

Editorial: Evidential Statistics, Model Identification, and Science

Mark L. Taper^{1*}, Jose M. Ponciano², Yukihiko Toquenaga³

¹Montana State University System, United States, ²University of Florida, United States, ³University of Tsukuba, Japan

Submitted to Journal:

Frontiers in Ecology and Evolution

Specialty Section:

Environmental Informatics and Remote Sensing

Article type:

Editorial Article

Manuscript ID:

883456

Received on:

25 Feb 2022

Journal website link:

www.frontiersin.org

Conflict of interest statement

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest

Author contribution statement

All authors contributed to the conception and writing of this editorial.

Keywords

replication crisis, Scientific Inference, Model misspecification, inference bias reduction, Statistical foundations

Contribution to the field

Statistics is arguably the most powerful of all scientific instruments. For the last century, statistics has been dominated by two alternative approaches: Error statistics and Bayesian statistics. Unfortunately, both approaches suffer from technical and philosophical problems. These problems create biases in scientific inference and also lead these approaches to misrepresent the uncertainty in scientific inference leading to the replication crisis in science. We believe that the evidential approach can provide a correction to statistics. Evidential statistics is a cluster of statistical methods and approaches being developed to meet a set of desiderata or meta-criteria that were selected so as to impose desirable inferential properties on those methods.

Editorial: Evidential Statistics, Model Identification, and Science.

1

2 **Mark L. Taper^{1,2*}, José M. Ponciano², Yukihiko Toquenaga³**

3

4 ¹Montana State University, Department of Ecology, Bozeman MT, 59717, USA.

5 ²University of Florida, Department of Biology, Gainesville FL, 32611-8525, USA.

6 ³Faculty of Life and Environmental Sciences, University of Tsukuba, Tsukuba, Japan

7

8 *** Correspondence:**

9 Mark L. Taper

10 MarkLTaper@gmail.com

11

12 **Keywords:** replication crisis, scientific inference, model misspecification, inference bias reduction

13 1 Why this research topic

14 We have undertaken this research topic for several reasons: First to promote and disseminate the ideas and techniques of evidential statistics
15 to ecologists and evolutionary biologists so that their research might benefit increased clarity that evidential thinking engenders. And, second
16 to encourage statisticians to think how their own work relates to this emerging approach to the fundamental problems of statistics.

17 2 How to read this volume

18 Selecting an optimal order to read the papers of this research topic requires decisions on the part of the reader. The papers are not
19 ordered in any developmental fashion, but simply by the order that they were first published. Another difficulty is that there are two target
20 audiences for this research topic: First, quantitative scientists, primarily ecologists and evolutionary biologists, who might wish to apply

21 evidential thinking to their own research; and second, statisticians who might be interested in furthering the technical development of
 22 evidential statistics.

23 Table 1 lays out the primary themes considered in each paper and identifies authorship abbreviations. Those readers who would like
 24 to begin with statistical principles, then move to applications, and then conclude with more philosophical considerations might read the topic
 25 in the order of DPT&L, P&T, L_b, TLPD&J, S&T, M&S, FTZJ&M, CC&H, T&G, S&B, JKECS&T, L_a, B&B, S&H. For readers who might
 26 prefer to begin with philosophy, and move to application, and then finish with technical details, a reasonable order might be: B&B, S&H,
 27 JKECS&T, T&G, L_a, S&B, FTZJ&M, CC&H, DPT&L, P&T, L_b, TLPD&J, M&S, S&T.

28 <Table 1 near here>

29 **3 What is evidential statistics**

30 Statistics is arguably the most powerful of all scientific instruments. For the last century, statistics has been dominated by two alternative
 31 approaches: Error statistics¹ and Bayesian statistics. Unfortunately, both approaches suffer from technical and philosophical problems (see
 32 Taper and Ponciano, 2016 for discussion). These problems make the instrument of statistics like the Hubble telescope before its optics were
 33 corrected in 1993: A fantastic tool not living up to its full potential.

34 We believe that the evidential approach can provide a similar technical correction to statistics. Evidential statistics is a cluster of statistical
 35 methods and approaches being developed to meet a set of desiderata or meta-criteria that were selected so as to impose desirable inferential
 36 properties on those methods (see JKECS&T for a list of desiderata).

37 The central question for evidence is simple: Which of two models of reality is better supported by the data? More technically, evidence is a
 38 data-based estimate of the difference of the divergences of each of the distributions implicit in two models to the data distribution resulting
 39 from an unknown true generating process (see Lele, 2004; and TLPD&J). Several salient features of the evidentialist perspective are
 40 immediately obvious: First, evidence is comparative, second, neither model is given a favored status, and third, that a “true” model is not
 41 assumed to be in the model set.

42 These guiding principles allows evidential statistics to draw on and refine elements from error statistics, likelihoodism, Bayesian statistics,
 43 information criteria, and robust methods, evidential statistics to create an approach that smoothly incorporates model identification, model

¹ By error statistics we mean that subcategory of frequentist statistics that uses error probabilities as the primary inferential quantity including Fisherian significance, null hypothesis significance testing, Neyman-Pearson hypothesis testing, and severe testing. The term classical statistics is sometimes applied to this grouping, but this can be considered a misnomer as Bayesian statistics predates these methods considerably.

44 uncertainty, model comparison, parameter estimation, parameter uncertainty, pre-data control of error, post-data assessment of uncertainty,
45 and post-data strength of evidence into a single coherent framework.

46 **4 Some implications of evidential statistics for science**

47 The implications of evidential statistics for science are manifold. For brevity, we focus here on the impact an evidential approach could
48 have on the replication crisis (Pashler and Wagenmakers, 2012). The replication crisis presents a profound challenge to both statistics and
49 science. As more replication of scientific studies is attempted, it is being found that studies tend not to replicate at their nominal rates. This
50 is undermining both trust in statistics by scientists and trust in science by the general population.

51 Virtually all models are to some degree misspecified (see TLPD&J for a technical definition of “misspecified”). Misspecification in itself is
52 not a bad thing. A true model would be enormously complex and would be neither comprehensible nor estimable. What is dangerous is
53 inference that doesn’t acknowledge misspecification. With Neyman-Person Hypothesis testing (NPHT), error rates become distorted when
54 both models are misspecified. Error rates can be less than, equal to, or greater than their nominal rates (DPT&L) making nominal rate
55 replication extremely unlikely. Furthermore, under some reasonable model space geometries, a NPHT will select the wrong model with
56 probabilities that go to 1 as sample size increases (DPT&L). In contrast, evidential model selection reliability seems in simulation to be
57 estimated unbiasedly (Taper et al. 2019) and all evidential error rates go to 0 as sample size increases (DPT&L).

