

BFF Four: Are we Converging?

Nancy Reid

May 2, 2017

[HOME](#) / [RESEARCH](#) /

Fourth Bayesian, Fiducial, and Frequentist
Conference (BFF4)

Harvard University

May 1 - 3, 2017

Hilles Event Hall, 59 Shepard St, Cambridge, MA



HARVARD
Faculty of Arts and Sciences
DEPARTMENT OF STATISTICS

Classical Approaches: A Look Way Back

Nature of Probability

BFF one to three: a look back

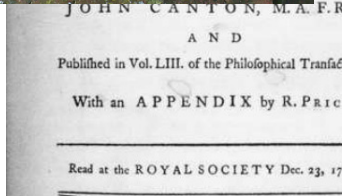
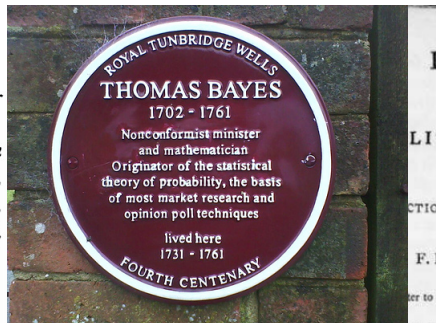
Comparisons

Are we getting there?

LII. *An Essay towards solving a Problem in the Doctrine of Chances. By the late Rev. Mr. Bayes, F. R. S. communicated by Mr. Price, in a Letter to John Canton, A. M. F. R. S.*

Dear Sir,

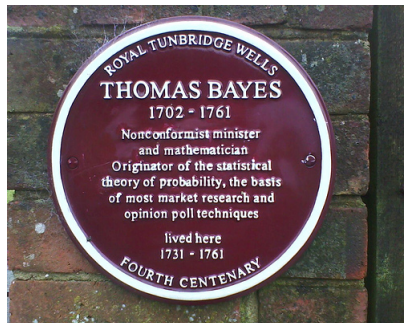
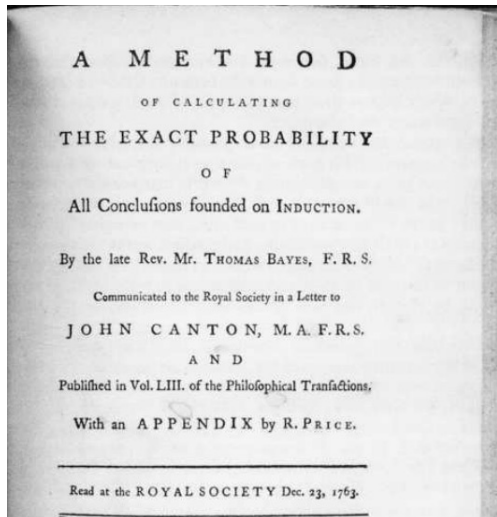
Read Dec. 23, 1763. **I** Now fend you an essay which I have found among the papers of our deceased friend Mr. Bayes, and which, in my opinion, has great merit, and well deserves to be preserved. Experimental philosophy, you will find, is nearly interested in the subject of it; and on this account there seems to be particular reason for thinking that a communication of it to the Royal Society cannot be improper.



Stigler 2013

Posterior Distribution

Bayes 1763



$$\pi(\theta | y^0) = f(y^0; \theta)\pi(\theta)/m(y^0)$$

probability distribution for θ
 y^0 is fixed

probability comes from $\pi(\theta)$

Fiducial Probability

Fisher 1930

528

Dr Fisher, Inverse probability

Inverse Probability. By R. A. FISHER, Sc.D., F.R.S., Gonville and Caius College; Statistical Dept., Rothamsted Experimental Station.

[Received 23 July, read 28 July 1930.]

$$df = -\frac{\partial}{\partial \theta} F(T, \theta) d\theta$$

fiducial probability density for θ , given statistic T

probability comes from (dist'n of) T

“It is not to be lightly supposed that men of the mental calibre of Laplace and Gauss ... could fall into error on a question of prime theoretical importance, without an uncommonly good reason”



SOME PROBLEMS CONNECTED WITH STATISTICAL INFERENCE

By D. R. Cox

*Birkbeck College, University of London*¹

1. Introduction. This paper is based on an invited address given to a joint meeting of the Institute of Mathematical Statistics and the Biometric Society at Princeton, N. J., 20th April, 1956. It consists of some general comments, few of them new, about statistical inference.

Since the address was given publications by Fisher [11], [12], [13], have produced a spirited discussion [7], [21], [24], [31] on the general nature of statistical methods. I have not attempted to revise the paper so as to comment point by point on the specific issues raised in this controversy, although I have, of course, checked that the literature of the controversy does not lead me to change the opinions expressed in the final form of the paper. Parts of the paper are controversial; these are not put forward in any dogmatic spirit.

