Topics in Likelihood Inference

STA2212H S LEC9101 4508 H

Nancy Reid University of Toronto

January 12, 2022



STA 4508: Topics in Likelihood Inference Winter 2022

Wednesdays 14.00-17.00, Jan 12 - Feb 16, SS 2120

Topics

- - 2. Likelihood for semi-parametric and non-parametric models: proportional hazards regression, partially linear models, penalized likelihood;
 - 3. Composite likelihood: definition, summary statistics, asymptotic theory; applications
 - 4. Likelihood inference for p > n;
 - 5. Simulated likelihoods, indirect inference and approximate Bayesian computation

Running list of references and background reading

Review Papers

- • Reid, N. (2013) Aspects of likelihood inference Bernoulli 19, 1404-1418.
- Reid, N. (2010) Likelihood Inference Wiley Interdisciplinary Reviews in Computational Statistics 5, 517-525. (I need to use Preview to view this, rather than Adobe.)

Likelihood Basics

- --->• Varin, C., Reid, N. and Yi, G. (20xx) (VRY) Ch 1
 - Davison, A.C. (2003) Statistical Models (SM) Cambridge University Press. -- Ch 4
 - Barndorff-Nielsen, O.E. and Cox, D.R. (1994) Inference and Asymptotics (BNC) Chapman and Hall. -- Ch 2.2
 - Cox, D.R. and Hinkley, D.V. (1974) Theoretical Statistics (CH) Chapman and Hall. -- Ch 2.1 (i), (ii)
 - Cox, D.R. (2006) Principles of Statistical Inference (Cox) -- Ch.2.1



- 1. likelihood, marginal and conditional likelihood, profile likelihood, adjusted profile
- 2. semi-parametric likelihood, partial likelihood
- 3. quasi-likelihood, composite likelihood

- 4. empirical likelihood, penalized likelihood
- 5. simulated likelihood, indirect inference
- 6. bootstrap likelihood, *h*-likelihood, weighted likelihood, pseudo-likelihood, local likelihood, sieve likelihood

- 1. likelihood, marginal and conditional likelihood, profile likelihood, adjusted profile
- 2. semi-parametric likelihood, partial likelihood
- 3. quasi-likelihood, composite likelihood

- 4. empirical likelihood, penalized likelihood
- 5. simulated likelihood, indirect inference
- 6. bootstrap likelihood, *h*-likelihood, weighted likelihood, pseudo-likelihood, local likelihood, sieve likelihood

- 1. likelihood, marginal and conditional likelihood, profile likelihood, adjusted profile
- 2. semi-parametric likelihood, partial likelihood
- 3. quasi-likelihood, composite likelihood

- 4. empirical likelihood, penalized likelihood
- 5. simulated likelihood, indirect inference
- 6. bootstrap likelihood, *h*-likelihood, weighted likelihood, pseudo-likelihood, local likelihood, sieve likelihood

- 1. likelihood, marginal and conditional likelihood, profile likelihood, adjusted profile
- 2. semi-parametric likelihood, partial likelihood
- 3. quasi-likelihood, composite likelihood

- 4. empirical likelihood, penalized likelihood
- 5. simulated likelihood, indirect inference
- 6. bootstrap likelihood, *h*-likelihood, weighted likelihood, pseudo-likelihood, local likelihood, sieve likelihood

- 1. likelihood, marginal and conditional likelihood, profile likelihood, adjusted profile
- 2. semi-parametric likelihood, partial likelihood
- 3. quasi-likelihood, composite likelihood

4. empirical likelihood, penalized likelihood

misspecified models

- Exercises (each week) ~ 207. R to react di ~
- 5. simulated likelihood, indirect inference
 6. bootstrap likelihood, h-likelihood, weighted likelihood, pseudo-likelihood, local likelihood, sieve likelihood

STA 4508 January 12 2022

• Principle: "The probability model and the choice of [parameter] serve to translate a subject-matter question into a mathematical and statistical one" Cox, 2006, p.3

- Principle: "The probability model and the choice of [parameter] serve to translate a subject-matter question into a mathematical and statistical one" Cox, 2006, p.3
- likelihood function is proportional to the probability model

- Principle: "The probability model and the choice of [parameter] serve to translate a subject-matter question into a mathematical and statistical one" Cox, 2006, p.3
- likelihood function is proportional to the probability model
- inference based on the likelihood function is widely accepted

- Principle: "The probability model and the choice of [parameter] serve to translate a subject-matter question into a mathematical and statistical one" Cox, 2006, p.3
- likelihood function is proportional to the probability model
- inference based on the likelihood function is widely accepted
- provides more than point estimate or test of point hypothesis

- Principle: "The probability model and the choice of [parameter] serve to translate a subject-matter question into a mathematical and statistical one" Cox, 2006, p.3
- likelihood function is proportional to the probability model
- inference based on the likelihood function is widely accepted
- provides more than point estimate or test of point hypothesis
- models needed for applications are more and more complex

- Principle: "The probability model and the choice of [parameter] serve to translate a subject-matter question into a mathematical and statistical one" Cox, 2006, p.3
- likelihood function is proportional to the probability model
- inference based on the likelihood function is widely accepted
- provides more than point estimate or test of point hypothesis
- models needed for applications are more and more complex
- need some analogues to the likelihood function for these complex settings

• Parametric model: $f(y; \theta), y \in \mathcal{Y}, \theta \in \Theta \subset \mathbb{R}^p$

$$\begin{array}{ll} \text{lel:} f(y;\theta), & y \in \mathcal{Y}, \theta \in \Theta \subset \mathbb{R}^p & f \text{ density for } \mathcal{Y} \text{ n.v. } \in \mathcal{Y} \text{ (s.s.)} \\ & \Theta e^{-\mathcal{Y}\Theta} & y = 0 \ ; \Theta > 0 \\ & (\overset{\circ}{y}) \Theta^{\mathcal{Y}} (1-\Theta)^{n-\mathcal{Y}} & y \in \{0, \dots, n\} \end{array}$$

- Parametric model: $f(y; \theta), y \in \mathcal{Y}, \theta \in \Theta \subset \mathbb{R}^p$
- Likelihood function

Court can't

- Parametric model: $f(y; \theta), y \in \mathcal{Y}, \theta \in \Theta \subset \mathbb{R}^p$
- Likelihood function

 $L(\theta; y) = f(y; \theta)$, or $L(\theta; y) = c(y)f(y; \theta)$, or $L(\theta; y) \propto f(y; \theta)$

• typically, $y = (y_1, ..., y_n)$ $X_1, ..., X_n$ i = 1, ..., n

- Parametric model: $f(y; \theta), y \in \mathcal{Y}, \theta \in \Theta \subset \mathbb{R}^p$
- Likelihood function

- typically, $y = (y_1, \ldots, y_n)$ x_1, \ldots, x_n $i = 1, \ldots, n$
- $f(y;\theta)$ or $f(y \mid x;\theta)$ is joint density $f(y_i) \approx id^{-1} + f(y_i;\theta) = \pi f(y_i;\theta \mid x_i)$

- Parametric model: $f(y; \theta), y \in \mathcal{Y}, \theta \in \Theta \subset \mathbb{R}^p$
- Likelihood function

- typically, $y = (y_1, ..., y_n)$ $x_1, ..., x_n$ i = 1, ..., n
- $f(y; \theta)$ or $f(y \mid x; \theta)$ is joint density
- under independence $L(\theta; y) \propto \prod f(y_i \mid x_i; \theta)$

- Parametric model: $f(y; \theta), y \in \mathcal{Y}, \theta \in \Theta \subset \mathbb{R}^p$
- Likelihood function

- typically, $y = (y_1, ..., y_n)$ $x_1, ..., x_n$ i = 1, ..., n
- $f(y; \theta)$ or $f(y \mid x; \theta)$ is joint density
- under independence $L(\theta; y) \propto \prod f(y_i \mid x_i; \theta)$
- log-likelihood $\ell(\theta; y) = \log L(\theta; y) = \sum \log f(y_i \mid x_i; \theta)$