58 None of Fisherian significance (FS), null hypothesis significance tests (NHST), or NPHT can produce evidence for the null model
59 (DPT&L). This is problematic because often it is the null which of scientific interest. Statisticians teach that “absence of evidence is not
60 evidence of absence”, but the need of scientists to say something about the null model forces this warning to be often ignored. In evidential
61 statistics reference and alternative models are always correctly treated symmetrically (DPT&L, TLPD&J, JKECS&T) for inference,
62 although this does not imply that decision thresholds need to be symmetric.

63 When scientists, reviewers, and journals that do recognize that FS, NHST, and NPHT do not produce evidence for the null, a common
64 response is publication bias, the tendency not to publish studies with attained P-values less than 0.05 (Franco et al. 2014). This “file drawer
65 problem” creates several biases in the literature. First, of course, is the lack of studies showing evidence for the null. More insidiously,
66 because all tests are stochastic, a number of studies are published falsely showing significant evidence for the alternative (Type I errors).
67 These are not balanced in the literature by the many studies in the file drawer.

68 The immense pressure on scientists to publish leads many, intentionally or unintentionally, into questionable research practices to avoid the
69 file drawer problem. One of these is “cherry picking”, the retroactive selection of data and/or statistics so as to achieve significance
70 (Ioannides, 2019). Another is HARKing, Hypothesizing After Results are Known (Kerr, 1998). Both have drastic effects on the replication
71 crisis.

72 Evidential analysis gives scientists statistically correct language (TLPD&J) to speak about strong evidence for the null versus the alternative,
 73 strong evidence for the alternative versus the null, and evidence that doesn't clearly distinguish between the two models. All of which are of
 74 scientific interest. Even results that can't distinguish between models tell us where more data is needed. The results of any well-designed
 75 scientific study now have meaning and could potentially be publishable—regardless of significance.

76 Undertaken in an evidential statistics context, HARKing is a legitimate and even beneficial practice (Taper and Gogan, 2002). The evidence
 77 in HARKing has always been clear, although estimation of the uncertainty remained a problem (Taper and Lele, 2004). Bootstrapping of
 78 evidential comparisons now improves the understanding of the uncertainty of even HARKed results (Taper & Lele 2011, Taper et al. 2019,
 79 TLPD&J).

80 **5 Comments on the articles:**

81 **5.1 Shimodaira and Terada (S&T)**

82 At the heart of ecology is a search to better understand and characterize the relationship between species as well as that of a group of species
 83 and their environmental variables. On the other hand, a central topic in evolutionary studies is inferring the ancestral relationships of a set of
 84 extant species. In both cases, graph theory has become the theoretical foundation upon which the biological edifices in these two fields are
 85 constructed. In ecology, species are thought as nodes in a diagram and the relationships between species are represented as edges uniting
 86 any two nodes. In evolution, a phylogenetic binary tree is a diagram representing the evolutionary relationships among a set of extant
 87 species, which are shown as the tips (leaves) of the tree. Each interior node in the tree connects with three other nodes: two descendants and
 88 one ancestor.

89 The binary phylogenetic trees are called bifurcating trees because there are two branches leading out from each interior node. Proceeding
 90 from the present-day species of interest backwards in time under this binary framework eventually leads to a common ancestor, the root of
 91 the tree. In that context, one particular “tree topology” is one specific construction of the possible set of relationships among the species of
 92 interest and represents a single hypothesis about the ancestral relationships between these species, all the way back to their most recent
 93 common ancestor. How many such hypotheses can one posit with n species? With two species the answer is one, with three species the
 94 answer is three, with four it's fifteen, with five it's one hundred and five and in general, with n species it's $(2n - 3)!/(2^{n-2}(n - 2)!)$. For
 95 example, for six species, the number considered by S&T, one could posit 945 such trees.

96 In such setting, it quickly becomes obvious that good treatments of the statistical problems of multi-model selection and multiple hypotheses
 97 testing are key to making any progress in this area. Previously, the leading approach to deal with the problem of selecting among these
 98 models (hypotheses) the best representation of reality used NHST. This body of work was started by Kishino and Hasegawa (1989), and
 99 continued by Shimodaira (1998), Shimodaira and Hasegawa (1999) and Shimodaira (2002). S&T now goes one step further and provides a
 100 novel methodology of shifting the phylogenetics question away from: “is a newly estimated tree topology significantly similar to the
 101 unknown, true species topology?” and instead ask: “from this set of models, which tree topology and group of models are significantly

108 $H_0 : \mu \in R$ versus $H_1 : \mu \in R^c$), the tests are not standard NP tests. Truth does not lie in either
 109 hypothesis, but instead is being projected onto the manifold $R \cup R^c$. Further, the pseudo data being used to generate the distribution of the
 110 test statistic does not come from H_0 , but is generated by a non-parametric bootstrap. Thus, the difference between the inference in S&T and
 111 TLPD&J may be little more than the statistics they choose to present.

112 5.2 Scheiner & Holt (S&H)

113 This paper takes the readers out of the weeds and forces them to look simultaneously at the trees and the forest. Deeply informed by both
114 the history and the philosophy of science, the manuscript points out that evidential statistics formally only deals with the relationships
115 among models and data; S&H then ask how evidential statistics can inform either the generation or the support for general and constitutive
116 theories. Clearly it can because Peirce's abduction (Peirce, 1974) can be thought of as a conceptual adequacy measure for models,
117 hypotheses, or theories, while modern abduction, i.e. inference to the best explanation (Haig, 2009) can be thought of as conceptual
118 evidence for the same.

119 In an analogy to biological evolutionary theory, S&H discuss how model selection, an evidential process, can act as a selective force to
120 winnow the models included in constitutive theories. S&H further suggest that pattern matching as well as Whewell's consilience and
121 coherence (Forster and Wolfe, 1999) might possibly be utilized in formal procedures for quantifying the evidence supporting one theory
122 over another.

123 Despite the excellence of this article, S&H do sin against science in suggesting that sometimes statistics is not necessary². They claim for
124 instance that if something never occurs then no statistics is necessary. To which a statistician would query, “never occurs in how many
125 trials?” The evidential impact of something never occurring is very different in experiments of 1 trial, 4 trials, or 8 trials (see JKECS&T).
126 Because they are writing as theoreticians, S&H’s sin is only venal. For theoreticians, statistics and even data, are always optional. The job
127 of theoretical science is to construct alternative internally consistent possible worlds. The job of empirical science is to determine which of
128 those possible worlds best describes the real world—and for that, statistics is always needed.

² In prepublication conversations on this point, we told the authors that they could say whatever they wanted in their paper, but that the final word would belong to the editors.

