

Topics in Likelihood Inference

STA4508H

Nancy Reid
University of Toronto

January 19, 2022



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

Off

Mon
Tue

7-8 pm
5-6 pm



Various 'types' of likelihood



1. likelihood, marginal and conditional likelihood, profile likelihood, adjusted profile
2. semi-parametric likelihood, partial likelihood \leftarrow Cox 1972 proportional hazards regression misspecified models
3. quasi-likelihood, composite likelihood
4. empirical likelihood, penalized likelihood Cox regression
5. simulated likelihood, indirect inference
6. bootstrap likelihood, h -likelihood, weighted likelihood, pseudo-likelihood, local likelihood, sieve likelihood

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

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

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

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

Recap

- likelihood function is proportional to the probability *of the observed data*
density
or observable data
- need to assume a probability model in order to write down a likelihood function
- these models are usually parametric, i.e. a class of models that vary with a parameter
 $\theta \in \Theta \subseteq \mathbb{R}^p$
- but are sometimes non-parametric, in the sense that Θ might be an infinite-dimensional space
 - e.g. the class of all twice-differentiable functions
 - e.g. the intensity function for a Poisson process
- random effects model: why do we integrate out the random effects?

$$y_{ij} = \underbrace{\mu + x_i^T \beta}_{j=1, \dots, n_i} + \varepsilon_{ij} \quad \begin{matrix} \varepsilon_{ij} \sim N(0, \sigma^2) \\ i=1, \dots, n \end{matrix} \quad y_{ij} = \mu + x_i^T \beta + b_i + \varepsilon_{ij} \quad \begin{matrix} j=1, \dots, n_i \\ i=1, \dots, n \end{matrix} \quad \begin{matrix} b_i \sim N(0, \sigma_b^2) \\ \varepsilon_{ij} \sim N(0, \sigma^2) \end{matrix}$$

$$\ell(\beta, \sigma^2; \gamma) = \prod_{i=1}^n \prod_{j=1}^{n_i} f(y_{ij} | x_i, b_i; \beta, \sigma^2)$$

$$= \prod_{i=1}^n \prod_{j=1}^{n_i} \frac{e^{-\frac{1}{2\sigma^2} (y_{ij} - x_i^T \beta)^2}}{\sqrt{2\pi\sigma^2}}$$

$y_{ij} | b_i$ ind't over j
 $\sim N(\mu + x_i^T \beta, \sigma^2)$

$f(y_{ij} | x_i, b_i; \beta, \sigma^2)$

$$= \prod_{i=1}^n f(y_{ij} | x_i, b_i; \beta, \sigma^2)$$

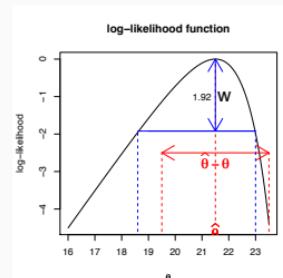
$$\prod_{i=1}^n f(y_{ij} | x_i, b_i; \beta, \sigma^2, \sigma_b^2) = \prod_{i=1}^n \int_{-\infty}^{\infty} f(y_{ij} | b_i; \beta, \sigma^2) \cdot f(b_i; \sigma_b^2) db_i$$

$y_{ij} \sim MVN(-, -)$ = $L(\beta, \sigma^2, \sigma_b^2; \gamma)$

$\text{var}(y_{ij}) = \sigma_b^2 + \sigma^2$
 $E(y_{ij}) = \mu + x_i^T \beta$
 $\text{cov}(y_{ij}, y_{ij'}) = \sigma_b^2$

$$= \text{cov}(b_{ij} + \varepsilon_{ij}, b_{jr} + \varepsilon_{jr})$$

- several examples: regression, time series, continuous time processes, correlated binary data, etc.
- several examples of likelihood functions that involve integration complicated
- an example where the likelihood function can't be written down completely Ising model
- these examples meant to motivate variations on the usual likelihood function to come
- notation and derived quantities: score function, observed and expected Fisher information, maximum likelihood estimate, likelihood ratio statistic



