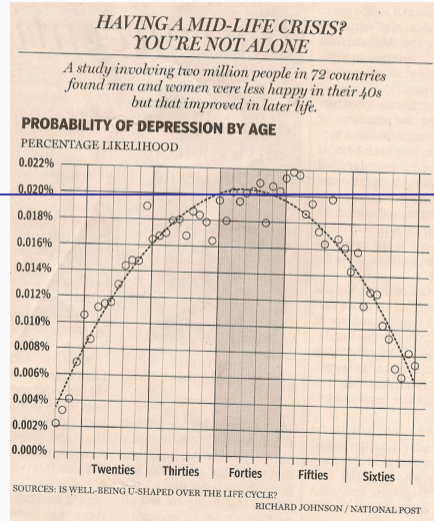


Topics in Likelihood Inference

STA2212H S LEC9101

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

1. Inference based on the likelihood function: derived quantities, limiting distributions, approximations to posterior distributions;
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

Various 'types' of likelihood

1. likelihood, marginal and conditional likelihood, profile likelihood, adjusted profile
2. semi-parametric likelihood, partial likelihood
3. quasi-likelihood, composite likelihood misspecified models
4. empirical likelihood, penalized likelihood
5. simulated likelihood, indirect inference
6. bootstrap likelihood, h -likelihood, weighted likelihood, pseudo-likelihood, local likelihood, sieve likelihood

Why so many?

- **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

The likelihood function

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

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

- typically, $\mathbf{y} = (y_1, \dots, y_n)$ $x_1, \dots, x_n \quad i = 1, \dots, n$
- $f(\mathbf{y}; \theta)$ or $f(\mathbf{y} | \mathbf{x}; \theta)$ is **joint density**
- under independence $L(\theta; \mathbf{y}) \propto \prod f(y_i | x_i; \theta)$
- **log-likelihood** $\ell(\theta; \mathbf{y}) = \log L(\theta; \mathbf{y}) = \sum \log f(y_i | x_i; \theta)$
- θ could have dimension $p > n$ (e.g. genetics), or $p \uparrow n$, or
- θ could have infinite dimension e.g.
- **regular model** $p < n$ and p fixed as n increases

Examples

- $y_i \sim N(\mu, \sigma^2)$:

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

- $E(y_i) = \mathbf{x}_i^T \boldsymbol{\beta}$:

$$L(\theta; \mathbf{y}) = \prod_{i=1}^n \sigma^{-n} \exp\left\{-\frac{1}{2\sigma^2}(y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2\right\}$$

- $E(y_i) = m(\mathbf{x}_i)$, $m(\mathbf{x}) = \sum_{j=1}^J \phi_j \mathbf{B}_j(\mathbf{x})$:

$$L(\theta; \mathbf{y}) = \prod_{i=1}^n \sigma^{-n} \exp\left\{-\frac{1}{2\sigma^2}(y_i - \sum_{j=1}^J \phi_j \mathbf{B}_j(\mathbf{x}_i))^2\right\}$$

- $y_i = \mu + \rho(y_{i-1} - \mu) + \epsilon_i$, $\epsilon_i \sim N(0, \sigma^2)$:

$$L(\theta; y) = \prod_{i=1}^n f(y_i | y_{i-1}; \theta) f_0(y_0; \theta)$$

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

$$L(\theta; y) = \prod_{i=1}^n \theta^{-n}, \quad 0 < y_{(1)} < \dots < y_{(n)} < \theta$$

... examples

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

$$\ell\{\lambda(\cdot); \mathbf{y}\} = \sum_{i=1}^n \log\{\lambda(y_i)\} - \int_0^\tau \lambda(u) du, \quad 0 < y_1 < \dots < y_n < \tau$$

Davison, §6.5

- multinomial: $\mathbf{y}_i = (y_{i1}, \dots, y_{ik})$, $y_{ic} = 1, y_{ic'} = 0, c' \neq c$