2. Inferences and decisions. A statistical inference will be defined for the



- "Much controversy has centred on the distinction between fiducial and confidence estimation"
- "... The fiducial approach leads to a distribution for the unknown parameter"
- "... the method of confidence intervals, as usually formulated, gives only one interval at some preselected level of probability"
- "... in ... simple cases ... there seems no reason why we should not work with confidence distributions for the unknown parameter"
- "These can either be defined directly, or ... introduced in terms of the set of all confidence intervals"

Confidence Distribution

Cox 1958; Efron 1993

μ , and that inferences are desired for $\theta = t(\mu)$, a real-valued function of μ . Let $\theta_x(\alpha)$ be the upper endpoint of an exact or approximate one-sided level- α confidence interval for θ . The standard intervals for example have

$$\theta_x(\alpha) = \hat{\theta} + \hat{\sigma} z^{(\alpha)}, \quad (1.1)$$

where $\hat{\theta}$ is the maximum likelihood estimate of θ , $\hat{\sigma}$ is the Fisher information estimate of standard error for $\hat{\theta}$, and $z^{(\alpha)}$ is the α -quantile of a standard normal distribution, $z^{(\alpha)} = \Phi^{-1}(\alpha)$. We write the inverse function of $\theta_x(\alpha)$ as $\alpha_x(\theta)$, meaning the value of α

This content downloaded from 142.150.190.39 on Sun, 09 Apr 2017 20:35:40 UTC
All use subject to <http://about.jstor.org/terms>



4

BRADLEY EFRON

corresponding to upper endpoint θ for the confidence interval, and assume that $\alpha_x(\theta)$ is smoothly increasing in θ . For the standard intervals, $\alpha_x(\theta) = \Phi((\theta - \hat{\theta})/\hat{\sigma})$, where Φ is the standard normal cumulative distribution function.

The confidence distribution for θ is defined to be the distribution having density

$$\pi_x^{\dagger}(\theta) = d\alpha_x(\theta)/d\theta. \quad (1.2)$$

- “assigns probability 0.05 to θ lying between the upper endpoints of the 0.90 and 0.95 confidence intervals, etc.”
- “Of course this is logically incorrect, but it has powerful intuitive appeal”
- “... no nuisance parameters [this] is exactly Fisher’s fiducial distribution”

Seidenfeld 1992; Zabell 1992

Biometrika (1966), 53, 1 and 2, p. 1
Printed in Great Britain

1

Structural probability and a generalization*

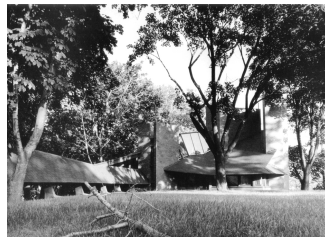
By D. A. S. FRASER
University of Toronto

SUMMARY

Structural probability, a reformulation of fiducial probability for transformation models, is discussed in terms of an error variable. A consistency condition is established concerning conditional distributions on the parameter space; this supplements the consistency under Bayesian manipulations found in Fraser (1961). An extension of structural probability for real-parameter models is developed; it provides an alternative to the local analysis in Fraser (1964*b*).

1. INTRODUCTION

Fiducial probability has been reformulated for location and transformation models (Fraser, 1961) and compared with the prescriptions in Fisher's papers (Fraser, 1963*b*). The transformation formulation leads to a frequency interpretation and to a variety of consistency conditions; the term *structural probability* will be used to distinguish it from Fisher's formulation.



- “a re-formulation of fiducial probability for transformation models”
- “This transformation re-formulation leads to a frequency interpretation”
- a change in the parameter value can be offset by a change in the sample
 $y \rightarrow y + a; \theta \rightarrow \theta - a$
- a local location version leads to:

$$df = -\frac{\partial}{\partial \theta} F(y, \theta) d\theta = -\frac{\partial}{\partial \theta} F(y, \theta) \frac{f(y^0, \theta)}{f(y^0, \theta)} = \overbrace{f(y^0, \theta)}^{\text{Likelihood}} \left. \frac{dy}{d\theta} \right|_{y^0}$$

Significance Function

Fraser 1991

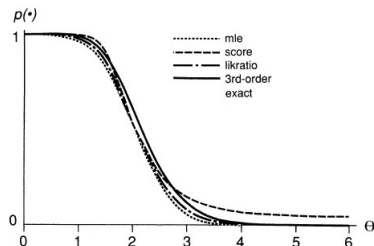
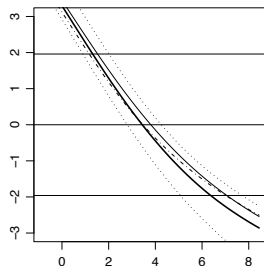


Figure 5. The Standardized Maximum Likelihood Estimate, Standardized Score, and Signed Likelihood Ratio Produce Three Approximations for the Significance Function. Model: location log gamma(3); data: $v^0 = 3.14$.