Don't forget the handout

STA 4508: Likelihood and derived quantities January 2022

Given a model for Y which assumes Y has a density $f(y; \theta)$, $\theta \in \Theta \subset \mathbb{R}^d$, we have the following definitions:

observed likelihood function	$L(\theta; y) = c(y)f(y; \theta)$
log-likelihood function	$\ell(\theta; y) = \log L(\theta; y) = \log f(y; \theta) + a(y)$
score function	$U(\theta) = \partial \ell(\theta; y) / \partial \theta$
observed information function	$j(\theta) = -\partial^2 \ell(\theta; y) / \partial \theta \partial \theta^T$
expected information (in one observation)	$i(\theta) = E_\theta U(\theta)U(\theta)^T$ (called $i_1(\theta)$ in CH)

When we have Y_i independent, identically distributed from $f(y_i; \theta)$, then, denoting the observed sample $y = (y_1, \dots, y_n)$ we have:

log-likelihood function	$\ell(\theta) = \ell(\theta; y) + a(y)$	$O_p(n)$
maximum likelihood estimate	$\hat{\theta} = \hat{\theta}(y) = \arg \sup_\theta \ell(\theta)$	$\theta + O_p(n^{-1/2})$
score function	$U(\theta) = \ell'(\theta) = \sum U_i(\theta) = U_+(\theta)$	$O_p(n^{1/2})$
observed information function	$j(\theta) = -\ell''(\theta) = -\ell(\theta; Y)$	$O_p(n)$
observed (Fisher) information	$j(\hat{\theta})$	
expected (Fisher) information	$i(\theta) = E_\theta \{U(\theta)U(\theta)^T\} = ni_1(\theta)$	$O(n)$,

STA 4508 January 19 2022 where with the risk of some confusion we use the same notation. Sometimes the expected Fisher information is defined instead as $i(\theta) = E_\theta \{-\partial U(\theta; Y) / \partial \theta^T\}$ (e.g.

... Recap: inference based on likelihood

- “pure likelihood”: values of θ are **plausible** if $L(\hat{\theta})/L(\theta)$ not too large
or $L(\theta)/L(\hat{\theta})$ not too small
- Bayesian inference: posterior \propto Likelihood \times prior
 $\int \pi(\theta|y) d\theta = 1$ $\pi(\theta | y) \propto L(\theta; y)\pi(\theta)$
 $= L(\theta, y)\pi(\theta)/m(y)$
- frequentist: quantities derived from the likelihood function have “good” properties
behave well when we have large samples from the model
- also frequentist: **pivotal quantities** derived from the likelihood function can be used to construct p -value functions
also called significance functions
- p -value functions provide nested sets of confidence intervals
if monotone in θ

A trio of limit results

y_1, \dots, y_n iid $f(y; \theta)$ $\theta \in \mathbb{R}$

1.

$$(*) \quad \frac{1}{\sqrt{n}} U(\theta) \xrightarrow{d} N\{0, i_1(\theta)\} \quad \text{CLT}$$

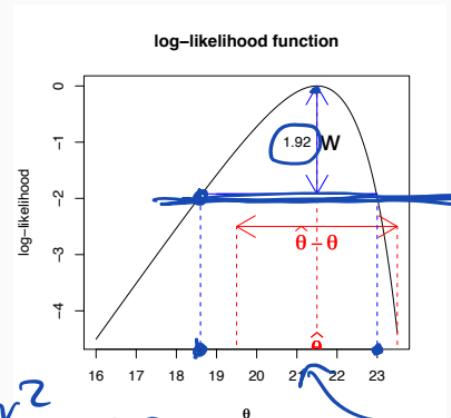
2.

