

Comment on ‘On the Birnbaum Argument for the Strong Likelihood Principle’ by Deborah G. Mayo

Michael Evans

University of Toronto

Abstract. We discuss Birnbaum’s result, its relevance to statistical reasoning, Mayo’s objections and the result in Evans (2013) that the proof of this result doesn’t establish what is commonly believed.

MSC 2010 subject classifications: Primary 62A01; secondary 62C10.

Key words and phrases: sufficiency, conditionality, likelihood, statistical evidence.

1. INTRODUCTION

The result established in Birnbaum (1962) that, if one accepts the frequentist principles of sufficiency (S) and conditionality (C), then one must accept the likelihood principle (L), has been an issue in the foundations of statistics for 50 years. Many statisticians and philosophers of science accept Birnbaum’s theorem as a logical fact because the proof is simple and, if they follow a pure likelihood or Bayesian prescription for inference, it doesn’t violate the way they think statistical analyses should be conducted. Many frequentist statisticians reject the result basically because they don’t like the consequence that frequentist evaluations of statistical methodologies are irrelevant.

In the end, an acceptable theory of inference has to be based on sound logic with no appeals to *ex cathedra* principles. Any principles used as part of forming such a theory have to have strong justifications and produce results that are free of paradoxes and contradictions. For example, the principle of conditional probability, which says we replace $P(A)$, as the measure of belief that event A is true, by $P(A|C)$ after being told that event C has occurred, seems like a basic principle of inference that, with careful application, is sound.

Does the likelihood principle carry the same weight in a theory of inference as the principle of conditional probability? We don’t think so and we will later argue that a somewhat weakened version is really just a consequence of the principle of conditional probability. Given that such principles can have a significant influence on what we view as correct statistical reasoning, it is important to examine the justifications for the likelihood principle, and Birnbaum’s theorem is commonly cited as such, to see if these are correct.

Department of Statistical Sciences, University of Toronto, Toronto, ON M5S 3G3, Canada

Another principle cited in Mayo's paper is the *principle of frequentism*. So what is the justification for this principle? Generally, this seems to be based on the belief that it typically produces sensible statistical methods although, as we will subsequently discuss, the story seems incomplete and unclear. If the principle of frequentism is correct, we need to have a good argument for it and a much more complete development of the theory.

The relevance of frequentism to Mayo's paper lies in the author's position that Birnbaum's argument is basically a violation of the principle. An argument is provided for why the joint application of S and C used in Birnbaum's proof constitute such a violation. I accept Mayo's reasoning. In fact, I think it is somewhat similar to the argument put forward in Evans, Fraser and Monette (1986) that the applications of S and C in the proof are incorrect because S discards as irrelevant precisely the information used by C to form the conditional model. So the justifications for S and C contradict one another in the proof and this doesn't seem right. This contradiction is avoided if one adopts the principle put forward in Durbin (1970), that we should restrict to ancillaries that are functions of a minimal sufficient statistic, and then Birnbaum's proof fails.

As we will discuss in Section 3, however, the issue in Birnbaum's argument is not really with S and C together, but rather with C itself and with what is actually proved. A very broad hint that this is the case is provided in Evans, Fraser and Monette (1986) where, using the same style of argument as Birnbaum, it is 'proved' that accepting C alone is equivalent to accepting L . So Durbin's point doesn't save the day even if we accept it. Actually, I don't think the arguments in Mayo's paper, or in Evans, Fraser and Monette (1986), completely dispense with Birnbaum's theorem either. They just reinforce the unsettling feeling that something is wrong somewhere.

Section 3 contains an outline of Evans (2013) that, for me at least, definitively settles the issue of what is wrong with Birnbaum's result and does this mathematically. As will be apparent from Section 2, however, it is clear that I believe that a proper prior is a necessary part of formally correct statistical reasoning. So why would a Bayesian want to invalidate Birnbaum's result? This is because the result, as usually stated, is not logically correct. Any valid theory has to have sound, logical foundations and so we don't want any faulty reasoning being used to justify that theory. In fact, Birnbaum's result even misleads as we've heard it said that model checking and checking for prior-data conflict violate the likelihood principle and so should not be carried out. Both of these activities are a necessary part of a statistical analysis. For this is how we deal, at least in part, with the subjectivity inherent in a statistical analysis due to the choices made by a statistician. This point, at least with respect to model checking, is also made in Mayo's paper and I think it is an excellent one.