• Parametric model:
$$f(y; \theta)$$
, $y \in \mathcal{Y}, \theta \in \Theta \subset \mathbb{R}^p$ \leftarrow { $f(y; \theta); \theta \in \Theta$

Likelihood function

- typically, $y = (y_1, \ldots, y_n)$ x_1, \ldots, x_n $i = 1, \ldots, n$
- $f(y; \theta)$ or $f(y \mid x; \theta)$ is joint density
- under independence $L(\theta; y) \propto \prod f(y_i \mid x_i; \theta)$
- log-likelihood $\ell(\theta; y) = \log L(\theta; y) = \sum \log f(y_i \mid x_i; \theta)$
- θ could have dimension p > n (e.g. genetics), or $p \uparrow n$, or



- Parametric model: $f(y; \theta), y \in \mathcal{Y}, \theta \in \Theta \subset \mathbb{R}^p$
- Likelihood function

 $L(\theta; y) = f(y; \theta), \text{ or } L(\theta; y) = c(y)f(y; \theta), \text{ or } L(\theta; y) \propto f(y; \theta)$

- typically, $y = (y_1, \dots, y_n)$ X_1, \dots, X_n $i = 1, \dots, n$
- $f(y; \theta)$ or $f(y | x; \theta)$ is joint density
- under independence $L(\theta; y) \propto \prod f(y_i \mid x_i; \theta)$
- log-likelihood $\ell(\theta; y) = \log L(\theta; y) = \sum \log f(y_i \mid x_i; \theta)$
- θ could have dimension p > n (e.g. genetics), or $p \uparrow n$, or
- θ could have infinite dimension e.g. m some applications STA 4508 January 12 2022 θ could have infinite dimension e.g. $f_i = m(\alpha_i) + \Xi_i$ $\xi = n((0, \sigma^2))$ $g_i = m(\alpha_i) + \Xi_i$ $\xi = n((0, \sigma^2))$ $m(\omega)$ is a function of f the $f = \sigma f(\alpha)$ $(\sim m f \omega \sigma m)$

4

- Parametric model: $f(y; \theta), y \in \mathcal{Y}, \theta \in \Theta \subset \mathbb{R}^p$
- Likelihood function

 $L(\theta; y) = f(y; \theta), \text{ or } L(\theta; y) = c(y)f(y; \theta), \text{ or } L(\theta; y) \propto f(y; \theta)$

- typically, $y = (y_1, ..., y_n)$ $x_1, ..., x_n$ i = 1, ..., n
- $f(y; \theta)$ or $f(y | x; \theta)$ is joint density
- under independence $L(\theta; y) \propto \prod f(y_i \mid x_i; \theta)$
- log-likelihood $\ell(\theta; y) = \log L(\theta; y) = \sum \log f(y_i \mid x_i; \theta)$

ዋ

- θ could have dimension p > n (e.g. genetics), or $p \uparrow n$, or
- + $\boldsymbol{\theta}$ could have infinite dimension e.g.
- regular model p < n and p fixed as n increases

most of the couse is in repulse

STA 4508 January 12 2022

Examples

•
$$y_i \sim N(\mu, \sigma^2)$$
:
 $y_i \stackrel{i}{\leftarrow} N(\mu, \sigma^2)$:
 $f_{i=1} \int \sigma^{-n} \exp\{-\frac{1}{2\sigma^2}(y_i - \mu)^2\}$
 $f_{i=1} \int \sigma^{-n} \exp\{-\frac{1}{2\sigma^2}(y_i - \mu)^2\}$

Examples

• $y_i \sim N(\mu, \sigma^2)$:

$$L(\theta; \mathbf{y}) = \prod_{i=1}^{n} \sigma^{-n} \exp\{-\frac{1}{2\sigma^2}(\mathbf{y}_i - \mu)^2\}$$

• $E(y_i) = x_i^T \beta$: $L(\theta; y) = \prod_{i=1}^n \sigma^{-n} \exp\{-\frac{1}{2\sigma^2}(y_i - x_i^T \beta)^2\} \qquad (f = 0, \sigma^2)$

Examples

•
$$y_i \sim N(\mu, \sigma^2)$$
:

$$L(\theta; y) = \prod_{i=1}^n \sigma^{-n} \exp\{-\frac{1}{2\sigma^2}(y_i - \mu)^2\}$$
• $E(y_i) = x_i^T \beta$:

$$L(\theta; y) = \prod_{i=1}^n \sigma^{-n} \exp\{-\frac{1}{2\sigma^2}(y_i - x_i^T \beta)^2\}$$
• $E(y_i) = m(x_i), \quad m(x) = \sum_{j=1}^J \phi_j B_j(x)$:

$$L(\theta; y) = \prod_{i=1}^n \sigma^{-n} \exp\{-\frac{1}{2\sigma^2}(y_i - \sum_{j=1}^J \phi_j B_j(x_i))^2\}$$

$$Q = (\varphi_i, \dots, \varphi_j, \sigma^2)$$

... examples

$$\begin{array}{ccc} & \underbrace{y_{i}} = \mu + \rho(y_{i-1} - \mu) + \epsilon_{i}, & \epsilon_{i} \sim \mathcal{N}(0, \sigma^{2}): & i = o_{1}, \dots, p \end{array} \\ & \underbrace{e(\psi) f(\psi; \theta)} = L(\theta; y) = \prod_{i=1}^{n} f(y_{i} \mid y_{i-1}; \theta) f_{0}(y_{0}; \theta) \qquad \begin{array}{c} & \underbrace{\theta = (e, \sigma^{2})} \\ & \underbrace{\theta = (e, \sigma^{2})} \end{array} \end{array}$$



•
$$\mathbf{y}_i = \mu + \rho(\mathbf{y}_{i-1} - \mu) + \epsilon_i, \quad \epsilon_i \sim N(\mathbf{0}, \sigma^2)$$
:

$$L(\theta; \mathbf{y}) = \prod_{i=1}^{n} f(\mathbf{y}_i \mid \mathbf{y}_{i-1}; \theta) f_0(\mathbf{y}_0; \theta)$$

•
$$y_1, \ldots, y_n$$
 i.i.d. observations from a $U(0, \theta)$ distribution:

$$f(y, \theta) = \frac{1}{\theta}, o(y) = 0$$

$$L(\theta; y) = \prod_{i=1}^{n} \theta(y) = -\eta \left[\frac{1}{\theta} + \frac{1}{\theta} \right] \left[\frac{1}{\theta} + \frac{1}{\theta} + \frac{1}{\theta} + \frac{1}{\theta} + \frac{1}{\theta} \right] \left[\frac{1}{\theta} + \frac{1$$

... examples

• y_1, \ldots, y_n are the times of jumps of a non-homogeneous Poisson process with rate function $\lambda(\cdot)$:

$$\ell\{\lambda(\cdot); y\} = \sum_{i=1}^{n} \log\{\lambda(y_i)\} - \int_{0}^{\tau} \lambda(u) du, \quad 0 < y_1 < \dots < y_n < \tau$$
parameter Θ
of f . -dim

... examples

• y_1, \ldots, y_n are the times of jumps of a non-homogeneous Poisson process with rate function $\lambda(\cdot)$:

$$\ell\{\lambda(\cdot); y\} = \sum_{i=1}^{n} \log\{\lambda(y_i)\} - \int_{0}^{\tau} \lambda(u) du, \quad 0 < y_1 < \cdots < y_n < \tau$$

Davison, §6.5

negative cross-entropy
$$p_{ic} = p(x_{ic}; \theta)$$
, as above $= P_{1}(y_{ic} = 1)$

Hastie et al., Ch. 7

4.1 · Likelihood



95

4 · Likelihood

Figure 4.2 Cauchy

likelihood for the spring

likelihood and log

failure data at stress

950N/mm2.