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

3.

$$= [2\ell(\hat{\theta}) - \ell(\theta)] \xrightarrow{d} \chi^2_1$$

$w(\theta) \sim \chi^2_1$ 95% pt of $\chi^2_1 = 3.84$



$$U(\theta) \sim N(0, i_n(\theta))$$



$$\hat{\theta} \sim N(\theta, \frac{1}{i_n(\theta)})$$

Leading to a trio of approximate confidence intervals:

$$1. \quad \{\theta : |U(\theta)i^{-1/2}(\theta)| \leq z_{1-\alpha/2}\}$$

$$\hat{\theta} \pm 1.96 \sqrt{\frac{1}{i_n(\theta)}}$$

$\sim 95\%$ CI for θ

$$2. \quad \{\theta : |(\hat{\theta} - \theta)i^{1/2}(\hat{\theta})| \leq z_{1-\alpha/2}\}$$



$$3. \quad \{\theta : 2[\ell(\hat{\theta}) - \ell(\theta)]\} \leq \chi^2_{1,1-\alpha}$$

symmetric about $\hat{\theta}$
captures the asymmetry

p-value functions of θ

- or leading to a trio of approximate pivotal quantities

$$\begin{aligned}
 1. \quad r_u(\theta) &= U(\theta)j^{-1/2}(\hat{\theta}) \sim N(0, 1) & U(\theta) &\sim N(0, \dots) \\
 2. \quad r_e(\theta) &= (\hat{\theta} - \theta)j^{1/2}(\hat{\theta}), \sim N(0, 1) & r_u(\theta) &\sim N(0, 1) \\
 3. \quad \underline{r(\theta)} &= \text{sign}(\hat{\theta} - \theta)[2\{\ell(\hat{\theta}) - \ell(\theta)\}]^{1/2} \sim N(0, 1) \\
 && w(\theta) &\sim \chi^2_1
 \end{aligned}$$

piv. q: $f^*(y; \theta)$ with a known dist

$$C_{95} = \left\{ \theta : -1.96 \leq \underline{r(\theta)} \leq 1.96 \right\} \Leftrightarrow \left\{ \theta : \underline{\theta_L} \leq \theta \leq \underline{\theta_U} \right\}$$

$$\left\{ \theta : -1.96 \leq \underline{r_e(\theta)} \leq 1.96 \right\} \stackrel{\text{pivoting}}{\Leftrightarrow} \left\{ \theta : -1.96 \leq (\hat{\theta} - \theta)j^{1/2}(\hat{\theta}) \leq 1.96 \right\},$$

p-value functions of θ

$$\hat{\theta} \pm (-.96, j^{1/2}(\hat{\theta}))$$

- or leading to a trio of approximate pivotal quantities

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

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

$$3. r(\theta) = \text{sign}(\hat{\theta} - \theta)[2\{\ell(\hat{\theta}) - \ell(\theta)\}]^{1/2}$$

}

- $\Pr\{r_u(\theta) \leq r_u^0(\theta)\} \doteq \Phi\{r_u^0(\theta)\}$

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

p-value functions of θ

- or leading to a trio of approximate pivotal quantities

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

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

$$3. \quad r(\theta) = \text{sign}(\hat{\theta} - \theta)[2\{\ell(\hat{\theta}) - \ell(\theta)\}]^{1/2}$$

- $\Pr\{r_u(\theta) \leq r_u^0(\theta)\} \doteq \Phi\{r_u^0(\theta)\}$

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

- and a trio of *p*-value functions

$$\underbrace{H_0: \theta = \theta_0}_{\text{H}_0} \quad p\text{-value} = P_{\theta_0} \left\{ (\hat{\theta} - \theta) j^{1/2}(\hat{\theta}) \geq (\hat{\theta}^0 - \theta_0) j^{1/2}(\hat{\theta}^0) \right\}$$