- “from likelihood to significance”
- “significance records probability left of the observed data point”
”likelihood records probability at the observed data point”
- the significance function is a plot of this probability
as a function of θ
- “the full spectrum of confidence intervals is obtained ...
suggesting the alternate name **confidence distribution function**”

- probability to describe physical haphazard variability **aleatory**
 - probabilities represent features of the “real” world in somewhat idealized form
 - subject to empirical test and improvement
 - conclusions of statistical analysis expressed in terms of interpretable parameters
 - enhanced understanding of the data generating process

- probability to describe the uncertainty of knowledge **epistemic**
 - measures rational, supposedly impersonal, degree of belief, given relevant information Jeffreys

 - measures a particular person’s degree of belief, subject typically to some constraints of self-consistency Ramsey, de Finetti, Savage

 - often linked with personal decision making necessarily?

... nature of probability

- Bayes posterior describes uncertainty of knowledge
- probability comes from the prior

- confidence intervals or p -values refer to empirical probabilities

- in what sense are confidence distribution functions, significance functions, structural or fiducial probabilities to be interpreted?

- empirically? degree of belief?
- literature is not very clear imho
- we may avoid the need for a different version of probability by appeal to a notion of calibration

BFF 1 – 4

- posterior distribution
- fiducial probability
- confidence distribution
- structural probability
- significance function
- belief functions
- objective Bayes
- generalized fiducial inference
- confidence distributions and confidence curves
- approximate significance functions
- inferential models

What has changed?

computation



- noninformative, default, matching, reference, ... priors
- we may avoid the need for a different version of probability by appeal to a notion of calibration

Cox 2006, R & Cox 2015

- as with other measuring devices within this scheme of repetition, probability is defined as a hypothetical frequency
- it is unacceptable if a procedure yielding high-probability regions in some non-frequency sense are poorly calibrated
- such procedures, used repeatedly, give misleading conclusions

Bayesian Analysis, V1(3) 2006

... objective Bayes

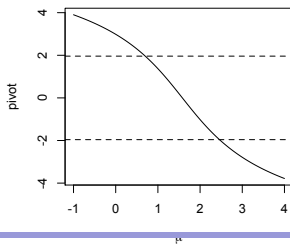
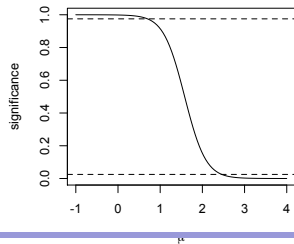
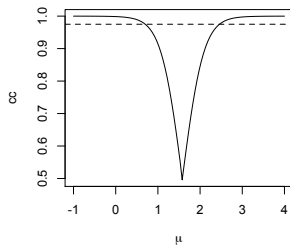
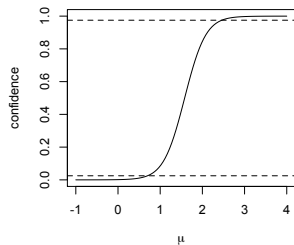
- pragmatic solution as a starting point
- some versions may not be correctly calibrated
- requires checking in each example
- calibrated versions must be targetted on the parameter of interest
- only in very special cases can calibration be achieved for more than one parameter in the model, from the same prior
- the simplicity of a fully Bayesian approach to inference is lost

Gelman 2008; PPM LW

- any function $H : \mathcal{Y} \times \Theta \rightarrow (0, 1)$ which is
- a cumulative distribution function of θ for any $y \in \mathcal{Y}$
- has correct coverage: $H(Y, \theta) \sim U(0, 1)$ $Y \sim f(\cdot; \theta)$
- CDs, or approximate CDs, are readily obtained from pivotal quantities
- pivotal quantity: $g(y, \theta)$ with sampling distribution known $\sqrt{n}(\bar{y} - \mu)/s$
- sufficiently general to encompass bootstrap distribution and many standard likelihood quantities
- recent examples include robust meta-analysis and identification of change-points