96

• example: clustered binary data

Renard et al. (2004)

• example: clustered binary data

Renard et al. (2004)

• latent variable: $z_{ir} = x_{ir}^T \beta + b_i + \epsilon_{ir}$, $b_i \sim N(0, \sigma_b^2)$, $\epsilon_{ir} \sim N(0, 1)$

• example: clustered binary data

Renard et al. (2004)

- latent variable: $z_{ir} = x_{ir}^T \beta + b_i + \epsilon_{ir}$, $b_i \sim N(0, \sigma_b^2)$, $\epsilon_{ir} \sim N(0, 1)$
- $r = 1, ..., n_i$: observations in a cluster/family/school... i = 1, ..., n clusters

• example: clustered binary data

Renard et al. (2004)

- latent variable: $z_{ir} = x_{ir}^T \beta + b_i + \epsilon_{ir}, \quad b_i \sim N(0, \sigma_b^2), \quad \epsilon_{ir} \sim N(0, 1)$
- $r = 1, ..., n_i$: observations in a cluster/family/school... i = 1, ..., n clusters
- random effect b_i introduces correlation between observations in a cluster
Complicated likelihoods

• example: clustered binary data

Renard et al. (2004)

- latent variable: $z_{ir} = x_{ir}^T \beta + b_i + \epsilon_{ir}$, $b_i \sim N(0, \sigma_b^2)$, $\epsilon_{ir} \sim N(0, 1)$
- $r = 1, ..., n_i$: observations in a cluster/family/school... i = 1, ..., n clusters
- random effect b_i introduces correlation between observations in a cluster
- observations: $y_{ir} = 1$ if $z_{ir} > 0$, else 0

$Y_i \sim N(x_i^T(s_j \sigma_{+\sigma}^2))$ **Complicated likelihoods** y:=x:B+E: +b: composite lik. • example: clustered binary data Renard et al. (2004) • latent variable: $z_{ir} = x_{ir}^T \beta + b_i + \epsilon_{ir}$, $b_i \sim N(0, \sigma_b^2)$, $\epsilon_{ir} \sim N(0, 1)$ • $r = 1, ..., n_i$: observations in a cluster/family/school... i = 1, ..., n clusters random effect b_i introduces correlation between observations in a cluster corrid binary data • observations: $y_{ir} = 1$ if $z_{ir} > 0$, else 0 $\Phi(z) = \int^z \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx$ • $Pr(y_{ir} = 1 \mid b_i) = \Phi(x_{ir}^T \beta + \lambda_i) = p_i$ $L(\theta; y) = \prod_{i=1}^{n} \int_{-\infty}^{\infty} \prod_{r=1}^{n_i} p_i^{y_{ir}} (1-p_i)^{1-y_{ir}} \phi(b_i, \sigma_b^2) db_i$ f(yirlbi) $f(y_{i1},..., y_{inc}; \beta, \sigma_b^2) = f(y_i; \beta, \sigma_b^2)$ $L(0; y) = f(y_1) f(y_2) \cdots f(y_n)$ STA 4508 January 12 2022 10

Complicated likelihoods

• example: clustered binary data

Renard et al. (2004)

- latent variable: $z_{ir} = x_{ir}^T \beta + b_i + \epsilon_{ir}$, $b_i \sim N(0, \sigma_b^2)$, $\epsilon_{ir} \sim N(0, 1)$
- $r = 1, ..., n_i$: observations in a cluster/family/school... i = 1, ..., n clusters
- random effect b_i introduces correlation between observations in a cluster
- observations: $y_{ir} = 1$ if $z_{ir} > 0$, else 0

•
$$Pr(y_{ir} = 1 \mid b_i) = \Phi(x_{ir}^T \beta + b_i) = p_i$$

$$\Phi(z) = \int^z \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx$$

$$L(\theta; y) = \prod_{i=1}^n \int_{-\infty}^{\infty} \prod_{r=1}^{n_i} p_i^{y_{ir}} (1 - p_i)^{1 - y_{ir}} \phi(b_i, \sigma_b^2) db_i$$

• more general: $z_{ir} = x_{ir}^T \beta + w'_{ir} b_i + \epsilon_{ir}$

Renard et al. (2004) Multi-level probit models CSDA

• generalized linear geostatistical models

 $\mathsf{E}\{\mathsf{Y}(\mathsf{s}) \mid u(\mathsf{s})\} = g\{\mathsf{x}(\mathsf{s})^{\mathsf{T}}\beta + u(\mathsf{s})\}, \quad \mathsf{s} \in \mathcal{S} \subset \mathbb{R}^{d}, d \geq 2$

Diggle & Ribeiro, 2007

• generalized linear geostatistical models

$$\mathsf{E}\{\mathsf{Y}(\mathsf{s}) \mid u(\mathsf{s})\} = g\{\mathsf{x}(\mathsf{s})^{\mathsf{T}}eta + u(\mathsf{s})\}, \quad \mathsf{s} \in \mathcal{S} \subset \mathbb{R}^{d}, d \geq \mathsf{2}\}$$

Diggle & Ribeiro, 2007

 random intercept u is a realization of a stationary GRF, expected value o, covariance
 Gaussian random field

$$\operatorname{COV}\{u(\mathbf{S}), u(\mathbf{S}')\} = \sigma^2 \rho(\mathbf{S} - \mathbf{S}'; \alpha)$$

• generalized linear geostatistical models

$$\mathsf{E}{Y(s) \mid u(s)} = g{x(s)^T eta + u(s)}, \quad s \in \mathcal{S} \subset \mathbb{R}^d, d \geq 2$$

Diggle & Ribeiro, 2007

 random intercept u is a realization of a stationary GRF, expected value o, covariance
 Gaussian random field

$$\operatorname{cov}\{u(\mathbf{s}), u(\mathbf{s}')\} = \sigma^2 \rho(\mathbf{s} - \mathbf{s}'; \alpha)$$

• *n* observed locations $y = (y_1, \ldots, y_n)$ with $y_i = y(s_i)$

• generalized linear geostatistical models

$$\mathsf{E}\{\underbrace{Y(s)} \mid \underbrace{u(s)}_{\frown}\} = g\{\underbrace{x(s)^{\mathsf{T}}\beta}_{\uparrow} + \underbrace{u(s)}_{\uparrow}\}, \quad s \in \mathcal{S} \subset \mathbb{R}^{d}, d \geq 2$$

Diggle & Ribeiro, 2007

random intercept u is a realization of a stationary GRF, expected value o, Gaussian random field

$$\operatorname{cov}\{u(s), u(s')\} = \sigma^2 \rho(s - s'; \alpha)$$

 $\circ n \text{ observed locations } y = (y_1, \dots, y_n) \text{ with } y_i = y(s_i)$
 $\operatorname{corr} f = \hat{f}$
 space

• likelihood function \checkmark

$$L(\theta; y) = \int_{\mathbb{R}^n} \prod_{i=1}^n f(y_i \mid u_i; \theta) \underbrace{f(u_i; \theta)}_{N_d(o, \Sigma)} du_1 \dots du_n$$

STA 4508 January 12 2022

• generalized linear geostatistical models

$$\mathsf{E}\{\mathsf{Y}(\mathsf{s}) \mid u(\mathsf{s})\} = g\{\mathsf{x}(\mathsf{s})^{\mathsf{T}}eta + u(\mathsf{s})\}, \quad \mathsf{s} \in \mathcal{S} \subset \mathbb{R}^{d}, d \geq \mathsf{2}\}$$

Diggle & Ribeiro, 2007

 random intercept u is a realization of a stationary GRF, expected value o, covariance
 Gaussian random field

$$\operatorname{cov}\{u(\mathbf{s}), u(\mathbf{s}')\} = \sigma^2 \rho(\mathbf{s} - \mathbf{s}'; \alpha)$$

- *n* observed locations $y = (y_1, \ldots, y_n)$ with $y_i = y(s_i)$
- likelihood function

$$L(\theta; \mathbf{y}) = \int_{\mathbb{R}^n} \prod_{i=1}^n f(\mathbf{y}_i \mid u_i; \theta) \underbrace{f(u_i; \theta)}_{N_d(\mathbf{o}, \Sigma)} du_1 \dots du_n$$

• no factorization into lower dimensional integrals, as with previous example

• generalized linear geostatistical models

$$\mathsf{E}{Y(s) \mid u(s)} = g{x(s)^T eta + u(s)}, \quad s \in \mathcal{S} \subset \mathbb{R}^d, d \geq 2$$