↑
a.v.
↑
obs value
of θ , for fixed data

$$= P_{\theta_0} \left\{ r_e(\theta_0) \geq r_e^{\text{obs}}(\theta_0) \right\}$$

$$= P_{\theta_0} \left\{ \hat{\theta} > \hat{\theta}^{\text{obs}} \right\}$$

p-value functions of θ

- or leading to a trio of approximate pivotal quantities

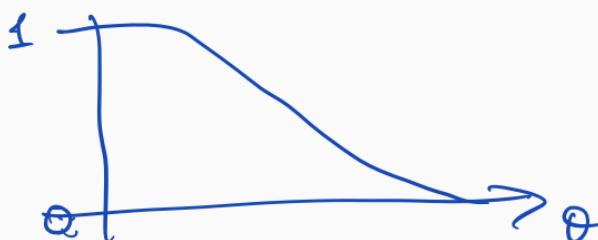
Efron & Hartley '78

$$\begin{aligned} \textcircled{1}. \quad r_u(\theta) &= U(\theta) j^{-1/2}(\hat{\theta}) \sim N(0, 1) && \xleftarrow{\text{Sometimes use}} i(\theta) \text{ inst.} \\ \rightarrow 2. \quad r_e(\theta) &= (\hat{\theta} - \theta) j^{1/2}(\hat{\theta}), \sim N(0, 1) \\ 3. \quad r(\theta) &= \text{sign}(\hat{\theta} - \theta)[2\{\ell(\hat{\theta}) - \ell(\theta)\}]^{1/2} \sim N(0, 1) \end{aligned}$$

- $\Pr\{r_u(\theta) \leq r_u^0(\theta)\} \doteq \Phi\{r_u^0(\theta)\}$

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

- and a trio of *p*-value functions
- similarly



$$\begin{aligned} 1. \quad p_u(\theta) &= \Phi\{r_u^0(\theta)\}, && \leftarrow \Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} dz \\ 2. \quad p_e(\theta) &= \Phi\{r_e(\theta)\} && \leftarrow \\ 3. \quad p_r(\theta) &= \Phi\{r(\theta)\} && \leftarrow \end{aligned}$$

Observed and expected Fisher information

$\frac{1}{\sqrt{n}} U(\theta) \xrightarrow{d} N(0, i_1(\theta))$ under $f(y; \theta)$
iid sample

app. 1 $\frac{U(\theta)}{\sqrt{i_n(\theta)}}$ ~ $N(0, 1)$

2 $\frac{U(\theta)}{\sqrt{i_n(\hat{\theta})}}$ ~ $N(0, 1)$

$\overbrace{x_n \xrightarrow{d} x}$
 $y_n \xrightarrow{P} a$
 $x \xrightarrow{P} a$

3. $\frac{U(\theta)}{\sqrt{j_n(\hat{\theta})}} \sim N(0, 1)$

$$\hat{i}_n(\theta) = n i_1(\theta)$$

$$\frac{U(\theta)}{\sqrt{i_n(\theta)}} \cdot \frac{\sqrt{i_n(\theta)}}{\sqrt{i_n(\hat{\theta})}} \xrightarrow{P} 1$$

Slutsky

$$= -\ell''(\hat{\theta}) = \hat{i}(\theta) = E\{-\ell''(\theta)\}$$

Example: Exponential

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

$$\log(1 + x) = x - \frac{1}{2}x^2 + \frac{1}{3}x^3 - \frac{1}{4}x^4 + \dots$$

- $\ell(\theta) = n \log \theta - \theta \sum y_i$

$$\Rightarrow \frac{\partial \ell}{\partial \theta} = \sum y_i \Rightarrow \hat{\theta} = \frac{n}{\sum y_i} = \bar{y}$$

- $\ell''(\theta) = -n/\theta^2 \quad i_i(\theta) = n/\theta^2 \quad i_i(\theta) = 1/\theta^2 \quad j(\theta) = n/\theta^2 \quad j(\hat{\theta}) = n/\hat{\theta}^2$

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

- $r_e(\theta) = (\hat{\theta} - \theta) j^{+1/2}(\hat{\theta}) = (\hat{\theta} - \theta) \frac{\sqrt{n}}{\hat{\theta}} = \sqrt{n} \left(1 - \frac{\theta}{\hat{\theta}}\right)$