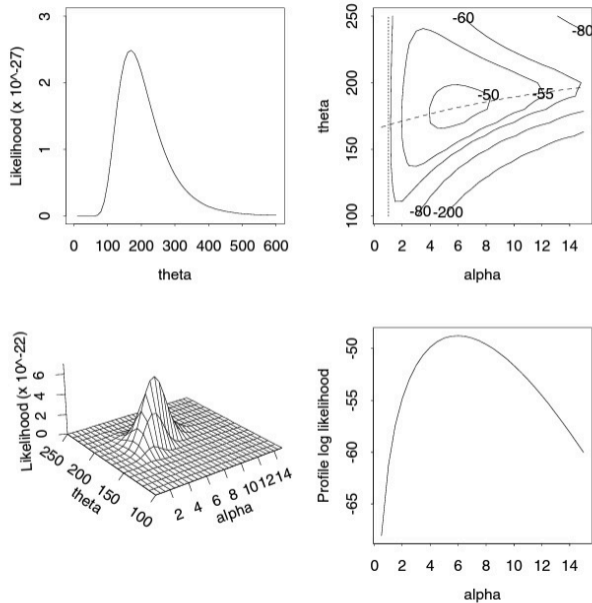
$$\ell(\theta; \mathbf{y}) = \sum_{i=1}^n \sum_{c=1}^k y_{ic} \log(p_{ic})$$

negative cross-entropy

$p_{ic} = p(x_{ic}; \theta)$, as above

Hastie et al., Ch. 7

Figure 4.1 Likelihoods for the spring failure data at stress 950 N/mm^2 . The upper left panel is the likelihood for the exponential model, and below it is a perspective plot of the likelihood for the Weibull model. The upper right panel shows contours of the log likelihood for the Weibull model; the exponential likelihood is obtained by setting $\alpha = 1$, that is, slicing L along the vertical dotted line. The lower right panel shows the profile log likelihood for α , which corresponds to the log likelihood values along the dashed line in the panel above, plotted against α .



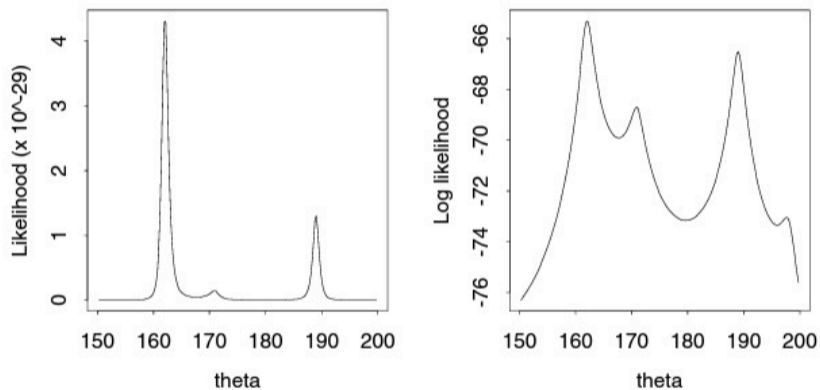


Figure 4.2 Cauchy likelihood and log likelihood for the spring failure data at stress 950N/mm^2 .

Data: times of failure of a spring under stress
 225, 171, 198, 189, 189, 135, 162, 135, 117, 162

Complicated likelihoods

- example: clustered binary data Renard et al. (2004)
- latent variable: $z_{ir} = \mathbf{x}_{ir}^T \boldsymbol{\beta} + \mathbf{b}_i + \epsilon_{ir}$, $\mathbf{b}_i \sim N(\mathbf{0}, \sigma_b^2)$, $\epsilon_{ir} \sim N(0, 1)$
- $r = 1, \dots, n_i$: observations in a cluster/family/school... $i = 1, \dots, n$ clusters
- random effect \mathbf{b}_i introduces correlation between observations in a cluster
- observations: $y_{ir} = 1$ if $z_{ir} > 0$, else 0

- $Pr(y_{ir} = 1 \mid \mathbf{b}_i) = \Phi(\mathbf{x}_{ir}^T \boldsymbol{\beta} + \mathbf{b}_i) = p_i$

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

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

- more general: $z_{ir} = \mathbf{x}_{ir}^T \boldsymbol{\beta} + \mathbf{w}'_{ir} \mathbf{b}_i + \epsilon_{ir}$

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

- generalized linear geostatistical models

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

Diggle & Ribeiro, 2007

- random intercept u is a realization of a stationary GRF, expected value 0, covariance

Gaussian random field

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

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

$$L(\theta; y) = \int_{\mathbb{R}^n} \prod_{i=1}^n f(y_i \mid u_i; \theta) \underbrace{f(u_i; \theta)}_{N_d(0, \Sigma)} du_1 \dots du_n$$

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

- Ising model:

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

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

Davison §6.2

Ravikumar et al. (2010).

