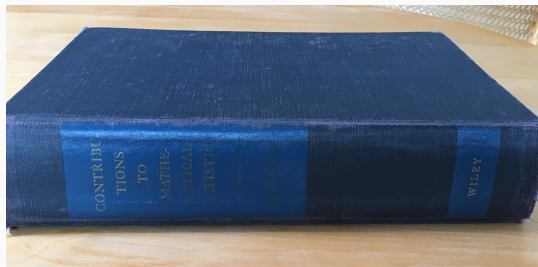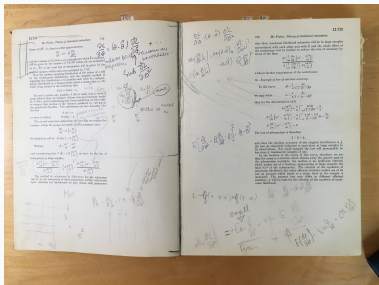# Fisher's contributions to mathematical statistics
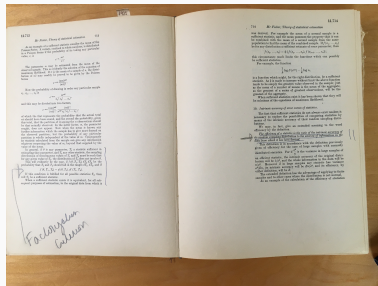
Nancy Reid
University of Toronto

April 21 2022
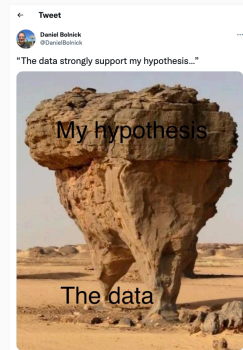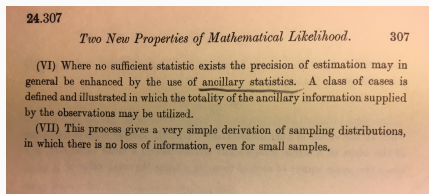
Lecture Notes in Statistics

Edited by D. Brillinger, S. Fienberg, J. Gani,
J. Hartigan, and K. Krickeberg

1

R. A. Fisher
An Appreciation

Edited by S. E. Fienberg and D. V. Hinkley

Springer-Verlag
New York  Heidelberg  Berlin

- Statistics needs a healthy interplay between theory and applications
  - theory meaning <span style="color:red">foundations</span>, rather than theoretical analysis of specific techniques

- Foundation**s**?
- "A solid base, on which rests a large structure"

OED

  - must be continually tested against new applications

  - "the practical application of general theorems is a different art from their establishment by mathematical proof"







Tweet

Daniel Bolnick
@DanielBolnick

"The data strongly support my hypothesis..."

My hypothesis

The data

- 1922: On the mathematical foundations of theoretical statistics:
statistic/parameter, estimation, consistent, sufficient, efficient, likelihood, maximum likelihood estimate, information, intrinsic accuracy

- 1925: Theory of statistical estimation:
all of the above, scoring algorithm, loss of information, ancillary

- 1934: Two new properties of mathematical likelihood:
conditional inference, location model, ancillary configuration, recovery of information, exponential family, distribution of sufficient estimate, uniformly most powerful tests

$\frac{n}{\partial \theta_2} \log f$

for all values of $x$ from $-\infty$ to $\infty$, and so to the average amount of information contained in a sample of $n$ observations.

*Summary.*

(I) Reasons are given for the use of mathematical likelihood in problems of inductive inference.

(II) When a statistic exists, satisfying the criterion of sufficiency, the likelihood function involves only that statistic.

(III) An example is given of a sufficient statistic, and its sampling distribution is expressed in terms of the likelihood function.

(IV) This property is generalized for all cases of simple estimation, where a sufficient statistic exists.

**24.307**

*Two New Properties of Mathematical Likelihood.* 307

(VI) Where no sufficient statistic exists the precision of estimation may in general be enhanced by the use of ancillary statistics. A class of cases is defined and illustrated in which the totality of the ancillary information supplied by the observations may be utilized.

(VII) This process gives a very simple derivation of sampling distributions, in which there is no loss of information, even for small samples.

