cohen (1994).'

In other words, the best we can do (and what the Proponent should have done from the outset) is to first adequately determine the spatial distributions of the objects of interest (in this study, all living components), in order to properly sample them. We must then decide whether or not our study is truly experimental in nature (which this study is not), and if it is, and we opt for the frequentist statistical framework in much of the data treatment, we must properly enunciate the hypotheses to be tested prior to gathering the data, present the evidence (results), properly assess the probability of these results arising if the null hypothesis were true (Type 1 error), reduce as much as possible the probability of Type 2 error, and situate any statistical result in the context of effect size and biological/ecological significance. **As detailed in the preceding, none of these conditions were met in the present study. In view of the problems outlined above regarding statistical approach, treatment, and interpretation, the validity of the study conclusions are considerably weakened, and should not be**

**considered 'solid science'.** For these reasons alone, the Panel cannot formulate a recommendation based on the Proponent's report.


## 2. Other serious flaws

### 2.1. The siren call of models

The Proponent's report makes extensive use of modelling (see Appendix E). We must never lose sight of the famous quote by George Box, one of the greatest statistical minds in history: 'All models are wrong, but some are useful' (Box 1976). Modelling works best when it is confined to past events; it can be egregiously wrong when it is used to predict the future. As I write this, I have seen the 48-hour weather predictions for Nantes change 3 times in as many days. Weather is probably the most-studied system on Earth, yet the models elaborated to predict even this purely abiotic system can fail over time periods of as little as 48 or even 24 hours. The reliability of model predictions obviously decreases with the future time interval. Those of us who have worked in the fisheries sector know how tenuous are our models in that domain, where biological variables and long time intervals considerably complicate the modelling landscape. Models which require large numbers of parameters are not better than simple models; they are in fact much worse (Box, 1976). Yet the complexity of ecological systems defies their description and prediction using few parameters. If this sounds paradoxical, it is, and that is why attempts to predict the ecological future are largely akin to crystal ball gazing. It is not something we should be doing when the fate of an entire population of shorebirds is at stake.


### 2.2. The unexpected

Predictions based on data gathered in the past and present are contingent upon (i) the certainty that we have adequate data from all possible sources which could influence the future of the system, and (ii) the certainty that no any unexpected or idiosyncratic event will impact the system. Obviously, while the first of these conditions is partly under our control, the second is not. Yet even the first condition is not satisfied in this report. The authors concentrate on the fatty acids as the crucial variable, and while they are undoubtedly important, they may not be the only specific substances required by WESA at this point in their migration. The fact that their diet changes so radically over the course of the migration (Beninger et al 2011) tells us that their nutritional requirements have changed (there are ample small invertebrates available, just as there were at the beginning of the migration, but these are not targeted at Roberts Bank). Epibenthic microbial mats are temporarily-stabilized microbial soups, containing a vast array of molecules. WESA may require inputs of specific vitamins, cofactors, cation ratios, trace elements, potential hormone precursors, etc. at this penultimate stop prior to nesting, and these may be more abundant at specific sites

than others, since the mudflat is a microbial mosaic.  At present, we simply lack the information needed to answer this question, and **we have no idea how an RBT2 would affect the supply of such resources.**

*2.3. Worrying lapses*

At several points in the report, one is confronted with statements which are such basic errors, that confidence in the rest of the report is greatly diminished.  At the outset (p 1), we read that 'The Fraser River Estuary (FRE), comprised of Roberts Bank, Sturgeon Bank, and Boundary Bay, is the largest estuary on the Pacific coast of North America…'  It is painfully obvious that this is in fact the FRE *outfall*; the estuary actually comprises all of the Fraser River inland with salinities greater than that of freshwater.  It is perplexing to read such an erroneous basic statement, for which even an undergraduate student would be sanctioned.

Similarly, on p 5 we read 'Sampling was initiated…providing enough time for biofilm to migrate to the surface'.  We are surprised to discover that these authors believe that biofilm 'migrates'.  While certain elements of the microphytobenthos can and do migrate as a function of the tidal cycle, one cannot say that the *biofilm* migrates.  This betrays a basic lack of understanding about what is meant by 'biofilm'; **a worrying state of affairs, since virtually the entire report is concerned with biofilm.**

A third such lapse is found in Table D8, which purports to present 'Summary Statistics of All Benthic Invertebrates Components', yet actually presents results for the 'Biofilm component'.  The error brings into question whether the authors understand that invertebrates are not equivalent to biofilm.