For proper Bayesians, a form of the likelihood principle is a consequence of the principle of conditional probability, a far more important principle. Applying the principle of conditional probability to the joint probability model for the model parameter and data after observing the data, we have that *probability statements about the model parameter* depend on the sampling model and data only through the likelihood (note the emphasis). Of course, the likelihood map is minimal sufficient so there is nothing surprising in this.

2. BIRNBAUM AND EVIDENCE

There is an aspect of Birnbaum's work in this area that is particularly noteworthy. This is his emphasis on trying to characterize statistical evidence concerning the true value of the model parameter as expressed by the function Ev . Consider the pairs (M, x) , where $M = \{f_\theta : \theta \in \Theta\}$ is a set of probability distributions indexed by parameter $\theta \in \Theta$ and x is observed data coming from a distribution in M . Then Birnbaum (1962) writes $Ev(M_1, x_1) = Ev(M_2, x_2)$ to mean that the evidence in (M_1, x_1) is the same as the evidence in (M_2, x_2) whenever certain conditions are satisfied. We require here that M_1 and M_2 have the same parameter space but this can be weakened to include models with parameter spaces that are bijectively equivalent.

The principles S, C and L are considered as possible partial characterizations of statistical evidence. For example, if (M_1, x_1) and (M_2, x_2) are related via S , then Birnbaum says that, for frequentist statisticians, $Ev(M_1, x_1) = Ev(M_2, x_2)$ and similarly for C . Birnbaum is careful to say that Ev does not characterize what statistical evidence is, it is a kind of 'equivalence relation' (see Section 3).

In essence Birnbaum brings us to the heart of the matter in statistical inference. What is statistical evidence or more appropriately, how do we measure it? It seems collectively we talk about it, but we rarely get down to details and really spell out how we are supposed to handle this concept. Perhaps the closest to doing this is the pure likelihood theory, as discussed for example in Royall (1997), but this is only a definition of relative evidence when comparing two values of the full model parameter. For marginal parameters, this approach uses the profile likelihood as the only general way to compare the evidence for different values and this is unsatisfactory from many points of view. For example, a profile likelihood function is not generally a likelihood function.

For a frequentist theory of statistical inference, as opposed to a theory of statistical decision, it seems essential that a general method for measuring statistical evidence be provided that can be applied in any particular problem. The p-value is often used as a frequentist measure of evidence against a hypothesis, but, for a variety of reasons, it does not seem to be appropriate. For example, we need a measure that can also provide evidence *for* something being true and not just evidence against, given that we have assumed that the true distribution *is* in M .

If we add a proper prior to the ingredients, then it seems we can come up with sensible measures of evidence. For evidence, as expressed by observed data in statistical problems, is what causes beliefs to change and so we can measure evidence by measuring change in belief. For example, if we are interested in the truth of the event A , and this has prior probability $P(A) > 0$, then after observing C , the principle of conditional probability leads to the posterior probability $P(A|C)$ as the appropriate expression of beliefs about A . Accordingly, we measure evidence by the change in belief from $P(A)$ to $P(A|C)$. A simple *principle of evidence* says that we have evidence for the truth of A when $P(A|C) > P(A)$, evidence against the truth of A when $P(A|C) < P(A)$ and no evidence one way or the other when $P(A|C) = P(A)$. This principle is common in discussions about evidence in the philosophy of science and it seems obviously correct.