- $r(\theta) = \left[2 \left\{ n \log \hat{\theta} - \frac{\hat{\theta} n}{\theta} - n \log \theta + \theta n / \hat{\theta} \right\} \right]^{1/2}$

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

$$= \left[2n \left\{ \log\left(\frac{\hat{\theta}}{\theta}\right) \frac{\hat{\theta}}{\theta} \left(1 - \frac{\theta}{\hat{\theta}}\right) \right\} \right]^{1/2} \quad \log\left(\frac{\hat{\theta}}{\theta}\right) = \log\left(\frac{\hat{\theta}}{\theta} - 1 + 1\right)$$

Example: Exponential

- $f(y_i; \theta) = \theta e^{-y_i \theta}, \quad i = 1, \dots, n$
- $\ell(\theta) = n \log \theta - n\theta \bar{y}$
- $\ell'(\theta) = \frac{n}{\theta} - n\bar{y}$ $\hat{\theta} = \bar{y}^{-1}$
- $\ell''(\theta) = -\frac{n}{\theta^2}$ $\frac{1}{\theta} - 1 \quad y \sim \text{exp}(\theta)$
- $r_u(\theta) = \frac{1}{\sqrt{n}} \ell'(\theta) j^{-1/2}(\hat{\theta}) = \sqrt{n} \left(\frac{1}{\theta \bar{y}} - 1 \right)$ $\hat{\theta} = \frac{n}{\sum y_i} = \frac{1}{\bar{y}}$
- $r_e(\theta) = (\hat{\theta} - \theta) j^{1/2}(\hat{\theta}) = \sqrt{n} (1 - \bar{y} \theta) \quad 1 - \bar{y} \theta$
- $r(\theta) = \sqrt{(2n)} \{ \theta \bar{y} - 1 - \log(\theta \bar{y}) \}^{1/2}$
expand $\log(\theta \bar{y})$ around 1 to get asymptotic equivalence to r_e, r_u
 $\sqrt{2} \left(\theta \bar{y} - 1 - \log \theta \bar{y} \right)^{1/2}$

... Example: Exponential

14

CHAPTER 2. UNCERTAINTY AND APPROXIMATION

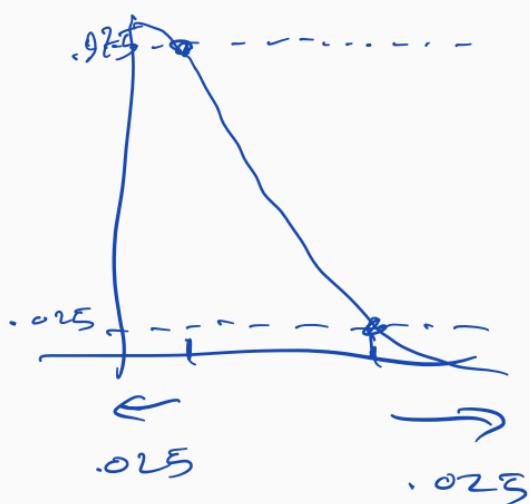
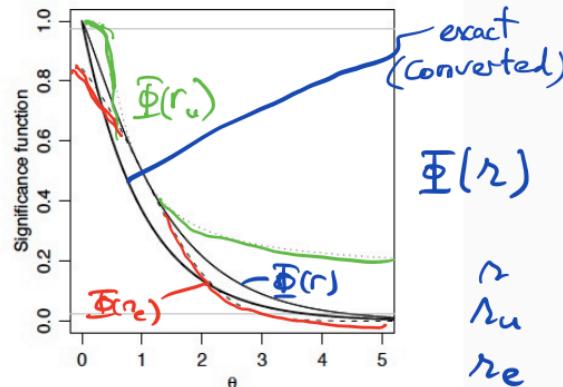
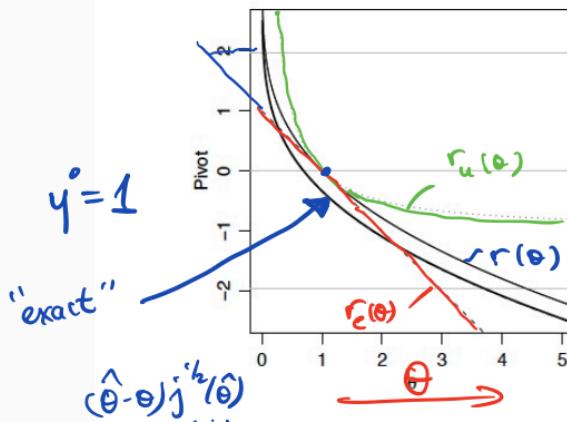
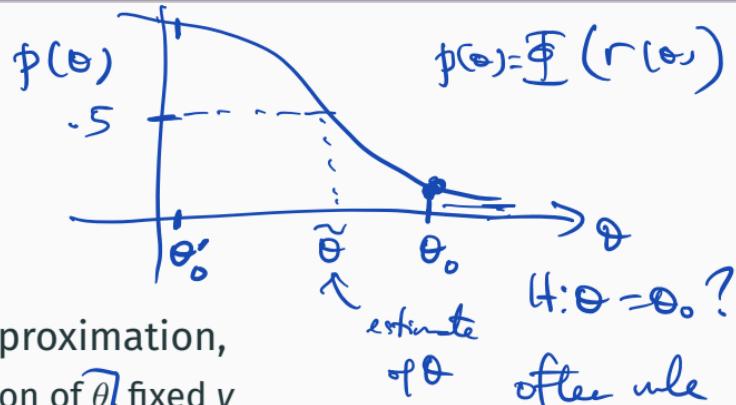