Diggle & Ribeiro, 2007

 random intercept u is a realization of a stationary GRF, expected value o, covariance
 Gaussian random field

$$\operatorname{cov}\{\mathbf{u}(\mathbf{s}),\mathbf{u}(\mathbf{s}')\}=\sigma^{2}\rho(\mathbf{s}-\mathbf{s}';\alpha)$$

- *n* observed locations $y = (y_1, \ldots, y_n)$ with $y_i = y(s_i)$
- likelihood function

$$L(\theta; \mathbf{y}) = \int_{\mathbb{R}^n} \prod_{i=1}^n f(\mathbf{y}_i \mid u_i; \theta) \underbrace{f(u_i; \theta)}_{N_d(\mathbf{o}, \Sigma)} du_1 \dots du_n$$

• no factorization into lower dimensional integrals, as with previous example

Diggle & Ribeiro (2007) Model-based Geostatistics Springer ¹¹

STA 4508 January 12 2022

• Ising model:

$$f(y; \theta) = \exp(\sum_{(i,j)\in E} \theta_{ij} y_i y_j) \frac{1}{Z(\theta)}$$

• Ising model:

$$f(\mathbf{y};\theta) = \exp(\sum_{(i,j)\in E} \theta_{ij} \mathbf{y}_i \mathbf{y}_j) \frac{1}{Z(\theta)}$$

• $y_i = \pm 1$; binary property of a node *i* in a graph with *n* nodes

• Ising model:

$$f(y;\theta) = \exp(\sum_{(i,j)\in E} \theta_{ij} y_i y_j) \frac{1}{Z(\theta)}$$

- $y_i = \pm 1$; binary property of a node *i* in a graph with *n* nodes
- θ_{ij} measures strength of interaction between nodes *i* and *j*

• Ising model:

 $f(\mathbf{y};\theta) \neq \exp($

Ji i=1,...,1 not indit

- $y_i = \pm 1$; binary property of a node *i* in a graph with *n* nodes
- θ_{ij} measures strength of interaction between nodes *i* and *j*
- *E* is the set of edges between nodes

 $Z_{(0)} = Z_{ij}^{(0)}$



• Ising model:

$$f(\mathbf{y};\theta) = \exp(\sum_{(i,j)\in E} \theta_{ij} \mathbf{y}_i \mathbf{y}_j) \frac{1}{Z(\theta)}$$

- $y_i = \pm 1$; binary property of a node *i* in a graph with *n* nodes
- θ_{ij} measures strength of interaction between nodes *i* and *j*
- *E* is the set of edges between nodes
- partition function $Z(\theta) = \sum_{y} \exp(\sum_{(i,j) \in E} \theta_{ij} y_i y_j)$

Davison §6.2

Ravikumar et al. (2010).

High-dimensional Ising model selection... Ann. Statist. p.1287

STA 4508 January 12

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.

CONTENTS.

S	ectic	on a state of the	Page
	1.	The Neglect of Theoretical Statistics	310
	2.	The Purpose of Statistical Methods	311
	3.	The Problems of Statistics	313
	4.	Criteria of Estimation	316
	5.	Examples of the Use of Criterion of Consistency	317
	6.	Formal Solution of Problems of Estimation	323
	7.	Satisfaction of the Criterion of Sufficiency	330
	8.	The Efficiency of the Method of Moments in Fitting Curves of the Pearsonian Type III	332
	9.	Location and Scaling of Frequency Curves in general	338
2022	10.	The Efficiency of the Method of Moments in Fitting Pearsonian Curves	342

History



STA 4508 January 12 2022

know nothing whatever. We must return to the actual fact that one value of p, of the frequency of which we know nothing, would yield the observed result three times as frequently as would another value of p. If we need a word to characterise this relative property of different values of p, I suggest that we may speak without confusion of the *likelihood* of one value of p being thrice the likelihood of another, bearing always in mind that likelihood is not here used loosely as a synonym of probability, but simply to express the relative frequencies with which such values of the hypothetical quantity p would in fact yield the observed sample. $(y)f(x; \theta)$

$$\begin{array}{ccc} \mathcal{A}^{\circ} & \underline{\mathcal{P}}(\mathcal{Y}^{\circ}; \mathbf{0}) \\ & \underline{\mathcal{P}}(\mathcal{Y}^{\circ}; \mathbf{0}^{\circ}) \end{array} = 3 &= \frac{\mathcal{L}(\mathbf{0}^{\circ}; \mathcal{Y}^{\circ})}{\mathcal{L}(\mathbf{0}^{\circ}; \mathcal{Y}^{\circ})} \end{array}$$

• makes probability modelling central

- makes probability modelling central $f(\gamma; \circ)$
- emphasizes the inverse problem of reasoning from y^{o} to θ or $f(\cdot)$

 $L(0,\gamma)$

- makes probability modelling central
- emphasizes the inverse problem of reasoning from y^o to θ or $f(\cdot)$
- suggested by Fisher as a measure of plausibility

 $L(\hat{\theta})/L(\theta) \in (1,3)$ very plausible; $\widehat{L(\hat{\theta})}/L(\theta) \in (3, 10)$ implausible; $L(\hat{\theta})/L(\theta) \in (10,\infty)$ very implausible

Royall, 1994

Statistical Evidence: A likelihood paradigm

all
$$\Theta \in \Theta$$
 for chick $L(\hat{\Theta})$ 23 $\hat{\Theta} = \arg \sup L(\Theta)$
 $I \qquad 1 < L(\Theta) \qquad 0$
 $\subseteq \mathbb{R}$
muary 12 2022 Surfevel of "plansible" walkes

STA 4508 Jan

16

makes probability modelling central

 $\pi(\Theta(\gamma))$

- emphasizes the inverse problem of reasoning from y^{o} to θ or $f(\cdot)$
- suggested by Fisher as a measure of plausibility

Royall, 1994

$L(\hat{ heta})/L(heta)\in(extsf{1}, extsf{3})$	very plausible;
$L(\hat{ heta})/L(heta)\in (3,10)$	implausible;
$L(\hat{ heta})/L(heta)\in(extsf{10},\infty)$	very implausible

Statistical Evidence: A likelihood paradigm

• converts a 'prior' probability $\pi(\theta)$ to a posterior $\pi(\theta \mid y)$ via Bayes' formula

 $\frac{L(0;y)\pi(0)}{\int L(0;y)\pi(0)d\theta}$

$$\int \pi(o|y) do = 1$$

- makes probability modelling central
- emphasizes the inverse problem of reasoning from y^{o} to θ or $f(\cdot)$
- suggested by Fisher as a measure of plausibility

Royall, 1994

$L(\hat{ heta})/L(heta)\in(extsf{1}, extsf{3})$	very plausible;
$L(\hat{ heta})/L(heta)\in (3,10)$	implausible;
$L(\hat{ heta})/L(heta)\in(extsf{10},\infty)$	very implausible

Statistical Evidence: A likelihood paradigm

- converts a 'prior' probability $\pi(\theta)$ to a posterior $\pi(\theta \mid y)$ via Bayes' formula
- provides a conventional set of summary quantities for inference based on properties Æ mle mox. lit. ert. ô estid stal error of ô of the postulated model lop-tik ropro statistic for testorp hypotheres 16

STA 4508 January 12 2022

• likelihood function depends on data only through sufficient statistics

- likelihood function depends on data only through sufficient statistics
- "likelihood map is sufficient"

Fraser & Naderi, 2006

- likelihood function depends on data only through sufficient statistics
- "likelihood map is sufficient"

Fraser & Naderi, 2006

• gives exact inference in transformation models

- likelihood function depends on data only through sufficient statistics
- "likelihood map is sufficient"

Fraser & Naderi, 2006

- gives exact inference in transformation models
- "likelihood function as pivotal"

Hinkley, 1980

... why likelihood?

- likelihood function depends on data only through sufficient statistics
- "likelihood map is sufficient"
- gives exact inference in transformation models
- "likelihood function as pivotal"
- provides summary statistics with known limiting distribution



Hinkley, 1980

... why likelihood?