## Likelihood Inference

- Model: $Y \sim f(y; \theta), \theta \in \mathbb{R}^p, y \in \mathbb{R}^n$

- Likelihood function: $L(\theta; y) \propto f(y; \theta)$

- Maximum likelihood estimator $\hat{\theta} = \hat{\theta}(y) = \arg \sup_\theta L(\theta; y) = \arg \sup_\theta \log\{L(\theta; y)\}$

- Observed and expected Fisher information
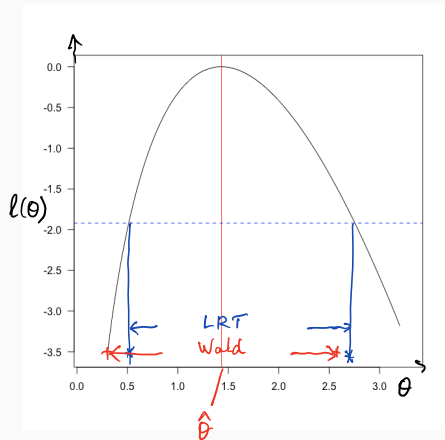  $$j(\theta) = -\partial^2 \ell(\theta; y)/\partial \theta^2, \quad i(\theta) = \mathrm{E}\{-\partial^2 \ell(\theta; Y)/\partial \theta^2\}$$
  $\ell(\theta; y) = \log L(\theta; y)$

- in "large samples" $\hat{\theta} \stackrel{.}{\sim} N_p\{\theta, j^{-1}(\hat{\theta})\}$    equivalently $N_p\{\theta, i^{-1}(\theta)\}$

- in "large samples" $\ell'(\theta) \stackrel{.}{\sim} N_p\{0, i^{-1}(\theta)\}$

- in "large samples" $\hat{\theta} \overset{\cdot}{\sim} N_p\{\theta, j^{-1}(\hat{\theta})\}$
- in "large samples" $\ell'(\theta) \overset{\cdot}{\sim} N_p\{0, i^{-1}(\theta)\}$
- in "large samples" $2\{\ell(\hat{\theta}) - \ell(\theta)\} \overset{\cdot}{\sim} \chi_p^2$

- $y = (y_1, \ldots, y_n); \quad y_i \sim Gamma(\theta, 1)$
- $L(\theta) = \prod_{i=1}^{n} y_i^{\theta-1} e^{-y_i} / \Gamma(\theta)$

- $\psi(\hat{\theta}) = \frac{1}{n} \sum \log(y_i)$ $\qquad \psi = \log \Gamma'$
- $j(\hat{\theta}) = n\psi'(\hat{\theta})$

## Small samples

- can we find the exact distribution of the maximum likelihood estimator?

- special case 1: location model $Y_i \sim f(y_i - \theta), i = 1, \ldots, n; \theta \in \mathbb{R}$      Fisher 1934

- ancillary statistic    $a = (y_1 - \hat{\theta}, \ldots, y_n - \hat{\theta})$      $\sum(\partial/\partial\theta)\log\{f(y_i; \hat{\theta})\} = 0$

- special case 2: exponential family model $Y_i \sim \exp\{s(y)^T\theta - nc(\theta)\}h(y)$      Fisher 1925

- sufficient statistic $s = s(y)$ is 'matched' to $\theta$      same dimension
- maximum likelihood estimate is sufficient      likelihood map is sufficient

- can we find the exact distribution of the maximum likelihood estimator?

- special case 1: $y \to (\hat{\theta}, a)$                                              Fisher 1934

$$f(\hat{\theta} \mid a; \theta) = \frac{L(\theta; \hat{\theta}, a)}{\int L(\theta; \hat{\theta}, a) d\theta} = \frac{\exp\{\ell(\theta; \hat{\theta}, a)\}}{\int \exp\{\ell(\theta; \hat{\theta}, a)\} d\theta}$$

- special case 2: $y \to s$

$$f(s; \theta) = \exp\{s^T \theta - nc(\theta)\}\tilde{h}(s) \qquad s = nc'(\hat{\theta}), \quad \tilde{h}(s) = \int ...dy$$

- general case                                              $1 + O(n^{-3/2})$

$$f(\hat{\theta}; \theta \mid a) \doteq c|j(\hat{\theta})|^{1/2} \exp\{\ell(\theta; \hat{\theta}, a) - \ell(\hat{\theta}; \hat{\theta}, a)\}$$

- we know the distribution of the maximum likelihood estimator

  to a high order of approximation

- how do we use this for inference?