Figure 2.2: Approximate pivots and P-values based on an exponential sample of size $n = 1$. Left: likelihood root $r(\theta)$ (solid), score pivot $s(\theta)$ (dots), Wald pivot $t(\theta)$ (dashes), modified likelihood root $r^*(\theta)$ (heavy), and exact pivot $\theta \sum y_j$ (dot-dash). The modified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

Aside

- 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



]
← we fix θ , see if
under sampling y_1, \dots, y_n
this r.v. has 'good'
properties

Vector parameter limit theorems and approximations

1. $\cdot \underline{U}(\theta)$

$$\frac{1}{\sqrt{n}} \underline{U}(\underline{\theta}) \xrightarrow{d} N_p(0, i_n(\underline{\theta}))$$

$$\begin{pmatrix} u_1(\theta) \\ \vdots \\ u_p(\theta) \end{pmatrix} = \begin{pmatrix} \partial \ell / \partial \theta_1 \\ \vdots \\ \partial \ell / \partial \theta_p \end{pmatrix}$$

$\underline{U}(\theta) \sim N_p(0, j_n(\hat{\theta})) \sim N(0, i_n(\hat{\theta})) \sim N(0, i_n(\theta))$

$$i(\theta) = E\left(-\frac{\partial^2 \ell}{\partial \theta \partial \theta^T}\right)_{p \times p}$$

anti-sym

2. $\cdot \hat{\theta}$

$$\sqrt{n} (\hat{\theta} - \theta) \xrightarrow{d} N_p(0, i_n^{-1}(\theta))$$

$$\hat{\theta} \sim N_p(\theta, i_n^{-1}(\theta)) \sim N_p(\theta, i_n^{-1}(\hat{\theta})) \sim N_p(\theta, i_n^{-1}(\hat{\theta}))$$

3. $\cdot 2\{\ell(\hat{\theta}) - \ell(\theta)\}$



Parameter of interest and nuisance parameter

- $\theta = (\psi, \lambda) = (\psi_1, \dots, \psi_q, \lambda_1, \dots, \lambda_{p-q})$: $\psi \in \mathbb{R}$ usually
 $q=1$