## 3. Conclusion

The preceding brief critique of the Biofilm Report is not a total rejection of the science contained therein; rather, I expect that some good primary scientific publications might be retrievable from this work, suitably amended.  However, for the reasons outlined above, I do not believe that the Proponents have presented a scientifically credible case, or in fact any case at all, for the innocuity of the proposed RBT2 project with respect to the Roberts Bank mudflat ecosystem, and the biofilm-shorebird dimension in particular.  I recommend that the Panel reject the Proponent's proposal for RBT2.

## References

Amrhein V, Greenland S, McShane B (2019)  Scientists rise up against statistical significance.  Nature 567: 305-307.

Beninger PG, Boldina I. 2018.  Quantitative considerations in mudflat ecology.  In PG Beninger (ed) Mudflat Ecology, Springer Nature, Cham, Switzerland, 389-419.

Beninger PG, Decottignies P, Elner RW. 2011.  Downward trophic shift during breeding migration in the shorebird *Calidris mauri* (Western sandpiper).  Mar. Ecol. Prog. Ser. 428: 259-269

Box GE. 1976. Science and statistics. J. Am. Stat. Assoc. 71 (356), 791–799.

Carver RP. 1978. The case against statistical significance testing. Harv. Educ. Rev. 48: 378–399.

Cohen J. 1994. The Earth is round (pb.05). Am. Psychol. 49: 997–1003.

Fairweather PG. 1991. Statistical power and design requirements for environmental monitoring. Aust. J. Mar. Fresh. Res. 42: 555–567.

Fisher RA. 1959. Statistical Methods and Scientific Inference, 2nd ed. Oliver and Boyd, Edinburgh

Fortin MJ, Dale M. 2009. Spatial autocorrelation in ecological studies: a legacy of solutions and myths. Geogr Anal 41: 392–397

Gelman A, Stern H. 2006. The difference between 'significant' and 'not significant' is not itself statistically significant. Am. Stat. 60: 328–331

Gigerenzer G. 2004. Mindless statistics. J. Socio-Econom. 33: 587–606.

Gigerenzer G, Krauss S, Vitouch O. 2004. The null ritual: what you always wanted to know about significance testing but were afraid to ask. In: Kaplan, D. (Ed.), The Sage Handbook of Quantitative Methodology for the Social Sciences, Sage Ltd.,Thousand Oaks, CA, pp. 391–408.

Green RH 1989. Power analysis and practical strategies for environmental monitoring. Environ. Res. 50: 195–205.

Greenland S. 1990. Randomization, statistics, and causal inference. Epidemiology 1: 421–429.

Hubbard R, Bayarri MJ. 2003. Confusion over measures of evidence (p's) versus errors (α's) in classical statistical testing. Am. Stat. 57: 171–182.

Johnson DH. 1999. The insignificance of statistical significance testing. J. Wildl. Manage. 63: 763–772

Mapstone BD, 1995. Scalable decision criteria in environmental impact assessment: effect size, type I, and type II errors. In: Schmitt, R.J., Osenberg, C.W. (Eds.),

Detecting Ecological Impacts: Concepts and Applications in Coastal Habitats. Academic Press, New York, pp. 67–80

Peterman R. 1990. Statistical power analysis can improve fisheries research and management. Can. J. Fish. Aquat. Sci. 47: 2–15.

Poole C. 2001. Low P-values or narrow confidence intervals: which are more durable? Epidemiology 12: 291–294.

Rothman KJ. 1990b. Statistics in nonrandomized studies. Epidemiology 1: 417–418.

Stefano JD, Fidler F, Cumming G. 2005. Effect size estimates and confidence intervals: an alternative focus for the presentation and interpretation of ecological data. In: Burke, A.R. (Ed.), New Trends in Ecology Research. Nova Science Publishers, New York, pp. 71–102

Stephens PA, Buskirk SW, Hayward G, Del Rio, CM. 2005. Information theory and hypothesis testing: a call for pluralism. J. Appl. Ecol. 42: 4–12.

Stephens PA, Buskirk SW, Hayward GD, Del Rio CM. 2007. Inference in ecology and evolution. Trends Ecol. Evol. 22: 192–197.

Stewart-Oaten A. 1995. Rules and judgments in statistics: three examples. Ecology 76: 2001–2009.

Streiner DL. 1990. Sample size and power in psychiatric research. Can. J. Psychiatry 35: 616–620.

Yoccoz NG. 1991. Use, overuse, and misuse of significance tests in evolutionary biology and ecology. Bull. Ecol. Soc. Am. 72: 106–111.