- find values of $\theta$ that are consistent with our observed data     confidence intervals

- find the probability for a given $\theta_o$ of observing a result
  "as or more extreme than our observed data" (F 1925)     *p*-values

- computationally feasible?     $f(\hat{\theta}; \theta \mid a) \doteq c|j(\hat{\theta})|^{1/2} \exp\{\ell(\theta; \hat{\theta}, a) - \ell(\hat{\theta}; \hat{\theta}, a)\}$

- models with many parameters: $\theta = (\psi, \lambda)$,     $\ell_{\mathsf{p}}(\psi) = \ell(\psi, \hat{\lambda}_\psi)$

  profile log-likelihood

## Models with many parameters

- statistical models with many parameters $\theta = (\psi, \lambda)$       parameter of interest
                                                        nuisance parameter(s)

- profile or concentrated log-likelihood function $\ell_{\text{prof}}(\psi) \equiv \ell(\psi, \hat{\lambda}_\psi; y)$    $\hat{\lambda}_\psi = \arg\sup_\lambda L(\psi, \lambda)$

- now use "large samples" theory on $\ell_{prof}(\psi)$        approximation can be very poor

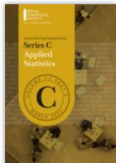- can sometimes isolate parameters in a marginal or conditional distribution

$$\text{e.g.} \qquad f(y; \psi, \lambda) \propto f_c(s \mid t; \psi) f(t; \psi, \lambda)$$

                                                                   Fisher's exact test

- can approximate this conditional likelihood with relatively simple adjustments

                                                         B-N 1983; Cox R 1987

$$\ell_{\text{mod}}(\psi) = \underbrace{\ell_{\text{prof}}(\psi)}_{O(n)} - \underbrace{\frac{1}{2}\log|j_{\lambda\lambda}(\psi, \hat{\lambda}_\psi)|}_{O(1)} + \underbrace{M(\psi)}_{O(1/\sqrt{n})}$$

**Volume 71, Issue 2**

Pages: **269-492**
**March 2022**

Series C
Applied Statistics
C

< Previous Issue

## Articles

🔓 Open Access

**Assessing seismic origin of geological features by fitting equidistant parallel lines**

P.E. Jupp, I.B.J. Goudie, R.A. Batchelor, R.J.B. Goudie

First Published: 17 April 2022

Abstract | Full text | PDF | References | Request permissions

🔓 Open Access

**Modelling the extremes of seasonal viruses and hospital congestion: The example of flu in a Swiss hospital**

Setareh Ranjbar, Eva Cantoni, Valérie Chavez-Demoulin, Giampiero Marra, Rosalba Radice, Katia Jaton

First Published: 13 April 2022

Abstract | Full text | PDF | References | Request permissions

🔓 Open Access

**The modelling of movement of multiple animals that share behavioural features**

Gianluca Mastrantonio

"you are not thinking in the right way"

- "investigate the use of the marginal likelihood function for model specification"

  Fong & Holmes 2020

- "maximum likelihood estimates are obtained from a multivariate Poisson regression model" Muñoz-Pichardo et al. 2021

- "a penalized likelihood approach to integrate high-dimension subject-level information along low-dimensional aggregate information" Sheng, Huang & Kim 2021

- "explores retrospective and prospective likelihood in terms of power of the score tests" Liu et al. 2020

- "likelihood ratio test for sequential change-point detection" Dette & Gössman 2020

- "proposed a modified profile likelihood method" for genetic association studies

  Zhang et al. 2020

- "a variant of the maximum likelihood estimator using a subset of the data …resulting estimator is still consistent" Ekvall & Jones 2021

- "small-sample bias correction for [the variance of] the maximum likelihood estimator"

  Ozenne et al. 2020

- …

- spatial dependence, nested sampling designs, high-dimensional parametrization

- pseudo-likelihood builds distribution from local characteristics

$$L_{\text{pseudo}}(\theta) = \prod_{j \in \mathcal{N}(y_i)} f(y_i \mid y_j; \theta)$$
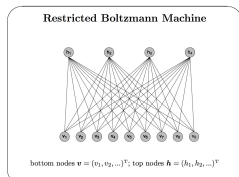
- example spatial modelling

Besag 74

```
.   ×   .   ×   .   ×   .   ×   .   ×
×   .   ×   .   ×   .   ×   .   ×   .
.   ×   .   ×   .   ×   .   ×   .   ×
×   .   ×   .   ×   .   ×   .   ×   .
```
Fig. 1. Coding pattern for a first-order scheme.