- likelihood function depends on data only through sufficient statistics
- "likelihood map is sufficient"
- gives exact inference in transformation models
- "likelihood function as pivotal"
- provides summary statistics with known limiting distribution
- leading to approximate pivotal functions, based on normal distribution

Hinkley, 1980

Fraser & Naderi, 2006

- likelihood function depends on data only through sufficient statistics f_{ac} the worker() "likelihood map is sufficient" $(f_{ac}t^{-1}t_{bc})$ Fraser & Naderi, 2006

 - gives exact inference in transformation models
 "likelihood function as pivotal"

 - provides summary statistics with known limiting distribution
- leading to approximate pivotal functions,
 based on normal distribution
 - likelihood function + sample space derivative gives better approximate inference



Hinkley, 1980

• direct use of likelihood function

Likelihood inference

- direct use of likelihood function
- note that only relative values are well-defined

Likelihood inference

- direct use of likelihood function
- note that only relative values are well-defined
- define relative likelihood $RL(\theta) = \frac{L(\theta)}{\sup_{\theta'} L(\theta')} = \frac{L(\theta)}{L(\hat{\theta})}$

Likelihood inference

- direct use of likelihood function
- note that only relative values are well-defined

• define relative likelihood
$$RL(\theta) = \frac{L(\theta)}{\sup_{\theta'} L(\theta')} = \frac{L(\theta)}{L(\hat{\theta})}$$

$$1 \ge RL(\theta) > \frac{1}{3}, \qquad \theta \text{ strongly supported,}$$

$$\frac{1}{3} \ge RL(\theta) > \frac{1}{10}, \qquad \theta \text{ supported,}$$

$$\frac{1}{10} \ge RL(\theta) > \frac{1}{100}, \qquad \theta \text{ weakly supported,}$$

$$\frac{1}{100} \ge RL(\theta) > \frac{1}{1000}, \qquad \theta \text{ poorly supported,}$$

$$\frac{1}{1000} \ge RL(\theta) > 0, \qquad \theta \text{ very poorly supported.}$$

$$\frac{(\hat{\theta})}{(\theta)} \in [1, 3]$$

- combine with a probability density for $\boldsymbol{\theta}$

... likelihood inference

٠

- combine with a probability density for θ

$$\pi(\theta \mid \mathbf{y}) = \frac{f(\mathbf{y}; \theta)\pi(\theta)}{\int f(\mathbf{y}; \theta)\pi(\theta)d\theta}$$

... likelihood inference

٠

- combine with a probability density for $\boldsymbol{\theta}$

$$\pi(\theta \mid \mathbf{y}) = \frac{f(\mathbf{y}; \theta)\pi(\theta)}{\int f(\mathbf{y}; \theta)\pi(\theta) d\theta}$$

• inference for θ via probability statements from $\pi(\theta \mid y)$
... likelihood inference

٠

- combine with a probability density for $\boldsymbol{\theta}$

$$\pi(\theta \mid \mathbf{y}) = \frac{f(\mathbf{y}; \theta)\pi(\theta)}{\int f(\mathbf{y}; \theta)\pi(\theta) d\theta}$$

- inference for θ via probability statements from $\pi(\theta \mid y)$
- e.g., "Probability ($\theta > 0 \mid y$) = 0.23", etc.

.

- combine with a probability density for $\boldsymbol{\theta}$

$$\pi(\theta \mid \mathbf{y}) = \frac{f(\mathbf{y}; \theta)\pi(\theta)}{\int f(\mathbf{y}; \theta)\pi(\theta) d\theta}$$

- inference for θ via probability statements from $\pi(\theta \mid y)$
- e.g., "Probability ($\theta > 0 \mid y$) = 0.23", etc.
- any other use of likelihood function for inference relies on derived quantities and their distribution under the model

۰

• i

• combine with a probability density for θ

$$\frac{\pi(\theta \mid y)}{\int f(y;\theta)\pi(\theta)d\theta}$$
• inference for θ via probability statements from $\pi(\theta \mid y)$
• e.g., "Probability ($\theta > 0 \mid y$) = 0.23", etc.
$$\int_{0}^{\infty} \pi(\theta \mid y) d\theta$$

- any other use of likelihood function for inference relies on derived quantities and their distribution under the model
- the Likelihood Principle states two experiments with proportional likelihood functions lead to the same inference about the same parameter C& H, 1974, p.39 (strong likelihood)

STA 4508 January 12 2022

only Bayeria of "respects" he isk pr.

observed likelihood
$$L(\theta; y) = c(y)f(y; \theta)$$

log-likelihood $\ell(\theta; y) = \log L(\theta; y) = \log f(y; \theta) + a(y)$
score function $U(\theta) = \frac{\partial \ell(\theta; y)}{\partial \theta}$
observed information $j(\theta) = \frac{-\partial^2 \ell(\theta; y)}{\partial \theta \partial \theta}$
First $= \int_{0}^{\infty} expected information i(\theta) = E_{\theta}U(\theta)U(\theta)$ called $(i,(\theta)$ in CH
STA 4508 January 12 2022 $\sum M_{exp} = C(y)f(y; \theta)$
 $D = C(y)f(y; \theta)$

... derived quantities, i.i.d. sample

observed likelihood $L(\theta; y) \propto \prod_{i=1}^{n} f(y_i; \theta)$

 $\ell(\theta; \mathbf{y}) = \sum_{i=1}^{n} \log f(\mathbf{y}; \theta) + a(\mathbf{y})$ log-likelihood $\mathcal{P} = \sum_{i=1}^{n} \frac{\partial \ell(\theta; y_i)}{\partial \theta} = O_p(\sqrt{n})$ $\frac{1}{\sqrt{n}}\mathcal{U}(\Theta) = \mathcal{G}(i)$ bold in prob. score maximum likelihood estimate $\hat{\theta} = \hat{\theta}(\mathbf{v}) = \arg \sup_{\theta} \ell(\theta; \mathbf{v})$ served j(8;y) $j(\hat{\theta}) = -\partial^2 \ell(\hat{\theta}; \mathbf{y}) / \partial \theta \partial \theta^{\mathsf{T}} = O_p(n)$ **Fisher information** expected information $i(\theta) = \mathrm{E}_{\theta} \left(\mathcal{U}(\theta) \mathcal{U}(\theta) \right)^{2} = \mathcal{O}(n) \simeq \operatorname{cov} \left(\mathcal{U}(0) \right)^{2}$ exp Find - ying is nx i, (0) STA 4508 January 12 2022 21

Bartlett identities

$$F_{0}U(\Phi) = 0$$

$$I = \int f(y;\theta)dy$$

$$I = \int f(y;\theta)dy$$

$$I = \int f(y;\theta)dy$$

$$I = \int \frac{\partial}{\partial\theta} \int f(y;\theta)dy = \int \frac{\partial}{\partial\theta} f(y;\theta)dy$$

$$I = \int \frac{\partial}{\partial\theta} \int f(y;\theta)dy = F_{\theta}\{U(\theta;Y)\}$$

$$I = \int \frac{\partial}{\partial\theta} \ell(\theta;y)f(y;\theta)dy = F_{\theta}\{U(\theta;Y)\}$$

$$I = \int \frac{\partial}{\partial\theta} \int \frac{\partial}{\partial\theta} \ell(\theta;y)f(y;\theta)dy$$

$$I = \int \frac{\partial}{\partial\theta} \int \frac{\partial}{\partial\theta} \frac{\partial}{\partial\theta} \ell(\theta;y)f(y;\theta)dy$$

$$I = \int \frac{\partial}{\partial\theta} \int \frac{\partial}{\partial\theta} \frac{\partial}{\partial\theta} \frac{\partial}{\partial\theta} \frac{\partial}{\partial\theta} \ell(\theta;y)f(y;\theta)dy$$

$$I = \int \frac{\partial}{\partial\theta} \int \frac{\partial}{\partial\theta} \frac{$$

You can keep going, as long as the endpoints don't depend on θ , the log-density is differentiable, and the required moments exist.