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

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.

Section	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
10. The Efficiency of the Method of Moments in Fitting Pearsonian Curves	342



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.

Why likelihood?

- 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

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

$L(\hat{\theta})/L(\theta) \in (3, 10)$ implausible;

$L(\hat{\theta})/L(\theta) \in (10, \infty)$ very implausible

Statistical Evidence: A likelihood paradigm

- converts a 'prior' probability $\pi(\theta)$ to a posterior $\pi(\theta | y)$ via Bayes' formula
- provides a conventional set of summary quantities for inference based on properties of the postulated model

... why likelihood?

- 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
- 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

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

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

$$\begin{aligned} 1 &\geq RL(\theta) > \frac{1}{3}, && \theta \text{ strongly supported,} \\ \frac{1}{3} &\geq RL(\theta) > \frac{1}{10}, && \theta \text{ supported,} \\ \frac{1}{10} &\geq RL(\theta) > \frac{1}{100}, && \theta \text{ weakly supported,} \\ \frac{1}{100} &\geq RL(\theta) > \frac{1}{1000}, && \theta \text{ poorly supported,} \\ \frac{1}{1000} &\geq RL(\theta) > 0, && \theta \text{ very poorly supported.} \end{aligned}$$

SM (4.11)

- combine with a probability density for θ

-

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

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

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

log-likelihood $\ell(\theta; \mathbf{y}) = \log L(\theta; \mathbf{y}) = \log f(\mathbf{y}; \theta) + a(\mathbf{y})$

score $U(\theta) = \partial \ell(\theta; \mathbf{y}) / \partial \theta$

observed information $j(\theta) = -\partial^2 \ell(\theta; \mathbf{y}) / \partial \theta \partial \theta^T$

expected information $i(\theta) = \mathbb{E}_\theta U(\theta)U(\theta)^T$ called $i_1(\theta)$ in CH

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

observed likelihood

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

log-likelihood

$$\ell(\theta; \mathbf{y}) = \sum_{i=1}^n \log f(y_i; \theta) + a(\mathbf{y})$$

score

$$U(\theta) = \partial \ell(\theta; \mathbf{y}) / \partial \theta = O_p(\sqrt{n})$$

maximum likelihood estimate

$$\hat{\theta} = \hat{\theta}(\mathbf{y}) = \arg \sup_{\theta} \ell(\theta; \mathbf{y})$$

Fisher information

$$j(\hat{\theta}) = -\partial^2 \ell(\hat{\theta}; \mathbf{y}) / \partial \theta \partial \theta^T = O_p(n)$$

expected information

$$i(\theta) = E_{\theta} U(\theta) U(\theta)^T = O(n)$$

$$1 = \int f(\mathbf{y}; \theta) d\mathbf{y}$$

$$\begin{aligned} 0 &= \frac{\partial}{\partial \theta} \int f(\mathbf{y}; \theta) d\mathbf{y} = \int \frac{\partial}{\partial \theta} f(\mathbf{y}; \theta) d\mathbf{y} \\ &= \int \frac{\partial}{\partial \theta} \ell(\theta; \mathbf{y}) f(\mathbf{y}; \theta) d\mathbf{y} = \mathbb{E}_\theta \{ \mathbf{U}(\theta; \mathbf{Y}) \} \end{aligned}$$

$$\begin{aligned} 0 &= \frac{\partial}{\partial \theta} \int \frac{\partial}{\partial \theta} \ell(\theta; \mathbf{y}) f(\mathbf{y}; \theta) d\mathbf{y} \\ &= \int \left[\frac{\partial^2}{\partial \theta \partial \theta^\top} \ell(\theta; \mathbf{y}) + \left\{ \frac{\partial}{\partial \theta} \ell(\theta; \mathbf{y}) \right\} \left\{ \frac{\partial}{\partial \theta} \ell(\theta; \mathbf{y}) \right\}^\top \right] f(\mathbf{y}; \theta) d\mathbf{y} \end{aligned}$$