Of course, we also want to know how much evidence we have and this has lead to a variety of different measures based on $(P(A), P(A|C))$. The Bayes factor $BF(A|C) = P(A^c)P(A|C)P(A^c)/P(A)P(A^c|C)$ is one such measure, as

$BF(A|C) > 1$ if and only if $P(A|C) > P(A)$ and bigger values mean more evidence in favor of A being true. A central question associated with this, and other measures of evidence, is how to calibrate its values as in when is $BF(A|C)$ big and when is it small. Actually, we prefer measuring evidence via the relative belief ratio $RB(A|C) = P(A|C)/P(A)$, as the associated mathematics and the calibration of its values are both simpler. The generalization to continuous contexts is effected by taking limits and then both measures agree. A full theory of inference, both estimation and hypothesis assessment, can be built based on this measure of evidence together with a very natural calibration. This is discussed in Baskurt and Evans (2013). Of course, many will not like this because it involves proper priors, and so is subjective and supposedly not scientific. Alternatively, some may complain that priors are somehow hard to come up with.

In reality all of statistics, excepting the data when it is properly collected, is subjective. We choose models and we choose priors. What is important is that any choice we make, as part of a statistical analysis, be checkable against the objective data to ensure the choice at least make sense. We check the model by asking whether or not the data is surprising for each distribution in the model, and there are many well-known procedures for doing this. Perhaps not so familiar is that we can also check a proper prior by asking whether or not the true value is in a region of relatively low prior probability. Procedures for doing this consistently are developed in Evans and Moshonov (2006) and Evans and Jang (2011a). In fact, there are even logical approaches to modifying priors when prior-data conflict is found, as discussed in Evans and Jang (2011b). Moreover, with a suitable definition of evidence, we can measure *a priori* whether or not a prior is inducing bias into a problem, see Baskurt and Evans (2013). So subjectivity is not really the issue. We do our best to assess and control its effects, and maybe that is part of the role of statistics in science, but in the end it is an unavoidable aspect of any statistical investigation.

It is undoubtedly true that it is possible to write down complicated models for which it is extremely difficult, if not impossible, to prescribe an elicitation procedure in an application that leads to a sensible choice of a prior. But what does this say about our *choice* of model? It seems that we do not understand the effects of parameters in the model on the measurements we are taking sufficiently well to develop such a procedure. That is certainly possible, and perhaps even common, but it doesn't speak well for the modeling process and it shouldn't be held up as a criticism of what should be the gold standard for inference. An analogous situation arises with data collection where we know the gold standard is random sampling from the population(s) to which our inferences are to apply and, when we are interested in relationships among variables, controlled allocation of the values of predictors to sampled units. The fact that this is rarely, if ever, achieved doesn't cause us to throw out the baby with the dirty bath water. Gold standards serve as guides that we strive to attain and analyses that don't just need to be suitably qualified.

Our main point in this section is that the problem of measuring statistical evidence is the central issue in developing a theory of statistical inference. It seems that Birnbaum realized this and was searching for a way to accomplish this goal when he came upon what appeared to be a remarkable result.

3. WHAT'S WRONG WITH BIRNBAUM'S RESULT?

Perhaps everybody who has read the proof of Birnbaum's theorem is surprised at its simplicity. In fact, this is one of the reasons it is so convincing as there does not appear to be a logical flaw in the proof. As Mayo has noted, however, there are reasons to be doubtful of, if not even reject, the result as being valid within the domain of any sensible theory of statistical inference. Still suspicions linger as the formulation seems so simple.

As we will now explain, the result proved is not really the result claimed. If we want to treat Birnbaum's theorem and its proof as a piece of mathematics, then we have to be precise about the ingredients going into it. It is the imprecision in Birnbaum's formulation that leads to a faulty impression of exactly what is proved. This is more carefully explained in Evans (2013) but we can give a broad outline here.

Suppose we have a set D . A relation R on D is any subset $R \subset D \times D$. Meaningful relations express something and $(d_1, d_2) \in R$ means that d_1 and d_2 share some relevant property. Let \mathcal{I} denote the set of all model-data pairs (M, x) . So, for example, we can consider S as a relation on \mathcal{I} by saying the pair $((M_1, x_1), (M_2, x_2)) \in S \subset \mathcal{I} \times \mathcal{I}$ whenever (M_1, x_1) and (M_2, x_2) have equivalent minimal sufficient statistics. Similarly, C and L are relations on \mathcal{I} .