129 **5.3 Jerde, Kraskura, Eliason, Csik, Stier, and Taper (JKECS&T)**

130 KECS&T describe the motivation for, and the logic of, scientific inference using evidential statistics and demonstrate the utility of the
 131 evidential approach by tackling a long-standing controversial question in ecological physiology: How does standard metabolic rate (SMR)
 132 scale (intra-specifically) with individual body mass, and is this scaling similar among species? For fish, theoretical scaling rates of 0.67,
 133 0.75, and 1.00 have been proposed. Empirical estimates of scaling coefficients vary tremendously among studies and generally all have
 134 large uncertainties leaving the theoretical question unprobed. JKECS&T curate a large data set composed of a total of 1456 observations in
 135 55 separate trials on 12 species, all using current state of the art techniques for measuring SMR. The use of linear mixed effect models
 136 allowed JKECS&T to combine all of these trials for inference.

137 Four suites of four models using random and fixed effects carefully explore the impacts of species, trial (within species), and temperature on
 138 the scaling of SMR with body mass. Model families were evaluated using the Schwarz information criterion (SIC, also known as the
 139 BIC). The SIC is a consistent criterion and the comparison of SIC values is an evidential procedure. Within and between model suites,
 140 evidence for specific values of the scaling coefficient were compared using profile Δ SIC curves. A Δ SIC value comparing two models
 141 indicates strong evidence for the model with lower SIC.

142 Two model suites with a free parameter estimate of the metabolic scaling, separated themselves only by a Δ SIC of 1.5, were strongly
 143 differentiated from all others. Both had fixed effects for temperature and random effects (intercepts) for species. The best model had the
 144 log(weight) slope vary randomly across species (with modest variation), while the second-best model had a common slope over all species.
 145 In the best model the ML estimate for the mean scaling coefficient is 0.89 with a strong evidence profile Δ SIC interval spanning 0.82-0.99.

146 The evidence strongly indicates that none of the *a priori* theoretical scaling coefficients describe the scaling behavior in real fish.

147 **5.4 Dennis, Ponciano, Taper, and Lele (DPT&L)**

148 Mathematics, and in particular probability, have long been intertwined with biology. The theoretician J. E. Cohen adroitly summarized the
 149 transcendence of the synergy between these fields with his essay “Mathematics is biology’s next microscope, only better; biology is
 150 mathematics’ next physics, only better” (Cohen, 2004). Key to the success of this interaction between these fields is the recognition that
 151 fundamental hypotheses in biology can be translated using the languages of mathematics, probability, and statistics into propositions that
 152 can be clearly probed. The increase in possibilities with such synergism is so dramatic that in some cases, it’s as if a new portal to a field of
 153 scientific inquiry becomes available. Yet, becoming enamored with model construction and the phrasing of novel explanations of biological
 154 phenomena can sometimes obscure the analyst’s vision and the realization that by its very human nature, mathematical models are limited
 155 constructs of biological processes. Mathematical models are indeed misspecifications of natural processes. Understanding the effects of
 156 model misspecification in our scientific inquiry should be paramount. This is the focus of DPT&L. These authors assess analytically and
 157 numerically the performance of Neyman-Person Hypothesis testing (NPHT), Fisher significance testing (NHST), information criteria and
 158 evidential statistics under model misspecification.

159 As mentioned above, evidential statistics seeks to quantify the strength of the evidence in the data for a reference model relative to another
 160 model. This goal is achieved through an evidence function, which is simply a statistic for comparing two models. Our evidence function of
 161 choice was Schwarz Information Criterion, or SIC (Schwarz 1978). The salient property of this and all evidence functions is that their
 162 associated probabilities of making a wrong model choice approach 0 as sample size increases. These probabilities, analogous to Type I and
 163 II errors in the Neyman-Pearson Hypothesis Testing (NPHT) framework are in fact pre-data error rates. Royall (2000) showed that these
 164 probabilities measure the chances of obtaining “weak misleading evidence” as well as strong misleading evidence. DPT&L shows that in a
 165 context where both models are in fact mathematical misspecifications of reality, making the wrong model choice refers to deeming as best a
 166 model that is not the closest to the true generating process model. By the same token, “misleading evidence” simply corresponds to
 167 obtaining observations that either weakly or strongly support a model other than the one that is the closest to the data-generating process.

168 Unlike the classic NPHT and Bayesian approaches, the Evidential Statistics paradigm provides sound guidelines to evaluate inferential
 169 errors when none of the proposed statistical models are a perfect representation of the natural, data-generating process. The NPHT
 170 framework depends critically on either the Null or the Alternative hypotheses being a perfect representation of the data generating
 171 mechanism and then fixes the Type I error probability irrespectively of sample size and thus problematically assesses the evidence *against*
 172 the null hypothesis and remains silent with respect to the evidence *for* the null hypothesis. The asymmetry of the NPHT error structure leads
 173 to difficulties in interpretation of hypotheses tests. The decision to pick an alternative model over a null hypothesis in and of itself is not
 174 controversial as it has some intuitively desirable statistical properties: for example, the probability to reject the null hypothesis given that the
 175 alternative is true converges to 0 as sample size increases. However, the probability of erroneously choosing the alternative when the null is
 176 true remains stuck at the chosen level alpha regardless of how large a sample size is collected. Matters get more complicated when it is
 177 considered that the original Neyman-Pearson theorem assumes that the data was generated under one of the two models but provides no
 178 guidance whatsoever in the event of model misspecification, a scenario commonly encountered in science. The fact that in scientific practice
 179 model comparison rarely stops at two models further muddying the interpretation of experimental results using the NPHT. To be fair,
 180 overconfidence in model selection procedures also results when the model misspecification is ignored in Bayesian Statistics (Yang and Zhu
 181 2018).

182 The evidential approach proposes fixing cutoff values for the evidence statistic, not the error probabilities. Under this concept of evidence,
 183 the value of a statistic like the likelihood ratio is evidence, not an error rate that is pre-set. Then, the evidential error probabilities both
 184 converge to 0 as sample size grows large. Finally, under this evidential statistics approach, the conclusion structure of say, a comparison
 185 between two models H_1 and H_2 has a trichotomy of outcomes: *i*) strong evidence for H_1 , *ii*) weak or inconclusive evidence and *iii*) strong
 186 evidence for H_2 .

187 Some, not all, information criteria commonly used for model selection are evidence functions. While the AIC only penalizes the likelihood
 188 function using the number of parameters, the SIC is also scaled by the sample size. As a result, as sample size increases, the error in deeming
 189 a model as “best” using the SIC statistics becomes vanishingly small. DPT&L show that this desirable property, called “Information

190 consistency" is lacking in the AIC. Inconsistent criteria, such as the AIC, tend to overfit at all sample sizes. Hence, the AIC is not an
 191 evidence function because it is not information consistent.