$$\Rightarrow \mathbb{E}_\theta \{ \mathbf{U}(\theta) \mathbf{U}^\top(\theta) \} = \mathbb{E}_\theta \left\{ -\frac{\partial^2}{\partial \theta \partial \theta^\top} \ell(\theta; \mathbf{y}) \right\} = \mathbf{i}(\theta) = \mathbb{E}_\theta \{ \mathbf{j}(\theta) \}$$

... Bartlett identities

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 $i_{10} = 0$,

$$i_{01} + i_{20} = 0,$$

$$i_{001} + 3i_{11} + i_{30} = 0,$$

$$i_{0001} + 4i_{101} + 3i_{02} + 6i_{21} + i_{40} = 0.$$

Similar identities apply to the cumulants, but we refrain from writing these down, in order to avoid further conflict of notation.

Or when θ is a vector:

202

LIKELIHOOD FUNCTIONS

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

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

Further differentiation gives

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

$$\mu_{[rst]} = \mu_{rst} + \mu_{r,[st]} + \mu_{[rs]t} = 0$$

Limiting distributions

- $U(\theta) = \sum_{i=1}^n U_i(\theta)$
- $E\{U(\theta)\} = 0$
- $\text{var}\{U(\theta)\} = ni_1(\theta)$
- $U(\theta)/\sqrt{n} \xrightarrow{d} N\{0, i_1(\theta)\}$ need $0 < i_1(\theta) < \infty$
- Note that could have not i.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

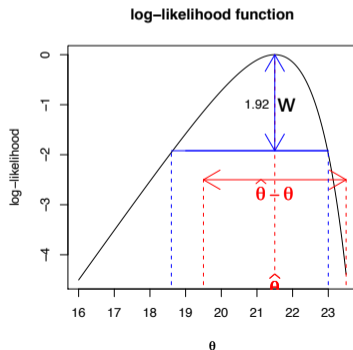
- $U(\theta)/\sqrt{n} \xrightarrow{d} N\{\mathbf{0}, i_1(\theta)\}$
- $U(\hat{\theta}) = \mathbf{0} = U(\theta) + (\hat{\theta} - \theta)U'(\theta) + R_n$
- $(\hat{\theta} - \theta) = \{U(\theta)/i(\theta)\}\{1 + o_p(1)\}$
- $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N\{\mathbf{0}, i_1^{-1}(\theta)\}$

- $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N\{\mathbf{0}, i_1^{-1}(\theta)\}$
- $\ell(\theta) = \ell(\hat{\theta}) + (\theta - \hat{\theta})\ell'(\hat{\theta}) + \frac{1}{2}(\theta - \hat{\theta})^2\ell''(\hat{\theta}) + R_n$
- $2\{\ell(\hat{\theta}) - \ell(\theta)\} = (\hat{\theta} - \theta)^2 i(\theta) \{1 + o_p(1)\}$
- $2\{\ell(\hat{\theta}) - \ell(\theta)\} \xrightarrow{d} \chi_d^2$

Inference from limiting distributions

- $\hat{\theta} \sim N_d\{\theta, j^{-1}(\hat{\theta})\}$
- “ θ is estimated to be 21.5 (95% CI 19.5 – 23.5)”
- $w(\theta) = 2\{\ell(\hat{\theta}) - \ell(\theta)\} \sim \chi_d^2$
- “likelihood based CI for θ with confidence level 95% is (18.6, 23.0)”

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



-

$$r_u(\theta) = U(\theta)j^{-1/2}(\hat{\theta}) \sim N(\mathbf{0}, \mathbf{1})$$

$$r_e(\theta) = (\hat{\theta} - \theta)j^{1/2}(\hat{\theta}),$$

$$r(\theta) = \text{sign}(\hat{\theta} - \theta)[2\{\ell(\hat{\theta}) - \ell(\theta)\}]^{1/2}$$

- approximate pivotal quantities

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

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

- p -value function (of θ , for fixed data)