From the book Tensor Methods by McCullagh:

sample space does not depend on θ . In the univariate case, power notation is often employed in the form $i_{rst} = E\left\{ \left(\frac{\partial l}{\partial \theta}\right)^r \left(\frac{\partial^2 l}{\partial \theta^2}\right)^s \left(\frac{\partial^3 l}{\partial \theta^3}\right)^t; \theta \right\}.$ The moment identities then become $i_{10} = 0$, $i_{01} + i_{20} = 0, \checkmark$ $i_{001} + 3i_{11} + i_{30} = 0$ $i_{0001} + 4i_{101} + 3i_{02} + 6i_{21} + i_{40} = 0.$ Similar identities apply to the cumulants, but we refrain from writ-

ing these down, in order to avoid further conflict of notation.

Or when θ is a vector:

STA 4508 January 12 2022

202

LIKELIHOOD FUNCTIONS

Differentiation with respect to θ and reversing the order of differentiation and integration gives

$$u_r = \kappa_r = \int u_r(\theta; y) f_Y(y; \theta) dy = 0.$$

Further differentiation gives

$$\mu_{[rs]} = \mu_{rs} + \mu_{r,s} = 0$$

• $U(\theta) = \sum_{i=1}^{n} U_i(\theta)$

- $U(\theta) = \sum_{i=1}^{n} U_i(\theta)$
- $E\{U(\theta)\} = O$

- $U(\theta) = \sum_{i=1}^{n} U_i(\theta)$
- $E\{U(\theta)\} = O$
- $var{U(\theta)} = ni_1(\theta)$

•
$$U(\theta) = \sum_{i=1}^{n} U_i(\theta)$$
 and quartities $l(\theta, y) = \sum_{i=1}^{n} U_i(\theta)$ and $Q_{i}(\theta, y) = \sum_{i=1}^{n} U_i(\theta)$

- $E\{U(\theta)\} = O$
- $\operatorname{var}\{U(\theta)\} = \underline{ni_1(\theta)}$ • $U(\theta)/\sqrt{n} \stackrel{d}{\rightarrow} N\{0, i_1(\theta)\}$ + some regularly need $0 < i_1(\theta) < \infty$ $\int_{V_T} \sum U_i \stackrel{d}{\rightarrow} N(0, \operatorname{var}(U_i))$ want CLT to be true true

- $U(\theta) = \sum_{i=1}^{n} U_i(\theta)$
- $E\{U(\theta)\} = O$
- $var{U(\theta)} = ni_1(\theta)$
- $U(\theta)/\sqrt{n} \stackrel{d}{\rightarrow} N\{o, i_1(\theta)\}$

need o < $i_1(\theta) < \infty$

• Note that could have not i.d., or not independent, if we can still prove the limiting normality of the sum. E.g. Lindeberg-Feller type conditions, or weak dependence

• $U(\theta)/\sqrt{n} \stackrel{d}{\rightarrow} N\{o, i_1(\theta)\}$

• $U(\theta)/\sqrt{n} \stackrel{d}{\rightarrow} N\{o, i_1(\theta)\}$

•
$$U(\hat{\theta}) = \mathbf{0} = U(\theta) + (\hat{\theta} - \theta)U'(\theta) + R_n$$

₹ Ø

$$(\hat{\theta} - \theta) = \frac{u(\theta)}{-u'(\theta)} + \frac{s_n}{\sqrt{\theta}}$$

$$\int_{d}^{d} \int_{d}^{d} \int_{d}^{\theta} \int_{0}^{\theta} \frac{1}{\sqrt{\theta}} \int_{0}^{\theta} \frac{1}{$$

=mle =
$$\arg\sup L(0; A)$$

= $\hat{o}(A)$ $\int \frac{1}{2}L(0; A) = 0$
 -30 $\int \frac{1}{0=0} = 0$
 $\int \frac{1}{10}$

- U(θ)/√n
- $U(\hat{\theta}) = 0$
- $(\hat{\theta} \theta) =$

$$\frac{\partial}{\partial n} = 0 = U(\theta) + (\hat{\theta} - \theta)U'(\theta) + R_n$$

$$\frac{\partial}{\partial \theta} = 0 = U(\theta) + (\hat{\theta} - \theta)U'(\theta) + R_n$$

$$\frac{\partial}{\partial \theta} = 0 = \frac{U(\theta)}{U(\theta)} + \frac{\partial}{\partial \theta} + \frac{\partial}{\partial \theta}$$

$$\begin{array}{c} \underbrace{\mathcal{C}} \underbrace{\mathcal{C}} \\ \cdot U(\theta) / \sqrt{n} \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}(\theta)}\} \\ \cdot U(\theta) = 0 = U(\theta) + (\hat{\theta} - \theta)U'(\theta) + R_{n} \\ \cdot U(\hat{\theta}) = 0 = U(\theta) + (\hat{\theta} - \theta)U'(\theta) + R_{n} \\ \cdot (\hat{\theta} - \theta) = \{U(\theta) / i(\theta)\}\{1 + 0_{p}(1)\} \\ \cdot \underbrace{\mathcal{C}} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{q, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{q, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{q, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{q, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{q, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{q, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, \underline{i_{1}^{-1}(\theta)}\}}_{\sqrt{n} i_{1}(\theta)} \\ \cdot \underbrace{\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0,$$

• $\sqrt{n(\hat{\theta} - \theta)} \xrightarrow{d} N\{0, i_1^{-1}(\theta)\}$

$$\sqrt{n(\hat{\theta} - \theta)} \xrightarrow{d} N\{0, i_1^{-1}(\theta)\}$$
 by $de^{j} \xrightarrow{-}$
$$\frac{f_{-}(\hat{\Theta})}{f_{-}(\theta)} \xrightarrow{\leq 3}$$

$$\ell(\theta) = \ell(\hat{\theta}) + (\theta - \hat{\theta})\ell'(\hat{\theta}) + \frac{1}{2}(\theta - \hat{\theta})^2\ell''(\hat{\theta}) + R_n$$

$$\ell(\theta) = \ell(\hat{\Theta}) = (\hat{\Theta} - \hat{\Theta})^2\ell''(\hat{\Theta}) + R_n$$

$$\ell(\theta) = \ell(\hat{\Theta}) = (\hat{\Theta} - \hat{\Theta})^2\ell''(\hat{\Theta}) + \dots + \frac{1}{2}(\theta - \hat{\Theta})^2\ell$$

STA 4508 January 12 2022

• $\sqrt{n(\hat{\theta} - \theta)} \xrightarrow{d} N\{\mathbf{0}, i_1^{-1}(\theta)\}$

- $\ell(\theta) = \ell(\hat{\theta}) + (\theta \hat{\theta})\ell'(\hat{\theta}) + \frac{1}{2}(\theta \hat{\theta})^2\ell''(\hat{\theta}) + R_n$
- $2\{\ell(\hat{\theta}) \ell(\theta)\} = (\hat{\theta} \theta)^2 i(\theta)\{1 + o_p(1)\}$

• $\sqrt{n(\hat{\theta} - \theta)} \xrightarrow{d} N\{\mathbf{0}, i_1^{-1}(\theta)\}$

- $\ell(\theta) = \ell(\hat{\theta}) + (\theta \hat{\theta})\ell'(\hat{\theta}) + \frac{1}{2}(\theta \hat{\theta})^2\ell''(\hat{\theta}) + R_n$
- $2\{\ell(\hat{\theta}) \ell(\theta)\} = (\hat{\theta} \theta)^2 i(\theta)\{1 + O_p(1)\}$
- $\mathbf{2}\{\ell(\hat{\theta}) \ell(\theta)\} \xrightarrow{d} \chi^2_d$

• $\hat{\theta} \sim N_d\{\theta, j^{-1}(\hat{\theta})\}$

$$j(\hat{\theta}) = -\ell''(\hat{\theta}; y)$$



STA 4508 January 12 2022

27

- $\hat{\theta} \sim N_d \{ \theta, j^{-1}(\hat{\theta}) \}$
- " θ is estimated to be 21.5 (95% Cl 19.5 23.5)"

$$\begin{split} j(\hat{\theta}) &= -\ell^{\prime\prime}(\hat{\theta}; \mathbf{y}) \\ \hat{\theta} \pm \mathbf{2}\hat{\sigma} \end{split}$$



log-likelihood function

- $\hat{\theta} \sim N_d \{ \theta, j^{-1}(\hat{\theta}) \}$
- " θ is estimated to be 21.5 (95% Cl 19.5 23.5)"

$$\begin{split} j(\hat{\theta}) &= -\ell^{\prime\prime}(\hat{\theta}; \mathbf{y}) \\ \hat{\theta} \pm \mathbf{2}\hat{\sigma} \end{split}$$