192 Although all paradigms of statistical science (NPHT, Bayesian statistics, Evidential Statistics) have flaws (reviewed in L_a and L_b), the
 193 Evidential Statistics paradigm possesses more desirable characteristics for the quantification of uncertainty and ultimately, for the design of
 194 inferential statements about the models' proximity to the true, generating process.

195 **5.5 Brittán and Bandyopadhyay (B&B)**

196 Written by a pair of philosophers of science, B&B provides a good entry into the research topic. Despite maintaining a high level of
 197 intellectual rigor, B&B avoids getting bogged down in technical statistical detail. The authors review the logical structures for scientific
 198 evidence: Hypothetico-deductive testing, Popperian falsification and corroboration, Fisherian significance, Neyman-Pearson hypothesis
 199 testing, the severe testing of Mayo, Bayesian confirmation, and statistical evidence.

200 The authors are equal opportunity balloon poppers pointing out the limitation of all methodological approaches. B&B focus on the strengths,
 201 weaknesses, and complementarity of statistical evidence and Bayesian confirmation. Contra the prevailing scientific mythos B&B
 202 demonstrate that Bayesian inference is "irreducibly personal". Bayesian methods do a good job of quantifying personal beliefs, and thus of
 203 informing personal decisions. Echoing L_a, B&B contend that non-informative priors are not objective and suffer from a variety of other
 204 problems. In contrast, statistical evidence does objectively quantify the relative support in data for specified pairs of models even though the
 205 models put forth for comparison may be generated subjectively.

206 Science is plagued by a suite of cognitive biases. Being aware of them can mitigate their impact. The authors note that each methodology
 207 works best to answer fairly narrow but different questions. Greater methodological self-consciousness on the part of scientists to match their
 208 choice of statistical approaches to match their scientific questions would promote scientific progress.

209 B&B close on the same hopeful note and metaphor as do S&H. Despite the undeniable subjectivity of individual scientists, Science itself
 210 may achieve a "Darwinian Objectivity" when the mutational force of subjective scientific creativity is filtered by objective evidential model
 211 selection.

212 **5.6 Ponciano and Taper (P&T)**

213 Information criteria have had a profound impact on modern science because they allow researchers to overcome the inadequacies of NPHT
 214 and tackle the multi-model selection process. Although model selection via information criteria gives the analyst an estimate of which
 215 probabilistic approximating models are closest to the generating process, information criterion comparison does not solve the problem of
 216 knowing how good the best model is. Indeed, the absolute distance to the generating process is not estimated through this process.

217 This caveat is all the more important when it is considered that in science, models are commonly misspecified. In this work, the authors
 218 resolve this shortcoming by designing a methodology to estimate a geometric representation of all the models under consideration along with
 219 the generating process. Such representation is a projection of all the models at hand into a two or three-dimensional space. As well, the
 220 location of the generating process in this representation is fully estimated. To estimate this model projection, the authors examined five key
 221 insights from Hirotugu Akaike's original work. These insights reveal the deep yet easy to grasp geometrical nature of Akaike's formulation
 222 of the AIC. P&T extend Akaike's geometrical interpretation and propose visualizing all models at hand into a reduced space. This reduced
 223 space representation applies ordination techniques to the models themselves so that the analyst may see and estimate the divergence between
 224 each model and every other model including the generating process itself.

225 P&T's solution starts from the observation that while standard information criterion analysis considers only the divergences of each model
 226 from the generating process, the divergences amongst all approximating models, typically ignored, are indeed estimable. As a test bed for
 227 their ideas, the authors consider two ecological scenarios, one of them involving an individual-based model simulation framework that
 228 generates data to which different abundance models can be fitted and the second one involving structural equation models.

229 The authors also compare their approach to model averaging and show that model projection is not as sensitive as model averaging to the
 230 composition of the set of candidate models being investigated. Model averaging artificially favors redundancy of model specification
 231 because the more models are developed in any given region of model space, the more heavily this particular region gets weighted.
 232 Furthermore, examining the resulting model space configuration can lead to an in-depth analysis of what are the model attributes that change
 233 from one model to the next that make it so that a model will get closer and closer to the generating process. This examination is the first step
 234 to explore models outside the bounds of the available model set, whereas by using model averaging, by definition, the analyst cannot do so.

235 Uncertainties around the estimation of model space estimation are yet not fully worked, but TLPD&J offers a first, non-parametric bootstrap
 236 approach to begin examining such question. Model projection methodology should be the starting point to do a science-based examination of
 237 critical model attributes that allow a model to get closer to the generating process (see also T&G). Finally, although P&T use the KL
 238 divergence as the fundamental distance measure, the model projections methodology could be extended or adapted to any other metric.

239 **5.7 Ferguson, Taper, Zenil-Ferguson, Jasieniuk and Maxwell (FTZJ&M)**

240 There are a vast number of information criteria. Academic arguments about which is best are intense and often vitriolic. FTZJ&M indicates
 241 that these arguments may be a tempest in teapot.

242 Seeking to improve model identification techniques for complex models with inter-dependent parameters, the authors modify Bozdogan's
 243 Information Complexity Criteria, ICC, to make them consistent and invariant to more kinds of transformations. To validate their suggested
 244 new criteria, FTZJ&M perform a vast array of performance comparisons. Twenty-five information criteria are investigated: Two classical
 245 efficient criteria (AIC and AICc), two classical consistent criteria (BIC and BIC*), three forms of Bozdogan's ICC, and 18 new
 246 modifications of the ICC. All of these criteria were compared for their ability in attaining three different model selection goals: Selecting

247 models with minimum prediction error, identifying the form of the generating model, and estimating the Kullback-Leibler distance to the
 248 generating process. All of this is done under 3 different classes of generating and approximating models, 3 different sample sizes, 3 different
 249 levels of process error, and 3 different levels of collinearity.

250 FTZJ&M recommend one of their combined forms ($BIC+2CvE(\Psi)$) as achieving all measures of quality well under a broad range of
 251 modeling frameworks and having the theoretical advantage of being both scale invariant and consistent. However, it is important to note that
 252 No IC was best for *any* goal over all conditions and that All IC performed generally well for all goals.

253 Two important lessons should be taken from FTZJ&M: First, much more attention needs to be paid to the uncertainty of model
 254 identification. And second, for these goals to be achieved sample sizes need to be larger in all model classes than is generally the case in
 255 ecology.