Xie et al. 2011; Cunen et al. 2017; Hannig & Xie 2012

... confidence distribution



- Fisher: $g(Y, \theta)$ has a known distribution; invert this to create distribution for θ when y^0 is obtained
- Fraser: use data-generating equation to make the inversion more direct, e.g. $Y_i = \mu + e_i, \quad e_i = y_i^0 - \mu$
- Hannig et al. $Y = G(U, \theta), \quad U$ has known distribution
- suppose we can invert this for any y^0 : $\theta = Q_{y^0}(U)$
- fiducial distribution of θ is $Q_{y^0}(U^*)$ U^* independent copy of U
- inverse only exists if θ and Y have same dimension
- might get this by reduction of a sample to sufficient statistics
- more generally, fiducial density takes the form

$$r(\theta; y^0) \propto f(y^0, \theta) J(y^0, \theta)$$

- focus on parameter of interest

many arguments point to the need for this

- $y \in \mathbb{R}^n$, $\theta \in \mathbb{R}^p$, $\psi \in \mathbb{R}$

- focus on dimension reduction

- $n \downarrow p$ using an approximation location model

$$df = -\frac{\partial}{\partial \theta} F(y, \theta) d\theta = -\frac{\partial}{\partial \theta} F(y, \theta) \frac{f(y^0, \theta)}{f(y^0, \theta)} = \overbrace{f(y^0, \theta)}^{\text{Likelihood}} \frac{dy}{d\theta} \Big|_{y^0}$$

- $p \downarrow 1$ using a tangent exponential model

- combination of structural model and likelihood asymptotics

- leads for example to construction of default priors

Fraser et al 2010

- Hannig et al. $Y = G(U, \theta)$, U has known distribution
- suppose we can invert this for any y^0 : $\theta = Q_{y^0}(U)$
- fiducial distribution of θ is $Q_{y^0}(U^*)$ U^* independent copy of U
- use a random set \mathcal{S} to predict U
- this random set is converted to a belief function about θ
- need to ensure the belief function is **valid** and **efficient**
- valid = calibrated ?

Comparisons: conditioning

- objective Bayes
- generalized fiducial inference
- confidence distributions and confidence curves
- approximate significance functions
- inferential models
- yes
- yes and no JASA '16
- needs to be built in ahead of time
- yes; via approximate location model
- needs to be built in ahead of time

Comparisons: Eliminating Nuisance Parameters

- objective Bayes
- generalized fiducial inference
- confidence distributions and confidence curves
- approximate significance functions
- inferential models
- marginalization
 - rarely works ...
- depends on the problem
 - ?
- use profile log-likelihood, or similar
 - focus parameter
- marginalization
 - via Laplace approximation
- marginalization
 - invoked ahead of time

Comparisons: Calibration

- objective Bayes
- generalized fiducial inference
- confidence distributions and confidence curves
- approximate significance functions
- inferential models
- often
- yes
- typically approximate
- typically approximate
- yes

Comparisons: Nature of Probability

- Bayes / objective Bayes
- generalized fiducial inference
- confidence distributions and confidence curves
- approximate significance functions
- inferential models
- epistemic / empirical
- empirical
- empirical
but not prescriptive
- empirical
- ?epistemic?

What's the end goal?

- Applications – something that works
 - gives 'sensible' answers
 - not too sensitive to model assumptions
 - computable in reasonable time
 - provides interpretable parameters

- Foundations – peeling back the layers
 - what does 'works' mean?
 - what probability do we mean
 - 'Goldilocks' conditioning
 - **how does this impact applied work?**

Meng & Liu, 2016



- avoid apparent discoveries based on spurious patterns
- to shed light on the structure of the problem
- calibrated inferences about interpretable parameters
- realistic assessment of precision
- understanding when/why methods work/fail

Well, are we?



objective Bayes

Larry W: “the perpetual motion machine of Bayesian inference”

confidence distributions

Min-ge, Regina: “everything fits”
Nils: “CDs are the ‘gold standard’ ”

generalized fiducial

Jan: “bring it on ... I’ll figure it out”

inferential models

Ryan: “it’s the only solution”
Chuanhai: “ it might take 100 years”

significance functions

Don: “it’s the best solution ...
you can’t solve everything at once”

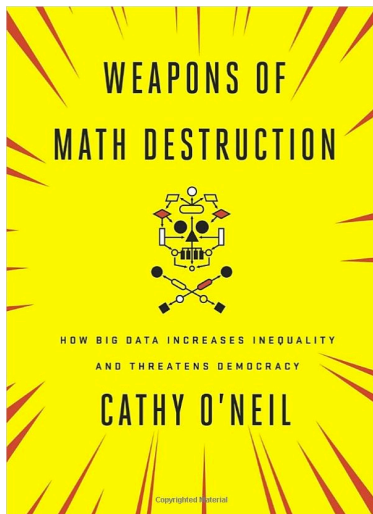
Some warning signs

Big data: are we making a big mistake?