An equivalence relation R on D is a relation that is reflexive: $(d, d) \in R$ for all $d \in D$, symmetric: $(d_1, d_2) \in R$ implies $(d_2, d_1) \in R$ and transitive: $(d_1, d_2), (d_2, d_3) \in R$ implies $(d_1, d_3) \in R$. It is reasonable to say that, whatever property is characterized by relation R , when R is an equivalence relation, then $(d_1, d_2) \in R$ means that d_1 and d_2 possess the property to the same degree. It is easy to prove that S and L are equivalence relations but C and $S \cup C$ are not equivalence relations, see Evans (2013).

Associated with an arbitrary relation R on D is the smallest equivalence relation on D containing R , which we will denote by \bar{R} . Clearly, \bar{R} is the intersection of all equivalence relations containing R . But \bar{R} can also be characterized in another way that is key to Birnbaum's proof.

Lemma If R is a reflexive relation on D , then $\bar{R} = \{(x, y) : \exists n, x_1, \dots, x_n \in D \text{ with } x = x_1, y = x_n \text{ and } (x_i, x_{i+1}) \in R \text{ or } (x_{i+1}, x_i) \in R\}$.

Note that S and C are both reflexive and thus $S \cup C$ is reflexive.

In Birnbaum's proof, he starts with $((M_1, x_1), (M_2, x_2)) \in L$, namely, these pairs have proportional likelihoods. Birnbaum constructs the mixture model (Birnbaumization) M^* and then argues that we have that $((M_1, x_1), (M^*, (1, x_1))) \in C$, $((M^*, (1, x_1)), (M^*, (2, x_2))) \in S$ and $((M^*, (2, x_2)), (M_2, x_2)) \in C$. Since $C \subset S \cup C$ and $S \subset S \cup C$, by the Lemma, this proves that $L \subset \overline{S \cup C}$ and this is all that Birnbaum's argument establishes. Since $S \cup C \subset L$ and L is an equivalence relation, we also have $L = \overline{S \cup C}$. As shown in Evans (2013), it is also true that $S \cup C$ is properly contained in L , so there is some content to the proof. In prose, Birnbaum's proof establishes the following: if we accept S , and we accept C , and we accept all the equivalences generated by these principles jointly, then we accept L . Certainly accepting S and C is not equivalent to accepting L since $S \cup C$ is a proper subset of L . We need the additional hypothesis and there doesn't appear to be any good reason why we should accept this as part of a theory of statistical inference. It is easy to construct relations R where \bar{R} is meaningless.

So we have to justify the additional pairs we add to a relation when completing it to be an equivalence relation.

It is interesting to note that the argument supposedly establishing the equivalence of C and L in Evans, Fraser and Monette (1986), also proceeds in the same way using the method of the Lemma. Since C is properly contained in L , this proof establishes that $\bar{C} = L$. So in fact, S is irrelevant in Birnbaum's proof. The problem with the principles S and C , as partial characterizations of statistical evidence, lies with C and the fact that it is not an equivalence relation. That C is not an equivalence relation is another way of expressing the well-known fact that, in general, a unique maximal ancillary doesn't exist.

The result $\bar{C} = L$ does have some content. To be a valid characterization of evidence in the context of \mathcal{I} , we will have to modify C so that it is an equivalence relation. The smallest equivalence relation containing C is L and this is unappealing, at least to frequentists, as it implies that repeated sampling properties are irrelevant for inference. Another natural candidate for a resolution is the largest equivalence relation contained in C that is compatible with all the equivalence relations based on maximal ancillaries. This is given by the equivalence relation based on the laminal ancillary. From Basu (1959), ancillary statistic a is a *laminal ancillary* if it is a function of every maximal ancillary and any other ancillary with this property is a function of a . The laminal ancillary is essentially unique. It is unclear how appealing this resolution would be to frequentists, but there don't seem to be any other natural candidates.

Many authors, including Mayo, refer to the weak conditionality principle which restricts attention to ancillaries that are physically part of the sampling. In such a case we would presumably write our models differently so as to reflect the fact that this sampling occurred in stages. In other words, the universe is different than \mathcal{I} , the one Birnbaum considered. There doesn't seem to be anything controversial about such a principle and it is well-motivated by the two measuring instruments example and many others.