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

- similarly $p_e(\theta) = \Phi\{r_e(\theta)\}$, $p_r(\theta) = \Phi\{r(\theta)\}$ are also p -value functions for θ

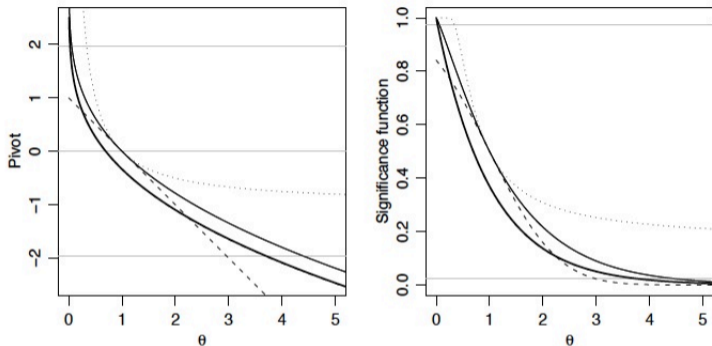
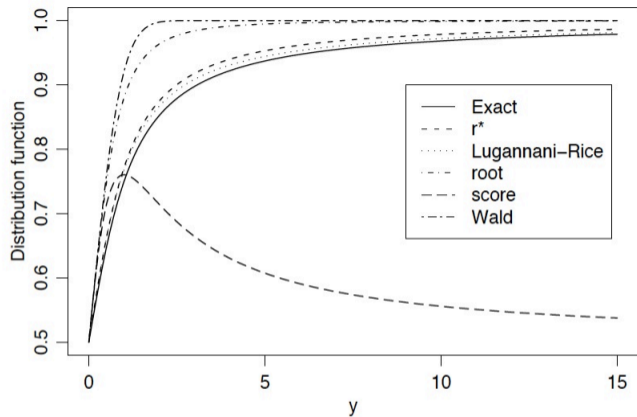


Figure 2.2: Approximate pivots and P-values based on an exponential sample of size $n = 1$. Left: likelihood root $r(\theta)$ (solid), score pivot $s(\theta)$ (dots), Wald pivot $t(\theta)$ (dashes), modified likelihood root $r^*(\theta)$ (heavy), and exact pivot $\theta \sum y_j$ (dot-dash). The modified 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

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

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

- $\ell'(\theta) = \frac{n}{\theta} - n\bar{y}$

$$\hat{\theta} = \bar{y}^{-1}$$

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

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

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

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

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

Example: Poisson

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

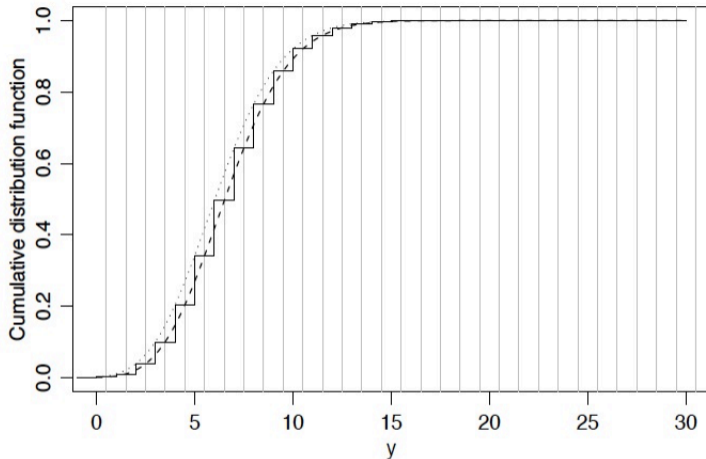


Figure 3.2: Cumulative distribution function for Poisson distribution with parameter 6.7 (solid), with approximations $\Phi\{r^*(y)\}$ (dashes) and $\Phi\{r^*(y + 1/2)\}$ (dots). The vertical lines are at 0.5, 1.5, 2.5, ...