256 **5.8 Claeskens, Cunen and Hjort (CC&H)**

257 Perhaps the most used statistical tools by ecologists are abundance count models. Simply counting the number of individuals of every
 258 species observed in a particular community is the point of entry to deeper studies aiming at understanding the generation and maintenance of
 259 organisms' diversity. Profound questions examining the processes driving ecological stability, resilience, resistance, invasion, and
 260 persistence all begin with being able to accurately ascertain organisms' abundances. In our (joint) decades of teaching and mentoring, time
 261 and again count models keep coming back as some of the main instruments of statistical inference sustaining masters' theses and PhD
 262 dissertations in biology, wildlife ecology and conservation. Ecologists are typically not only interested in estimating one or the other model
 263 parameters leading to particular predictions, but often see parameter estimation as the by-product of what they are typically after, which is
 264 understanding which hypothesized model components better represent the underlying natural processes generating the count data at hand.

265 CC&H propose and further elaborate on a methodology that may revolutionize the reaches of an ecology-driven statistical analyses and in
 266 particular, multi-model selection for models of count data. The main idea of the Focused Information Criterion (FIC) approach is to provide
 267 a model selection framework where the comparison and the ranking is formally defined according to the scientific quest at hand.

268 Recognizing that different scientific teams might ask different focused questions of the same data and list of candidate models, CC&H
 269 design a methodology to focus the model selection process using different functions of the parameters of interest. When mainstream model
 270 selection tools are used in ecology and in a given scenario a model is chosen as the "best" model, practitioners are often left wondering why,
 271 in a specific scientific sense, such model is indeed the best model. FIC offers a theoretically sound methodology to obtain better, more
 272 precise estimates of a quantity of interest. For count models, such quantity is often the probability of a rare event occurring. As arbitrary or
 273 stale as it may sound at first, understanding and estimating accurately rare events in ecology has always been at the center of key
 274 explanations of diversity. Rarity, or "rare counts", have been for a long time (e.g. Patil and Taillie, 1982) hypothesized to be a critical
 275 component of explanations of how hyper-diverse communities can be maintained. Such was also the conclusion of one of the most recent
 276 and cited explanations of the maintenance of diversity in tropical forests published by Levi et al in (2019). As it turns out, the Focused
 277 Information Criterion of CC&H, which seeks to minimize the bias and the variance of a quantity of interest, works particularly well for

278 estimating the probability of rare events. In line with the “rarity” comments above, CC&H show as examples a situation where the focus of
 279 the inference is estimation of the probability of observing counts of a species above an arbitrary number. Importantly, the authors show
 280 how other information criteria like the BIC, although they may address the problem of determining which model is the closest to the true
 281 data generating mechanism, may not point towards the models that do the best job at estimating for instance, the tail of a distribution of
 282 counts. By allowing for a flexible specification of different “foci” of interest, CC&H provide a welcome addition to the toolbox of the
 283 evidentialist. This tool is not only conceptual but is crystallized in a practical, easy to use library for R users, the “fic” library.

284 **5.9 Markatou and Sofikitou (M&S)**

285 Most of the papers summarized so far share a key point: a reliance on the Kullback-Leibler divergence as the main instrument to develop
 286 and exemplify the theory and practice of Evidential Statistics. A natural reaction of any statistician to such heavy reliance on a single metric
 287 should be to ponder what would happen if different metrics or distances are used. Can the desiderata of evidential statistics be kept under
 288 different measures of divergence between the generating process and any approximating model, or amongst models themselves? Would the
 289 theoretical and asymptotic warrants of evidential statistics hold under different distance measures? How can statisticians visualize the
 290 strength of evidence under different measures? How does a measure of “strong evidence” using the KL divergence translates to other scales
 291 of divergence? These and other questions are approached using philosophical and rigorous statistical techniques in the contribution by
 292 M&S. Importantly, M&S’s contribution builds upon the pioneering concepts of model adequacy by Lindsay (2004) and evidence functions
 293 by Lele (2004). Notably, the authors propose an explanatory analysis tool called a “standardized distance ratio plot” that can be used to
 294 visualize the strength of evidence provided for or against hypotheses of interest using different divergence measures. Hence, this paper
 295 represents itself growth in the field and marks a clear path for future research. Indeed, of all the contributions in this special issue, this one is
 296 perhaps the one topic that is most ripe for further research and study. An open direction that seems promising is shining light on the behavior
 297 of different statistical divergence measures under model misspecification. Whenever we give seminars in statistics departments about
 298 evidential statistics, the question of usage of other divergence measures invariably comes up. We therefore encourage both, a close reading
 299 of this paper and thinking about building extensions to these results using M&S’s work as the foundation.

300 **5.10 Stuart and Blume (S&B)**

301 New statistical approaches often face resistance from empirical scientists. It can help acceptance if a new technique seems familiar. Stuart
 302 and Blume (S&B) cleverly disguise an evidential procedure with the face of a p-value, something that virtually every working scientist is
 303 familiar with. It does look like a p-value in that the statistic can take on values of 0, 1, and everything in between. S&B even strengthen the
 304 familiarity by calling it a SGPV or second-generation p-value.

305 Of course, a SGPV is not a p-value, it is not even a probability. The SGPV is better than a p-value. The question of interest is whether an
 306 unknown, but estimated, parameter is in an interval null or is outside of the interval null. A p-value or a null hypothesis significance test
 307 (NHST) can indicate that the parameter is likely outside the null, but neither can give you support that it is inside the null. Conversely, an
 308 equivalence test can give you support for the parameter being inside the interval but not for being outside the interval.

309 Evidence like, the procedure divides the range of possible value for the SG_{PV} into 3 regions: The point SG_{PV}=0, which indicates strong
 310 evidence the parameter is in the interval null. The point SG_{PV}=1, which indicates strong evidence the parameter is not in the null. And, the
 311 region of all values in between, which indicate that the data are consistent with both hypotheses and which way the evidence is tipping.

312 S&B also demonstrate another important evidential property. The SG_{PV} is consistent; the probability of misleading evidence goes to 0 as
 313 sample size increases.

314 The SG_{PV} is very flexible and can be applied retroactively to any scientific literature in which a statistical interval is published. S&B claim
 315 that SG_{PV} is applicable to any type of interval confidence, support, or credible. The authors spend the bulk of the paper demonstrating good
 316 statistical properties for the SG_{PV} under a wide range of circumstances.