- example Boltzmann machine

Zhu

Restricted Boltzmann Machine

bottom nodes $\boldsymbol{v} = (v_1, v_2, \ldots)^{\mathsf{T}}$; top nodes $\boldsymbol{h} = (h_1, h_2, \ldots)^{\mathsf{T}}$

- spatial dependence, nested sampling designs, high-dimensional parametrization

- composite likelihood combines lower-dimensional marginal densities

-
$$L_{\text{composite}}(\theta) = \prod_{i=1}^{n} \prod_{j<k} f_2(y_{ij}, y_{ik}; \theta)$$

- example longitudinal data – each subject measured at several time points
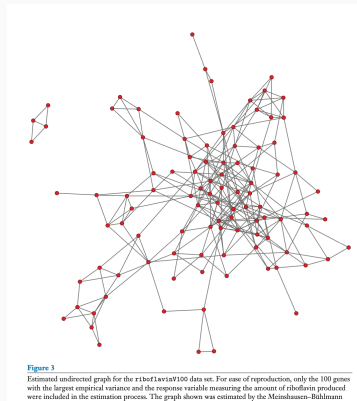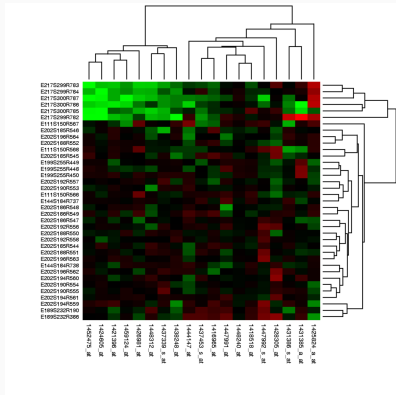
Renard 04; Kuk & Nott 03

- latent variable: $w_{ir} = x_{ir}'\beta + z_{ir}'b_i + \epsilon_{ir}, \quad \epsilon_{ir} \sim N(0, 1)$
- binary observations: $y_{ir} = \mathbb{1}(w_{ir} > 0); \quad r = 1, \dots m; i = 1, \dots n$      dependent binary vector $y_i$
- probit model: $Pr(y_{ir} = 1 \mid b_i) = \Phi(x_{ir}'\beta + b_i); \quad b_i \sim N_q(0, \Sigma_b)$      regression

$$L_{composite}(\beta, \Sigma_b; y) = \prod_{i=1}^{n} \prod_{j<k} P_{11,i}^{y_{ir}y_{is}} P_{10,i}^{y_{ir}(1-y_{is})} P_{01,i}^{(1-y_{ir})y_{is}} P_{00,i}^{(1-y_{ir})(1-y_{is})}$$

$P_{11,i}, P_{10,i}$, etc. evaluated using $\Phi_2(\cdot, \cdot; \rho_{irs})$

- spatial dependence, nested sampling designs, high-dimensional parametrization





**Figure 3**
Estimated undirected graph for the `riboflavinV100` data set. For ease of reproduction, only the 100 genes with the largest empirical variance and the response variable measuring the amount of riboflavin produced were included in the estimation process. The graph shown was estimated by the Meinshausen–Bühlmann

Andrade, Wikipedia

Buhlmann et al 14

## High-dimensional parameters

- new limit results, e.g. $\hat{\theta} \xrightarrow{d} N(\theta + \text{bias}, \sigma^2 \times \text{adjustment})$  <span style="float:right">Sur & Candès, Fan et al.</span>

- higher order approximations allows $p = O(n^{\alpha}), \alpha < 1/2$  <span style="float:right">Tang 22</span>

- sparsity – $\mathcal{S} \equiv \{j; \theta_j \neq 0\}; |\mathcal{S}| = s < n$

  - enforce sparsity, e.g. Lasso

  - discover sparsity, e.g. Battey 22

  - isolate parameter(s) of interest Battey & R 22; McCormack et al 19

## Various types of 'likelihood'

- likelihood, marginal likelihood, conditional likelihood, profile likelihood
  adjusted profile likelihood