Economist, journalist and broadcaster **Tim Harford** delivered the 2014 *Significance* lecture at the Royal Statistical Society International Conference. In this article, republished from the *Financial Times*, Harford warns us not to forget the statistical

“Big data” has arrived, but big insights have not

Some warning signs



Some warning signs

UNDERSTANDING DEEP LEARNING REQUIRES RE-THINKING GENERALIZATION

Chiyuan Zhang*

Massachusetts Institute of Technology
chiyuan@mit.edu

Samy Bengio

Google Brain
bengio@google.com

Moritz Hardt

Google Brain
mrtz@google.com

Benjamin Recht†

University of California, Berkeley
brecht@berkeley.edu

Oriol Vinyals

Google DeepMind
vinyals@google.com

ABSTRACT

Despite their massive size, successful deep artificial neural networks can exhibit a remarkably small difference between training and test performance. Conventional wisdom attributes small generalization error either to properties of the model family, or to the regularization techniques used during training.

Through extensive systematic experiments, we show how these traditional approaches fail to explain why large neural networks generalize well in practice. Specifically, our experiments establish that state-of-the-art convolutional networks for image classification trained with stochastic gradient methods easily fit a random labeling of the training data. This phenomenon is qualitatively unaffected by explicit regularization, and occurs even if we replace the true images by completely unstructured random noise. We corroborate these experimental findings with a theoretical construction showing that simple depth two neural networks already have perfect finite sample expressivity as soon as the number of parameters exceeds the number of data points as it usually does in practice.

Some warning signs

UNDERSTANDING DEEP LEARNING REQUIRES RE-THINKING GENERALIZATION

Chiyuan Zhang*
Massachusetts Institute of Technology
chiyuan@mit.edu

Samy Bengio
Google Brain
bengio@google.com

Moritz Hardt
Google Brain
mrtz@google.com

Benjamin Recht†
University of California, Berkeley
brecht@berkeley.edu

Oriol Vinyals
Google DeepMind
vinyals@google.com

“deep neural networks easily fit random labels”

“these observations rule out ... possible explanations for the generalization performance of state-of-the-art neural networks.”

Some warning signs

theguardian

Discrimination by algorithm: scientists devise test to detect AI bias

Researchers devise test to determine whether machine learning algorithms are introducing gender or racial biases into decision-making

theguardian.com

Facial recognition database used by FBI is out of control, House committee hears

Database contains photos of half of US adults without consent, and algorithm is wrong nearly 15% of time and is more likely to misidentify black people

Summary

- Bayes, fiducial, structural, confidence, belief
- BFF 1 - 4: Develop links to bridge gaps among different statistical paradigms
- targetting parameters
- limit distributions
- calibration in repeated sampling
- relevant repetitions for the data at hand

NR: Why is conditional inference so hard?

DRC: I expect we're all missing something, but I don't know what it is

StatSci Interview 1996

References

- Cox, D.R. (1958). *Ann. Math. Statist.*
- Cox, D.R. (2006). *Principles of Statistical Inference.*
- Cunen et al. (2017). *J. Statist. Plann. Infer.*
- Efron, B. (1993). *Biometrika*
- Fisher, R.A. (1930). *Proc. Cam. Phil. Soc.*
- Fraser, D.A.S. (1966). *Biometrika*
- Fraser, D.A.S. (1991). *J. Amer. Statist. Assoc.*
- Fraser, D.A.S. et al. (2010). *J. R. Statist. Soc. B*
- Fraser, D.A.S. and Reid, N. (1993). *Statist. Sinica*
- Gelman, A. (2008). *Ann. Appl. Statist.*
- Hannig, J. and Xie, M. (2012). *Elect. J. Statist.*
- Hannig, J. et al. (2016). *J. Amer. Statist. Assoc.*
- Hjort, N. and Schweder, T. *Confidence, Likelihood, Probability: Inference with Confidence Distributions*
- Meng, X.-L. and Liu, K. (2016). *Ann. Rev. Stat. and its Applic.*
- Norman, S. et al. (2016). *Ann. Rev. Stat. and its Applic.*
- Reid, N. and Cox, D.R. (2015). *Intern. Statist. Rev.*
- Stein, C. (1959). *Ann. Math. Statist.*
- Stigler, S. (2013). *Statistical Science*
- Xie, M. and Singh, K. (2013) *Internat. Statist. Rev.*
- Xie, M. et al. (2011). *J. Amer. Statist. Assoc.*
- Zabell, X. (1992). *Statistical Science*