• $W(\theta) = 2\{\ell(\hat{\theta}) - \ell(\theta)\} \sim \chi^2_d$



STA 4508 January 12 2022

$$\mathcal{F}(\hat{\phi} - \phi) \xrightarrow{d} \mathcal{N}(\phi, i, \mathcal{V}(\phi))$$

 $j(\hat{\theta}) = -\ell''(\hat{\theta}; \mathbf{y})$ $\hat{\theta} + 2\hat{\sigma}$

• $\hat{\theta} \sim N_{\varphi} \{\theta, j^{-1}(\hat{\theta})\}$ • " θ is estimated to be 21.5 (95% Cl 19.5 – 23.5)"

- $\mathsf{W}(heta) = 2\{\ell(\hat{ heta}) \ell(heta)\} \sim \chi^2_d$
- "likelihood based CI for θ with confidence level 95% is (18.6, 23.0)" \longrightarrow (θ)



p-value functions of θ

۰



$$\begin{aligned} -\ell''(\hat{b}) \text{ replaces } & F_0(\ell''(o)) \\ r_u(\theta) &= U(\theta)j^{-1/2}(\hat{\theta}) \sim N(0,1) \\ r_e(\theta) &= (\hat{\theta} - \theta)j^{1/2}(\hat{\theta}), \quad \neg \mathcal{N}(0,1) \\ r(\theta) &= \operatorname{sign}(\hat{\theta} - \theta)[2\{\ell(\hat{\theta}) - \ell(\theta)\}]^{1/2} \\ \mathcal{M}_{ben} \\ \mathcal{T}_n(\hat{b} - \theta) &= \mathcal{N}(0, i_1^{-1}(\theta)) \\ \mathcal{T}_n(0) &= \mathcal{N}(0, i_1^{-1}(\theta)) \\ \mathcal{T}_n(0) &= \operatorname{N}(0, i_1^{-1}(\theta)) \\ \mathcal{T}_n(1) &= \operatorname{N}(0, i_1^{-1}(\theta)) \\ \mathcal$$

p-value functions of θ

• approximate pivotal quantities

-

$$\Pr\{r_u(\theta) \le r_u^{\mathsf{o}}(\theta)\} \doteq \Phi\{r_u^{\mathsf{o}}(\theta)\}$$

under sampling from the model $f(y; \theta) = f(y_1, \dots, y_n; \theta)$

$$r_{u}(0) = r(0; u) \int as = r.v. r \sim N(0, 1) \text{ show } 0 \text{ is true value}$$

$$f$$

$$score f = pivotal q. f = of 0 \quad dist = si \quad fuen$$

$$r \sim N \qquad f = of Y$$

.

$$\begin{aligned} r_u(\theta) &= U(\theta)j^{-1/2}(\hat{\theta}) \sim N(0,1) \\ r_e(\theta) &= (\hat{\theta} - \theta)j^{1/2}(\hat{\theta}), \\ r(\theta) &= \operatorname{sign}(\hat{\theta} - \theta)[2\{\ell(\hat{\theta}) - \ell(\theta)\}]^{1/2} \end{aligned}$$

• approximate pivotal quantities

Pr{
$$r_u(\theta) \le r_u^o(\theta)$$
} $\doteq \Phi{r_u^o(\theta)}$
under sampling from the model $f(y; \theta) = f(y_1, \dots, y_n; \theta)$
 p -value function (of θ , for fixed data)
 $p_u(\theta) = \Phi{r_u^o(\theta)}$ for $fixed rest fixed rest$

STA 4508 January 12 2022

28

$$r_{u}(\theta) = U(\theta)j^{-1/2}(\hat{\theta}) \sim N(0,1) \quad \{\theta: | f_{u}(\theta)| \leq (.96)\}$$

$$r_{e}(\theta) = U(\theta)j^{-1/2}(\hat{\theta}) \sim N(0,1) \quad \{g: v \in 95\}, \quad v \in 95$$

$$p_u(\theta) = \Phi\{r_u^{\mathsf{o}}(\theta)\}$$

• similarly $p_e(\theta) = \Phi\{r_e(\theta)\}, \quad p_r(\theta) = \Phi\{r(\theta)\}$ are also *p*-value functions for θ based de-limiting disting

STA 4508 January 12 2022



Figure 2.2: Approximate pivots and P-values based on an exponential sample of size n = 1. Left: likelihood root $r(\theta)$ (solid), score pivot $s(\theta)$ (dots), Wald pivot $t(\theta)$ (dashes), modified likelihood root $r^*(\theta)$ (heavy), and exact pivot $\theta \sum y_j$ (dot-dash). The modified likelihood root is indistinguishable from the exact pivot. The horizontal lines are at $0, \pm 1.96$. Right: corresponding significance functions, with horizontal lines at 0.025 and 0.975.



BDR, Ch.3.2, Cauchy, distribution functions (y) at $\theta = 0$, n = 1

Example: Exponential

•
$$f(y_i; \theta) = \theta e^{-y_i \theta}, \quad i = 1, \dots, n$$

- $\ell(\theta) =$
- $\ell'(\theta) =$
- $\ell''(\theta) =$
- $r_u(\theta) =$
- $r_e(\theta) =$
- $r(\theta) =$

expand $\log(\theta \bar{y})$ around 1 to get asymptotic equivalence to r_e , r_u

Example: Exponential

•
$$f(y_i; \theta) = \theta e^{-y_i \theta}, \quad i = 1, \dots, n$$

•
$$\ell(\theta) = n \log \theta - n \theta \bar{y}$$

•
$$\ell'(\theta) = \frac{n}{\theta} - n\bar{y}$$
 $\hat{\theta} = \bar{y}^{-1}$

•
$$\ell''(\theta) = -\frac{n}{\theta^2}$$

•
$$r_u(\theta) = \frac{1}{\sqrt{n}}\ell'(\theta)j^{-1/2}(\hat{\theta}) = \sqrt{n}(\frac{1}{\theta\bar{y}}-1)$$

•
$$r_e(\theta) = (\hat{\theta} - \theta)j^{1/2}(\hat{\theta}) = \sqrt{n(1 - \bar{y}\theta)}$$

•
$$r(\theta) = \sqrt{(2n)} \{\theta \overline{y} - 1 - \log(\theta \overline{y})\}^{1/2}$$

expand $\log(\theta \bar{y})$ around 1 to get asymptotic equivalence to r_e , r_u

Example: Poisson

- $f(y_i; \theta) = \theta^{y_i} e^{-\theta} / y_i!$
- $\ell(\theta) =$
- $\ell'(\theta) =$
- $\ell''(\theta) =$
- $r_e(\theta) = (s n\theta)/\sqrt{s}$
- $Pr(S \le s) \ne 1 Pr(S \ge s)$
- upper and lower *p*-value functions: Pr(S < s), $Pr(S \le s)$
- mid *p*-value function: Pr(S < sr) + 0.5Pr(S = s)

STA 4508 January 12 2022



Figure 3.2: Cumulative distribution function for Poisson distribution with parameter 6.7 (solid), with approximations $\Phi\{r^*(y)\}$ (dashes) and $\Phi\{r^*(y+1/2)\}$ (dots). The vertical lines are at 0.5, 1.5, 2.5, ...