- pseudo-likelihood, composite likelihood

- semi-parametric likelihood, partial likelihood

- empirical likelihood, penalized likelihood

- bootstrap likelihood, *h*-likelihood, weighted likelihood, quasi-likelihood, local
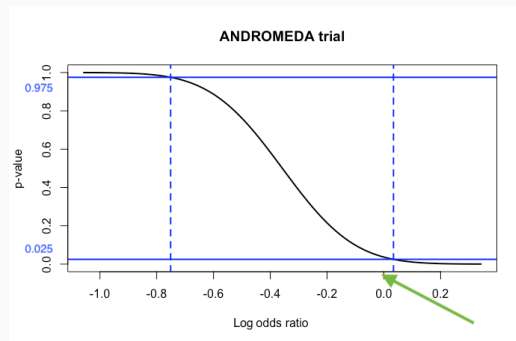  likelihood, sieve likelihood, simulated likelihood

- ANDROMEDA Trial: randomized clinical trial to compare two treatments for septic shock

  Hernandez et al 2019

- estimated hazard ratio 0.75 [0.55, 1.02]   after adjusting for confounders
- 2-sided p-value 0.06   34.9% vs 43.4% unadjusted

- Discussion: " a peripheral perfusion-targeted resuscitation strategy
  did not result in a significantly lower 28-day mortality
  when compared with a lactate level-targeted strategy"

- Abstract: "Among patients with septic shock, a resuscitation strategy targeting
  normalization of capillary refill time, compared with a strategy targeting serum lactate
  levels, did not reduce all-cause 28-day mortality."

## ANDROMEDA trial

|       | Died | Lived |     |
|-------|------|-------|-----|
| New   | 74   | 138   | 212 |
| Old   | 92   | 120   | 212 |
| Total | 166  | 258   | 424 |

2-sided *p*-value = 0.07

likelihood ratio test
no adjustment for covariates



90% confidence interval: [ −0.688, −0.030 ]
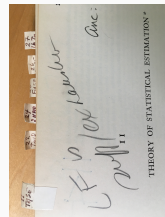95% confidence interval: [ −0.751,  0.034 ]
99% confidence interval: [ −0.825,  0.107 ]

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.]

$$\mathrm{d}f = -\frac{\partial}{\partial\theta}F(Y,\theta)\mathrm{d}\theta$$

fiducial probability density for $\theta$, given statistic $Y$

"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"

The First Workshop on BFF Inference and Statistical Foundations
(BFF 2014)

November 10 – November 14, 2014

BFF
CONFERENCE

TORONTO
MAY 2-4, 2022
HYBRID

SEVENTH BAYESIAN, FIDUCIAL &
FREQUENTIST (BFF) CONFERENCE

DEDICATED TO THE MEMORY
OF PROFESSOR DONALD A.S. FRASER

Distributions for parameters

- posterior distribution   Bayes 1763
- objective Bayes

- fiducial probability   Fisher 1930
- generalized fiducial inference, fiducial prediction, functional models, "slice-and-dice", …

- confidence distribution   Cox 1958
- confidence distributions / curves

- structural probability   Fraser 1964
- approximate significance functions

- belief functions   Dempster 1967
- inferential models

- study the structure of models which give 'valid' fiducial inference                Dawid 22

- change the modelling framework so fiducial arguments can be developed
  more cohesively                                                                     Lang 22

- generalized fiducial density                                                        Hannig 09 ff
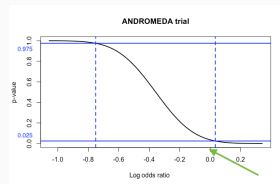
$$r(\theta; y^o) \propto \underbrace{f(y^o; \theta)}_{\textit{likelihood}} \underbrace{J(y^o, \theta)}_{\text{"prior"}}$$

- these methods all founder in models with many parameters

- unless each parameter of interest can be "measured separately"                      Fisher

- inference is intuitive
- combines easily with decision theory
- de-emphasizes point estimation and arbitrary cut-offs



- "it's tempting to conclude that $\mu$ is more likely to be near the middle of this interval, and if outside, not very far outside"  Cox 2006

- "assigns probability 0.05 to $\theta$ lying between the upper endpoints of the 0.90 and 0.95 confidence intervals, etc."  Efron 1993

- all inference statements become probability statements about unknowns  hmm…

- probability to describe physical haphazard variability                       aleatory/empirical
  - 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