317 **5.11 Lele a (L_a)**

318 It is undeniably true that State-Space Models (SMMs) or more generally, hierarchical statistical models, nowadays occupy a central role in
 319 ecology and evolution. SMMs are used to study the population dynamics of animals with complex life histories, to estimate abundances
 320 under detection limitations and heterogeneity (among individuals, across space and in time).. Entire statistical ecology books for graduate
 321 students and researchers alike with titles around “hierarchical models in ecology” now fill the electronic and physical bookshelves of modern
 322 ecologists and academicians. As well, social media with short instructionals, blogposts and even tweets by the authors of these books are
 323 consumed voraciously by graduate students needing to solve complex problems in the face of non-standard datasets. Software authors in
 324 turn, face the challenge of putting out for consumption accessible programs that can weather usage by anybody interested in applying a
 325 given hierarchical model. Over recent years, this high demand for accessible solutions to complex problems has facilitated the establishment
 326 of uncritical use of modern statistical machinery.

327 L_a approaches the consequences of such uncritical use head-on by clearly illustrating with real-life examples the predicaments brought
 328 about by using non-informative Bayesian analysis. Indeed, non-informative Bayesian analysis tends to be nowadays the default setting
 329 under which complex statistical models in ecology are fitted. In the name of pragmatism, it is often argued that in modern, extensive big
 330 data sets the sample size is so large that the likelihood information “swamps” any prior effect and that effectively, the data will “speak for
 331 itself”.

332 L_a carefully delineates the flaws in such reasoning and vividly details how and why wildlife management decisions can vastly suffer from
 333 such uncritical use of Bayesian techniques. In particular, he shows that because of the lack of parameterization invariance of non-
 334 informative Bayesian Analysis, all subjective Bayesian inferences can be disguised as “objective”, non-informative Bayesian inferences.
 335 Furthermore, cryptic biases can be introduced in the resulting analyses because the induced priors on functions of parameters are not non-
 336 informative.

337 Three other serious flaws are then discussed besides these two. However, even if the author had presented only these two problems,
338 practitioners, ecologists and wildlife managers should take note, because if the results of an uncritical non-informative Bayesian analysis is
339 subject to unstated and unqualified biases, it may be easily challenged in the legislature and in the Court of Law. For completeness,
340 professor Lele emphasizes that hierarchical models can be and are analyzed using the likelihood and frequentist methods. That is, **any**
341 Bayesian analysis can be transformed to a likelihood analysis by data cloning.

342 **5.12 Lele b (L_b)**

343 Uncertainty is a fundamental part of any inference, but the depth of its complexity is often not adequately appreciated. This paper, L_b, gives
344 a surprisingly readable review of many of the issues involved with statistical uncertainty. L_b begins with a short list, culled from the
345 literature, of desirable features for uncertainty quantification procedures: 1) transformation invariance, 2) uncertainty measure reflect data
346 informativeness, 3) ascertainability, and 4) diagnostic potential.

347 The first, transformation invariance, implies that the probability of an event occurring or not occurring is a reasonable measure of
348 uncertainty. This of course requires understanding what probability is and the paper next discusses the two major definitions of probability
349 used by statisticians and scientists alike: aleatory or frequency-based probability and epistemic or belief-based probability.

350 For adherents of frequentist statistics, data (*i.e.*, data sets) are random realizations from a stochastic generating process. Consequently,
351 estimates of parameters inherit stochasticity from the generating process through the stochasticity of data sets. The distribution of parameter
352 estimates over an infinite number of random data sets is called the true sampling distribution of the parameter. One can estimate a
353 parameters sampling distribution by bootstrap or analytic approximation. The estimated sampling distribution contains a great deal of
354 information about the uncertainty of the procedure. Much of this uncertainty is captured by confidence intervals. While arguing for the
355 utility of confidence intervals, L_b points out they are often misinterpreted.

356 L_b points out that the target of a confidence interval is to cover the true parameter, not to cover the parameter estimated in another
357 experiment. Another common way that confidence intervals are misinterpreted is by failing to distinguish between unconditional/pre-data
358 and conditional/post-data intervals. Both kinds of intervals are commonly used in the scientific literature. In separate sections L_b returns to
359 the questions of interval construction and interpretation from Bayesian and evidentialist perspectives.

360 As pointed out by B&B “any adequate (‘reliable’) hypothesis must be both explanatory and predictive.” It is only through the verification of
361 predictions that the ascertainment of models or hypotheses is possible. L_b takes this very seriously reviewing the representation of prediction
362 uncertainty in all three inferential paradigms. Further, a new flexible approach to the calculation of an evidential predictive density is
363 suggested and its advantages, both demonstrated and potential, are discussed.

364 The paper concludes by rehearsing the key features, strengths, and weaknesses of the characterization of uncertainty in the three paradigms
365 in the light of the four desiderata. None is perfect, but overall, the evidentialist most closely conforms. All three paradigms require scientists

366 to specify their models and whether inference should conditional or unconditional. Bayesian inference further requires the specification of
 367 priors, while evidence requires the specification of an evidence function. The last thing any reader wants to hear is that the quality of their
 368 scientific inference depends critically on the active choices they make—regardless of their statistical paradigm. Nevertheless, this is
 369 precisely the last thing that L_b says.

370 **5.13 Toquenaga and Gagné (T&G)**

371 Genetic sequencing is becoming an increasingly important tool in ecological and evolutionary studies. This trend has been accelerated by
 372 the new techniques of “next-generation sequencing”, NGS. These sequencing procedures work by digesting a genetic sequence into many
 373 small fragments (called reads), sequencing the fragments, and then inferring the original sequence computationally. This is like the spy
 374 novel trope of pasting a shredded letter back together.

375 With the scientific opportunities, come many statistical challenges. There are many programs that make these calculations. Unfortunately,
 376 they don’t agree—with each other and because many of the programs involve stochastic searches, even between multiple runs of the same
 377 program. T&G, use evidential principles to develop methods to choose among the many putative sequences offered by an array of
 378 sequencing software, to assess how good the proposed sequences are, and even to improve them.

379 The thinking in T&G is as follows: If multiple algorithms produce multiple sequences, each must be a model of the true sequence. If an
 380 appropriate function for measuring the divergence between these sequence models can be found, then the model projections in model space
 381 methods of P&T can be used to understand the relationships among the proposed models and even to a true sequence. The Levenshtein edit
 382 distance, as a measure of the minimum number of changes needed to equate two sequences from finite alphabets, offers itself as an
 383 appropriate divergence.

384 T&G test this proposition by taking a known genetic sequence and randomly breaking it into a number of fragments (with potential overlap).
 385 The number and distribution of fragment sizes are set to mimic typical digestion results. In their test case, T&G are able to construct, using
 386 non-metric dimensional scaling, a two-dimensional map of the sequence estimates produced by the various sequencing programs compared
 387 by the authors. Their map correctly identifies the best-proposed sequence.