- 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 $p_{\text{exact}}(\theta)$, 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 θ is a vector, choose one component at a time

Nuisance parameters

- $\theta = (\psi, \lambda) = (\psi_1, \dots, \psi_q, \lambda_1, \dots, \lambda_{d-q})$
- $U(\theta) = \begin{pmatrix} U_\psi(\theta) \\ U_\lambda(\theta) \end{pmatrix}, \quad U_\lambda(\psi, \hat{\lambda}_\psi) = \mathbf{0}$
- $i(\theta) = \begin{pmatrix} i_{\psi\psi} & i_{\psi\lambda} \\ i_{\lambda\psi} & i_{\lambda\lambda} \end{pmatrix} \quad j(\theta) = \begin{pmatrix} j_{\psi\psi} & j_{\psi\lambda} \\ j_{\lambda\psi} & j_{\lambda\lambda} \end{pmatrix}$
- $i^{-1}(\theta) = \begin{pmatrix} i^{\psi\psi} & i^{\psi\lambda} \\ i^{\lambda\psi} & i^{\lambda\lambda} \end{pmatrix} \quad j^{-1}(\theta) = \begin{pmatrix} j^{\psi\psi} & j^{\psi\lambda} \\ j^{\lambda\psi} & j^{\lambda\lambda} \end{pmatrix}.$
- $i^{\psi\psi}(\theta) = \{i_{\psi\psi}(\theta) - i_{\psi\lambda}(\theta)i_{\lambda\lambda}^{-1}(\theta)i_{\lambda\psi}(\theta)\}^{-1},$
- $\ell_P(\psi) = \ell(\psi, \hat{\lambda}_\psi), \quad j_P(\psi) = -\ell''_P(\psi)$

$$\begin{aligned}w_u(\psi) &= U_\psi(\psi, \hat{\lambda}_\psi)^T \{i^{\psi\psi}(\psi, \hat{\lambda}_\psi)\} U_\psi(\psi, \hat{\lambda}_\psi) \quad \sim \quad \chi_q^2 \\w_e(\psi) &= (\hat{\psi} - \psi) \{i^{\psi\psi}(\hat{\psi}, \hat{\lambda})\}^{-1} (\hat{\psi} - \psi) \quad \sim \quad \chi_q^2 \\w(\psi) &= 2\{\ell(\hat{\psi}, \hat{\lambda}) - \ell(\psi, \hat{\lambda}_\psi)\} = 2\{\ell_P(\hat{\psi}) - \ell_P(\psi)\} \quad \sim \quad \chi_q^2;\end{aligned}$$

Approximate Pivots, $q = 1$

$$\begin{aligned}r_u(\psi) &= \ell'_P(\psi) j_P(\hat{\psi})^{-1/2} \sim N(\mathbf{0}, \mathbf{1}), \\r_e(\psi) &= (\hat{\psi} - \psi) j_P(\hat{\psi})^{1/2} \sim N(\mathbf{0}, \mathbf{1}), \\r(\psi) &= \text{sign}(\hat{\psi} - \psi) [2\{\ell_P(\hat{\psi}) - \ell_P(\psi)\}]^{1/2} \sim N(\mathbf{0}, \mathbf{1})\end{aligned}$$

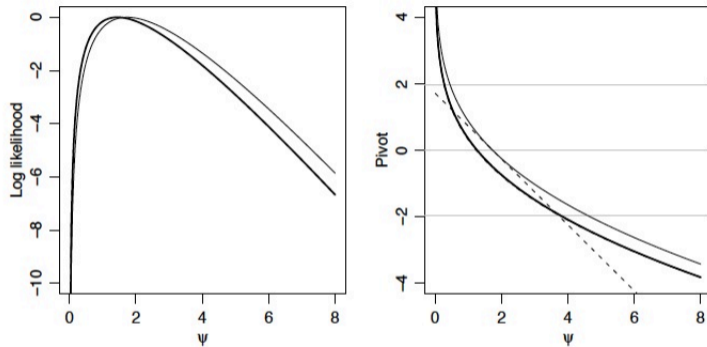


Figure 2.3: Inference for shape parameter ψ of gamma sample of size $n = 5$. Left: profile log likelihood ℓ_p (solid) and the log likelihood from the conditional density of u given v (heavy). Right: likelihood root $r(\psi)$ (solid), Wald pivot $t(\psi)$ (dashes), modified likelihood root $r^*(\psi)$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