- avoid apparent discoveries based on spurious patterns

- shed light on the structure of the problem

- obtain calibrated inferences about interpretable parameters

- provide a realistic assessment of precision

- understand when/why methods work/fail

- something that works

- gives 'sensible' answers

- not too sensitive to model assumptions

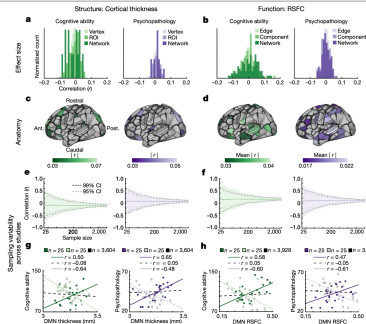- computable in reasonable time

- provides interpretable parameters
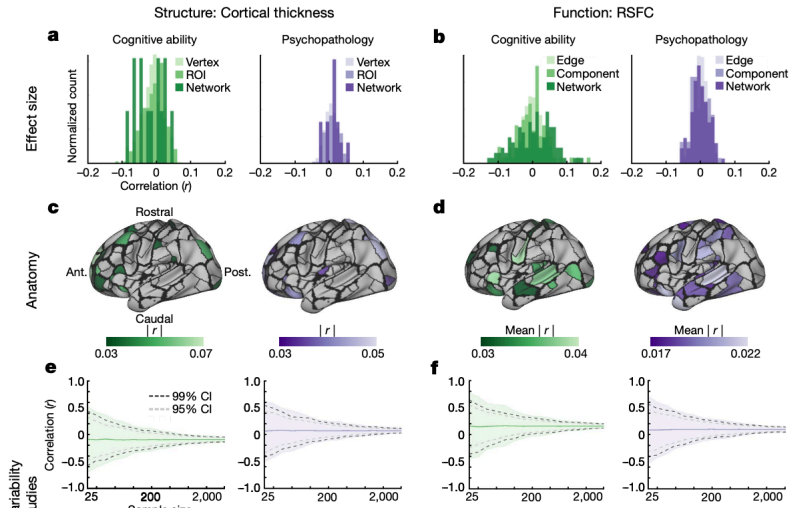


Fig. 1 | Effect sizes and sampling variability of univariate brain-wide associations. ABCD Study sample data (n = 3,928). a, b, Effect sizes were estimated using standard correlations 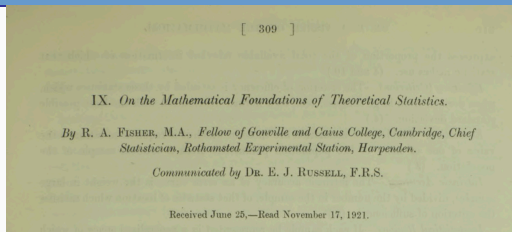(bivariate linear r). Brain-wide 135, 200, 265, 375, 525, 725, 1,000, 1,430, 2,000, 2,800 and 3,604 (3,928 for cortical thickness)) of the largest brain-wide association for each brain– behavioural phenotype pair, for cortical thickness with cognitive ability (left

## Article

**Reproducible brain-wide association studies require thousands of individuals**

- "relate population variability in brain features"                                    eg functional connectivity

- "and behavioural phenotypes"                                                         eg cognitive ability

- "across all univariate brain-wide associations, the largest correlation that replicated out-of-sample was $|r| = 0.16$"

- "at $n = 25$, the 99% confidence interval for univariate associations was $r \pm 0.52$."

- "Bias in favour of significant, larger BWAS effects has limited the publication of null results, perpetuating inflated effect sizes … "

[ 309 ]

IX. *On the Mathematical Foundations of Theoretical Statistics.*

By R. A. Fisher, M.A., *Fellow of Gonville and Caius College, Cambridge, Chief Statistician, Rothamsted Experimental Station, Harpenden.*

*Communicated by* Dr. E. J. Russell, F.R.S.

*Received June 25,—Read November 17, 1921.*

*The Annals of Statistics*
1976, Vol. 4, No. 3, 441–500

## ON REREADING R. A. FISHER

By Leonard J. Savage[1,2]

*Yale University*

「Fisher's contributions to statistics are surveyed. His background, skills, temperament, and style of thought and writing are sketched. His mathematical and methodological contributions are outlined. More atten-

Savage: "there is a world of R.A. Fisher at once very near to and very far from the world of modern statisticians … research for the fun of it is abundant and beautiful in Fisher's writings"

Fraser: "One important characteristic of Fisher was his ability to move into new areas of statistics, suggesting concepts and methods … left the theory open to modification and development"

Efron: "This paper makes me happy to be a statistician"

# THANK YOU