388 In this test case, one of the programs is able to correctly reconstruct the true sequence. However, such a felicitous occurrence may not be
 389 general. Usefully, T&G propose an approach that can suggest sequences likely to improve on the set of mistaken sequences. They do this by
 390 proposing new sequence models which are consensus sequences of existing models and seeing where they fit into the map.

391 T&G confirm their method with a parametric bootstrap based on a specified true sequence. Implicit in this is the potential to use similar
 392 bootstrapping to assess the uncertainty in sequence construction.

393 **5.14 Taper, Lele, Ponciano, Dennis, and Jerde (TLPD&J).**

394 TLPD&J, develops themes from two other papers in this research topic. DPT&L show that in the presence of model misspecification
 395 Royall's universal bound on the strength of misleading evidence does not hold. L_b reminds us that statical uncertainty comes in two forms:
 396 global/unconditional and local/conditional.

397 To Royall's regions of weak and strong evidence (Royall, 2000) the authors intersperse a third category, that of prognostic evidence. This is
 398 evidence not so weak as to be dismissed nor so strong as to be considered overwhelming. Thus, while evidence is itself continuous, useful
 399 descriptive categories for considering evidence are constructed.

400 TLPD&J show that even in the presence of model misspecification the uncertainty in model identification can be quantified in the form of
 401 nonparametric bootstrap confidence intervals on evidence. This decouples evidence and its uncertainty and allows scientists to consider both.
 402 The authors consider evidence (either prognostic or strong) for one model over another to be "secure" if the lower 5% confidence limit on
 403 the evidence is above the preset prognostic boundary, k_p .

404 To demonstrate the utility of this approach, TLPD&J make a detailed reanalysis of model selection in Grace and Keeley's (2006) classic
 405 structural equation modeling of post-fire diversity recovery in California shrublands. The use of evidence confidence intervals develops a
 406 much more nuanced understanding of which model components are likely to be robust and which are equivocal.

407 Technically, TLPD&J use an improved version of the EIC (see Kitagawa and Konishi, 2010). The improvements include: 1) bootstrapping
 408 of the Δ SIC rather than individual likelihoods to incorporate the effects of misspecification geometry. And 2) Identification of components
 409 of EIC that correspond to global and local inference.

410 The paper finishes with an extended discussion of the interpretation of global and local inference in science.

411 **6 Conflict of Interest**

412 The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a
 413 potential conflict of interest.

414 **7 Author Contributions**

415 All authors contributed to the conception and writing of this editorial.

416 **8 References**

417 Brittan, G., and P. S. Bandyopadhyay. 2019. Ecology, Evidence, and Objectivity: In Search of a Bias-Free Methodology. *Frontiers in*
 418 *Ecology and Evolution* 7.

419 Claeskens, G., C. Cunen, and N. L. Hjort. 2019. Model Selection via Focused Information Criteria for Complex Data in Ecology and
420 Evolution. *Frontiers in Ecology and Evolution* 7.

421 Cohen, J. E. 2004. Mathematics is biology's next microscope, only better; Biology is mathematics' next physics, only better. *PLoS Biology*
422 2:2017-2023.

423 Dennis, B., J. M. Ponciano, M. L. Taper, and S. R. Lele. 2019. Errors in Statistical Inference Under Model Misspecification: Evidence,
424 Hypothesis Testing, and AIC. *Frontiers in Ecology and Evolution* 7.

425 Ferguson, J. M., M. L. Taper, R. Zenil-Ferguson, M. Jasieniuk, and B. D. Maxwell. 2019. Incorporating parameter estimability into model
426 selection. *Frontiers in Ecology and the Environment*.

427 Forster, M. R., and A. Wolfe. 1999. Conceptual Innovation and the Relational Nature of Evidence: The Whewell-Mill Debate.
428 Электронный ресурс. Режим доступа: <http://philosophy.wisc.edu/forster/papers/Whewell1.pdf>.

429 Franco, A., N. Malhotra, and G. Simonovits. 2014. Publication bias in the social sciences: Unlocking the file drawer. *Science* 345:1502-
430 1505.

431 Grace, J. B., and J. E. Keeley. 2006. A structural equation model analysis of postfire plant diversity in California shrublands. *Ecological
432 Applications* 16:503-514.

433 Haig, B. D. 2009. Inference to the best explanation: A neglected approach to theory appraisal in psychology. *American Journal of
434 Psychology* 122:219-234.

435 Ioannidis, J. P. A. 2019. What Have We (Not) Learnt from Millions of Scientific Papers with P Values? *The American Statistician* 73:20-25.

436 Jerde, C. L., K. Kraskura, E. J. Eliason, S. R. Csik, A. C. Stier, and M. L. Taper. 2019. Strong Evidence for an Intraspecific Metabolic
437 Scaling Coefficient Near 0.89 in Fish. *Frontiers in Physiology* 10.

438 Kerr, N. L. 1998. HARKing: hypothesizing after the results are known. *Personality and social psychology review: an official journal of the
439 Society for Personality and Social Psychology, Inc* 2:196-217.

440 Kitagawa, G., and S. Konishi. 2010. Bias and variance reduction techniques for bootstrap information criteria. *Annals of the Institute of
441 Statistical Mathematics* 62:209-234.

442 Lele, S. R. 2004. Evidence Functions and the Optimality of the Law of Likelihood.in M. L. Taper and S. R. Lele, editors. *The Nature of
443 Scientific Evidence: Statistical, Philosophical and Empirical Considerations*. The University of Chicago Press, Chicago.

444 Lele, S. R. 2004. Evidence Functions and the Optimality of the Law of Likelihood.in M. L. Taper and S. R. Lele, editors. *The Nature of
445 Scientific Evidence: Statistical, Philosophical and Empirical Considerations*. The University of Chicago Press, Chicago.

446 Lele, S. R. 2020. Consequences of lack of parameterization invariance of non-informative Bayesian analysis for wildlife management:
447 Survival of San Joaquin kit fox and declines in amphibian populations. *Frontiers in Ecology and Evolution* 7:501.

448 Lele, S. R. 2020. How Should We Quantify Uncertainty in Statistical Inference? *Frontiers in Ecology and Evolution* 8.

449 Levenshtein, V.I. 1966. Binary codes capable of correcting deletions, insertions, and reversals. *Soviet Physics Doklady*. 10 (8):707–710.

450 Levi, T., M. Barfield, S. Barrantes, C. Sullivan, R. D. Holt, and J. Terborgh. 2019. Tropical forests can maintain hyperdiversity because of
451 enemies. *Proceedings of the National Academy of Sciences of the United States of America* 116:581-586.