We don't believe, however, that the weak conditionality principle resolves the problem with conditionality more generally. For example, how does weak conditionality deal with situations like Example 2-2 in Fraser (2004) and many others like it? Conditioning on an ancillary seems absolutely essential if we are to obtain sensible inferences in such examples, but there doesn't appear to be any physical aspect of the sampling that corresponds to the relevant ancillary.

Many frequentist statisticians ignore conditionality but, as noted in Fraser (2004), this is not logical. The theme in conditional inference is to find the right hypothetical sequence of repeated samples to compare the observed sample to. This takes us back to our question concerning the principle of frequentism: why are we considering repeated samples anyway? A successful frequentist theory of inference requires at least a resolution of the problems with conditionality. The lack of such a resolution leads to doubts as to the validity of the basic idea that underlies frequentism.

Issues concerning ancillaries are not irrelevant to Bayesians as they have uses in model checking and checking for prior-data conflict. Notice that the principle of conditional probability does not imply that these activities need refer to any kind of posterior probabilities and it is perfectly logical for these to be based on prior probabilities. For example, model checking can be based on the distribution of

an ancillary or the conditional distribution of the data given a minimal sufficient. Of course, C is not relevant for proper Bayesian probability statements about θ , as the principle of conditional probability implies that we condition on all of the data.

We acknowledge that it is possible that the problems with C might be fixable or even eliminated through a better understanding of what we are trying to accomplish in statistical analyses — these aren't just problems in mathematics. We can't resist noting, however, that the simple addition of a proper prior to the ingredients does the job, at least for inference.

4. CONCLUSIONS

Mayo's paper contains a number of insightful comments and more generally it helps to focus attention on what is the most important question in statistics, namely, what is the right way to formulate a statistical problem and carry out a statistical analysis. To a certain extent, Birnbaum's result has been an impediment in moving forward towards developing a theory of inference that has a solid foundation. It is good to have such underbrush removed from the discussion. We have to give great credit to Birnbaum, however, for his focus on what is important in achieving this goal, namely, the measurement of statistical evidence. That his theorem has lasted for so long is a testament to the difficulties involved in this task.

In general, we need a strong foundation for a theory of statistical inference rather than principles, often not clearly stated, that have only some vague, intuitive appeal. The only way we can determine whether or not an instance of statistical reasoning is correct, lies within the context of a sound theory. That two statistical analyses based on the same data and addressing the same question can be deemed to be correct and yet come to different conclusions, is not a contradiction. Statistics tells us that we simply must collect more data to resolve such differences. In our view, the role of statistics in science is to explain how to reason correctly in statistical contexts. Without a strong theory we can't do that.

REFERENCES

- BASKURT, Z. and EVANS, M. (2013). Hypothesis assessment and inequalities for Bayes factors and relative belief ratios. *Bayesian Analysis*, **8**, **3**, 569-590.
- BASU, D. (1959) The family of ancillary statistics. *Sankhya: The Indian Journal of Statistics*, **21**, **3/4**, 247-256. [MR0110115](#)
- EVANS, M. (2013) What does the proof of Birnbaum's theorem prove? *Electronic Journal of Statistics*, **7**, 2645-2655.
- EVANS, M. and JANG G. H. (2011a) A limit result for the prior predictive. *Statistics and Probability Letters*, **81**, 1034-1038. [MR2803740](#)
- EVANS, M. and JANG G. H. (2011b) Weak informativity and the information in one prior relative to another. *Statistical Science*, **26**, **3**, 423-439. [MR2917964](#)
- EVANS, M. and MOSHONOV, H. (2006) Checking for prior-data conflict. *Bayesian Analysis*, **1**, **4**, 893-914. [MR2282210](#)
- FRASER, D. A. S. (2004) Ancillaries and conditional inference (with discussion). *Statistical Science*, **19**, **2**, 330-369. [MR2140544](#)
- ROYALL, R. (1997). *Statistical Evidence: A likelihood paradigm*. Chapman & Hall. [MR1629481](#)