- for inference re θ , given y, plot $p(\theta)$ vs θ
- for *p*-value for $H_0: \theta = \theta_0$, compute $p(\theta_0)$
- for checking whether, e.g. $\Phi\{r_e(\theta)\}$ is a good approximation,
 - compare $p(\theta) = \Phi\{r_e(\theta)\}$ to $p_{exact}(\theta)$, as a function of θ , fixed y
 - or compare $p(\theta_0)$ to $p_{\text{exact}}(\theta_0)$ as a function of y
- if $p_{\text{exact}}(\theta)$ not available, simulate
- if $\boldsymbol{\theta}$ is a vector, choose one component at a time

•
$$\theta = (\psi, \lambda) = (\psi_1, \dots, \psi_q, \lambda_1, \dots, \lambda_{d-q})$$
•
$$\theta = (\psi, \lambda) = (\psi_1, \dots, \psi_q, \lambda_1, \dots, \lambda_{d-q})$$

•
$$U(\theta) = \begin{pmatrix} U_{\psi}(\theta) \\ U_{\lambda}(\theta) \end{pmatrix}, \qquad U_{\lambda}(\psi, \hat{\lambda}_{\psi}) = \mathbf{0}$$

•
$$\theta = (\psi, \lambda) = (\psi_1, \dots, \psi_q, \lambda_1, \dots, \lambda_{d-q})$$

•
$$U(\theta) = \begin{pmatrix} U_{\psi}(\theta) \\ U_{\lambda}(\theta) \end{pmatrix}, \quad U_{\lambda}(\psi, \hat{\lambda}_{\psi}) = \mathbf{0}$$

• $i(\theta) = \begin{pmatrix} i_{\psi\psi} & i_{\psi\lambda} \\ i_{\lambda\psi} & i_{\lambda\lambda} \end{pmatrix} \quad j(\theta) = \begin{pmatrix} j_{\psi\psi} & j_{\psi\lambda} \\ j_{\lambda\psi} & j_{\lambda\lambda} \end{pmatrix}$

•
$$\theta = (\psi, \lambda) = (\psi_1, \dots, \psi_q, \lambda_1, \dots, \lambda_{d-q})$$

•
$$U(\theta) = \begin{pmatrix} U_{\psi}(\theta) \\ U_{\lambda}(\theta) \end{pmatrix}, \quad U_{\lambda}(\psi, \hat{\lambda}_{\psi}) = 0$$

• $i(\theta) = \begin{pmatrix} i_{\psi\psi} & i_{\psi\lambda} \\ i_{\lambda\psi} & i_{\lambda\lambda} \end{pmatrix} \quad j(\theta) = \begin{pmatrix} j_{\psi\psi} & j_{\psi\lambda} \\ j_{\lambda\psi} & j_{\lambda\lambda} \end{pmatrix}$
• $i^{-1}(\theta) = \begin{pmatrix} i^{\psi\psi} & i^{\psi\lambda} \\ i^{\lambda\psi} & i^{\lambda\lambda} \end{pmatrix} \quad j^{-1}(\theta) = \begin{pmatrix} j^{\psi\psi} & j^{\psi\lambda} \\ j^{\lambda\psi} & j^{\lambda\lambda} \end{pmatrix}$

٠

•
$$\theta = (\psi, \lambda) = (\psi_1, \dots, \psi_q, \lambda_1, \dots, \lambda_{d-q})$$

•
$$U(\theta) = \begin{pmatrix} U_{\psi}(\theta) \\ U_{\lambda}(\theta) \end{pmatrix}, \quad U_{\lambda}(\psi, \hat{\lambda}_{\psi}) = 0$$

• $i(\theta) = \begin{pmatrix} i_{\psi\psi} & i_{\psi\lambda} \\ i_{\lambda\psi} & i_{\lambda\lambda} \end{pmatrix} \quad j(\theta) = \begin{pmatrix} j_{\psi\psi} & j_{\psi\lambda} \\ j_{\lambda\psi} & j_{\lambda\lambda} \end{pmatrix}$
• $i^{-1}(\theta) = \begin{pmatrix} i^{\psi\psi} & i^{\psi\lambda} \\ i^{\lambda\psi} & i^{\lambda\lambda} \end{pmatrix} \quad j^{-1}(\theta) = \begin{pmatrix} j^{\psi\psi} & j^{\psi\lambda} \\ j^{\lambda\psi} & j^{\lambda\lambda} \end{pmatrix}$

٠

•
$$i^{\psi\psi}(\theta) = \{i_{\psi\psi}(\theta) - i_{\psi\lambda}(\theta)i_{\lambda\lambda}^{-1}(\theta)i_{\lambda\psi}(\theta)\}^{-1},$$

•
$$\theta = (\psi, \lambda) = (\psi_1, \dots, \psi_q, \lambda_1, \dots, \lambda_{d-q})$$

•
$$U(\theta) = \begin{pmatrix} U_{\psi}(\theta) \\ U_{\lambda}(\theta) \end{pmatrix}, \quad U_{\lambda}(\psi, \hat{\lambda}_{\psi}) = 0$$

• $i(\theta) = \begin{pmatrix} i_{\psi\psi} & i_{\psi\lambda} \\ i_{\lambda\psi} & i_{\lambda\lambda} \end{pmatrix} \quad j(\theta) = \begin{pmatrix} j_{\psi\psi} & j_{\psi\lambda} \\ j_{\lambda\psi} & j_{\lambda\lambda} \end{pmatrix}$
• $i^{-1}(\theta) = \begin{pmatrix} j^{\psi\psi} & j^{\psi\lambda} \\ i^{\lambda\psi} & j^{\lambda\lambda} \end{pmatrix} \quad j^{-1}(\theta) = \begin{pmatrix} j^{\psi\psi} & j^{\psi\lambda} \\ j^{\lambda\psi} & j^{\lambda\lambda} \end{pmatrix}$

٠

•
$$\mathbf{i}^{\psi\psi}(\theta) = \{\mathbf{i}_{\psi\psi}(\theta) - \mathbf{i}_{\psi\lambda}(\theta)\mathbf{i}_{\lambda\lambda}^{-1}(\theta)\mathbf{i}_{\lambda\psi}(\theta)\}^{-1},$$

•
$$\ell_{\mathrm{P}}(\psi) = \ell(\psi, \hat{\lambda}_{\psi}), \qquad j_{\mathrm{P}}(\psi) = -\ell_{\mathrm{P}}''(\psi)$$

Inference from limiting distributions, nuisance parameters

$$W_{u}(\psi) = U_{\psi}(\psi, \hat{\lambda}_{\psi})^{\mathsf{T}} \{ i^{\psi\psi}(\psi, \hat{\lambda}_{\psi}) \} U_{\psi}(\psi, \hat{\lambda}_{\psi}) \quad \sim \quad \chi_{q}^{2}$$
$$W_{e}(\psi) = (\hat{\psi} - \psi) \{ i^{\psi\psi}(\hat{\psi}, \hat{\lambda}) \}^{-1} (\hat{\psi} - \psi) \quad \sim \quad \chi_{q}^{2}$$

$$\mathsf{W}(\psi) = \mathbf{2}\{\ell(\hat{\psi}, \hat{\lambda}) - \ell(\psi, \hat{\lambda}_{\psi})\} = \mathbf{2}\{\ell_{\mathrm{P}}(\hat{\psi}) - \ell_{\mathrm{P}}(\psi)\} \quad \sim \quad \chi_{q}^{2}$$

Approximate Pivots, q = 1

$$\begin{aligned} r_{u}(\psi) &= \ell_{\mathsf{P}}'(\psi) j_{\mathsf{P}}(\hat{\psi})^{-1/2} \sim \mathsf{N}(\mathsf{O},\mathsf{1}), \\ r_{e}(\psi) &= (\hat{\psi} - \psi) j_{\mathsf{P}}(\hat{\psi})^{1/2} \sim \mathsf{N}(\mathsf{O},\mathsf{1}), \\ r(\psi) &= \operatorname{sign}(\hat{\psi} - \psi) [2\{\ell_{\mathsf{P}}(\hat{\psi}) - \ell_{\mathsf{P}}(\psi)\}]^{1/2} \sim \mathsf{N}(\mathsf{O},\mathsf{1}), \end{aligned}$$



Figure 2.3: Inference for shape parameter ψ of gamma sample of size n = 5. Left: profile log likelihood ℓ_p (solid) and the log likelihood from the conditional density of u given v (heavy). Right: likelihood root $r(\psi)$ (solid), Wald pivot $t(\psi)$ (dashes), modified likelihood root $r^*(\psi)$ (heavy), and exact pivot overlying $r^*(\psi)$. The horizontal lines are at $0, \pm 1.96$.