452 Lindsay, B. G. 2004. Statistical distances as loss functions in assessing model adequacy. Pages 439-488 in M. L. Taper and S. R. Lele,
453 editors. *The Nature of Scientific Evidence: Statistical, Philosophical and Empirical Considerations*. The University of Chicago Press,
454 Chicago.

455 Markatou, M., and E. M. Sofikitou. 2019. Statistical Distances and the Construction of Evidence Functions for Model Adequacy. *Frontiers*
456 in Ecology and Evolution 7:447.

457 Pashler, H., and E. J. Wagenmakers. 2012. Editors' Introduction to the Special Section on Replicability in Psychological Science: A Crisis of
458 Confidence? *Perspectives on Psychological Science* 7:528-530.

459 Patil, G. P., and C. Taillie. 1982. Diversity as a concept and its measurement. *Journal of the American Statistical Association* 77:548-561.

460 Peirce, C. S. 1974. Harvard lectures on pragmatism 1903. Pages 188-189 in C. Hartshorne and P. Weiss, editors. *The Collected Papers of*
461 *Charles Sanders Peirce*, Volume 5. Harvard University Press.

462 Ponciano, J. M., and M. L. Taper. 2019. Model Projections in Model Space: A Geometric Interpretation of the AIC Allows Estimating the
463 Distance Between Truth and Approximating Models. *Frontiers in Ecology and Evolution* 7.

464 Royall, R. M. 2000. On the Probability of Observing Misleading Statistical Evidence. *Journal of the American Statistical Association*
465 95:760-780.

466 Scheiner, S. M., and R. D. Holt. 2019. Evidential Statistics in Model and Theory Development. *Frontiers in Ecology and Evolution* 7.

467 Schwarz, G. 1978. Estimating the dimension of a model. *Annals of Statistics* 6:461-464.

468 Shimodaira, H. 2002. An approximately unbiased test of phylogenetic tree selection. *Systematic Biology* 51:492-508.

469 Shimodaira, H., and M. Hasegawa. 1999. Multiple comparisons of log-likelihoods with applications to phylogenetic inference. *Molecular*
470 *Biology and Evolution* 16:1114-1116.

471 Shimodaira, H., and Y. Terada. 2019. Selective Inference for Testing Trees and Edges in Phylogenetics. *Frontiers in Ecology and Evolution*
472 7.

473 Stewart, T. G., and J. D. Blume. 2019. Second-Generation P-Values, Shrinkage, and Regularized Models. *Frontiers in Ecology and*
474 *Evolution* 7:Article 486.

475 Taper, M. L., and J. M. Ponciano. 2016. Evidential statistics as a statistical modern synthesis to support 21st century science. *Population*
476 *Ecology* 58:9-29.

477 Taper, M. L., and P. J. P. Gogan. 2002. The Northern Yellowstone elk: Density dependence and climatic conditions. *Journal of Wildlife*
478 *Management* 66:106-122.

479 Taper, M. L., S. R. Lele, J. M. Ponciano, B. Dennis, and C. L. Jerde. 2021. Assessing the Global and Local Uncertainty of Scientific
480 Evidence in the Presence of Model Misspecification. *Frontiers in Ecology and Evolution* 9.

481 Taper, M. L., S. R. Lele, J.-M. Ponciano, and B. Dennis. 2019. Assessing the uncertainty in statistical evidence with the possibility of model
482 misspecification using a non-parametric bootstrap. *arXiv e-prints:arXiv: 1911.06421*.

483 Toquenaga, Y., and T. Gagne. 2021. The Evidential Statistics of Genetic Assembly: Bootstrapping a Reference Sequence. *Frontiers in*
484 *Ecology and Evolution* 9.

485 Yang, Z. H., and T. Q. Zhu. 2018. Bayesian selection of misspecified models is overconfident and may cause spurious posterior probabilities
486 for phylogenetic trees. *Proceedings of the National Academy of Sciences of the United States of America* 115:1854-1859.

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509 **9 Funding**

510 MLT was not funded for this work. JMP was partially supported by two grants from the US National Institute of Health grant number to
511 University of Florida NIH R01 GM103604 to JP and NIH 1R01GM117617 to Profs. Jason K. Blackburn and J.M. Ponciano. JMP was also
512 partially supported by NSF Award Number 2052372 to University of Florida (PI JMP, Project title: "Collaborative Research: Scaling
513 Properties of Ecological Variation in Complex Dynamical Systems")YT was partly supported by JSPS KAKENHI Grant Numbers
514 17H04612 and 18K06410.

515 **10 Acknowledgments**

516 We thank Brian Dennis and Subhash R. Lele for critically reading a draft of this editorial. We are grateful to all of the authors who
517 contributed to this research topic and to the editors and reviewers who worked hard to clarify the presentations. Our understanding of
518 evidence has been sharpened by years of discussions with Prasanta S. Bandyopadhyay, Gordon Brittan Jr., Brian Dennis, Jake M. Ferguson,
519 Christopher L. Jerde, Subhash R. Lele, and Kenichiro SHIMATANI.

520

521

522

523 11 Tables and figures

	Paper publication order and authorship													
	1 S&T	2 S&H	3 JKECS&T	4 DPT&L	5 B&B	6 P&T	7 FTZJ&M	8 CC&H	9 M&S	10 S&B	11 L _a	12 L _b	13 T&G	14 TLPD&J
Thematic concern														
Building evidence functions							*	*	*	*			*	*
Quantifying the uncertainty of evidence	*											*		*
Logic of statistical scientific inference		*	*	*	*	*								*
Application	*		*		*	*		*					*	*
Model space geometry	*			*		*							*	*
Comparative statistical inference				*	*	*					*	*		
Multiple comparisons and combining data	*		*			*								*
Model set misspecification	*			*		*							*	*

524

525 Table 1: Thematic concerns present in each article. Article authorship by publication order is: 1) Shimodaira and Terada (S&T), 2) Scheiner
526 and Holt (S&H), 3) Jerde, Kraskura, Eliason, Csik, Stier, and Taper (JKECS&T), 4) Dennis, Ponciano, Taper, and Lele (DPT&L), 5) Brittan
527 and Bandyopadhyay (B&B), 6) Ponciano and Taper (P&T), 7) Ferguson, Taper, Zenil-Ferguson, Jasieniuk and Maxwell (FTZJ&M, 8)
528 Claeskens et al., 9) Markatou and Sofikitou (M&S), 10) Stewart and Blume (S&B), 11) Lele a (L_a), 12) Lele b (L_b), 13) Toquenaga and
529 Gagné (T&G), 14) Taper, Lele, Ponciano, Dennis and Jerde (TLPD&J).

530

531

532