Parameter of interest and nuisance parameter

- $\theta = (\psi, \lambda) =$

$$\ell(\theta; \gamma) = \ell(\psi, \lambda; \gamma)$$

- $U(\theta) = \begin{bmatrix} u_\psi(\theta) \\ u_\lambda(\theta) \end{bmatrix} = \begin{bmatrix} \frac{\partial \ell}{\partial \psi}(\psi, \lambda) \\ \frac{\partial \ell}{\partial \lambda}(\psi, \lambda) \end{bmatrix} \sim N_p \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}, i^*(\theta)$

$$i(\theta) = \begin{bmatrix} i_{\psi\psi}(\theta) & i_{\psi\lambda}(\theta) \\ i_{\lambda\psi}(\theta) & i_{\lambda\lambda}(\theta) \end{bmatrix}$$

$$i_{\psi\psi} = E\left(-\frac{\partial^2 \ell}{\partial \psi^2}\right) \quad i_{\psi\lambda} = E\left(-\frac{\partial^2 \ell}{\partial \psi \partial \lambda}\right)$$

Parameter of interest and nuisance parameter

- $\theta = (\psi, \lambda) =$

- $U(\theta) =$

- $i(\theta) =$

$$j_{44} = -\frac{\partial^2 l(\psi, \lambda)}{\partial \psi^2}$$

$$j(\theta) = \begin{bmatrix} j_{44}(\theta) & j_{42}(\theta) \\ j_{24}(\theta) & j_{22}(\theta) \end{bmatrix}$$

Parameter of interest and nuisance parameter

- $\theta = (\psi, \lambda) =$

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

- $U(\theta) =$

$$\hat{\psi} \sim N(\psi, (j^{-1}(\hat{\theta}))_{\psi\psi})$$

- $i(\theta) = \begin{pmatrix} i_{44} & i_{4\lambda} \\ i_{\lambda 4} & i_{\lambda\lambda} \end{pmatrix}$

$$j(\theta) = \begin{pmatrix} j_{44} & j_{4\lambda} \\ 0_{\lambda 4} & j_{\lambda\lambda} \end{pmatrix}$$

- $i^{-1}(\theta) = \begin{pmatrix} i^{44} & i^{4\lambda} \\ i^{\lambda 4} & i^{\lambda\lambda} \end{pmatrix}$

$$j^{-1}(\theta) = \begin{pmatrix} j^{44} & j^{4\lambda} \\ j^{\lambda 4} & j^{\lambda\lambda} \end{pmatrix}$$

$$j^{-1}(\theta) - j(\theta) = I$$

Parameter of interest and nuisance parameter

- $\theta = (\psi, \lambda) =$
- $U(\theta) =$
- $i(\theta) = \quad j(\theta) =$
- $i^{-1}(\theta) = \quad j^{-1}(\theta) =$
- $i^{\psi\psi}(\theta) =$

Parameter of interest and nuisance parameter

- $\theta = (\psi, \lambda) =$

$$\hat{\psi} \sim N(\psi, [i^{-1}(\theta)]_{\psi\psi})$$

- $U(\theta) =$

$$\sim N(\psi, i^{\psi\psi}(\theta))$$

- $i(\theta) =$

$$\begin{pmatrix} a & b \\ c & d \end{pmatrix}^{-1}$$

$$j(\theta) =$$

$$\frac{1}{a - bd^{-1}c}$$

- $i^{-1}(\theta) = \underbrace{\frac{1}{ad-bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}}$

$$j^{-1}(\theta) =$$

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

? ntbc?

$i^{\psi\psi}$ same

- $\ell_P(\psi) = \ell(\psi, \hat{\lambda}_\psi)$

~~$i_P(\theta)$~~

$$\ell'_P(\psi) = 0$$

$$\frac{\partial}{\partial \psi} \ell(\psi, \hat{\lambda}_\psi) = 0$$

↑

Profile

log-likelihood:

$$\hat{\lambda}_\psi \text{ maximizer } \ell(\psi, \lambda) \text{ wrt } \lambda$$

$$\frac{\partial}{\partial \psi} \ell(\psi, \lambda) + \frac{\partial}{\partial \lambda} \ell(\psi, \lambda) \frac{\partial \hat{\lambda}_\psi}{\partial \psi}$$

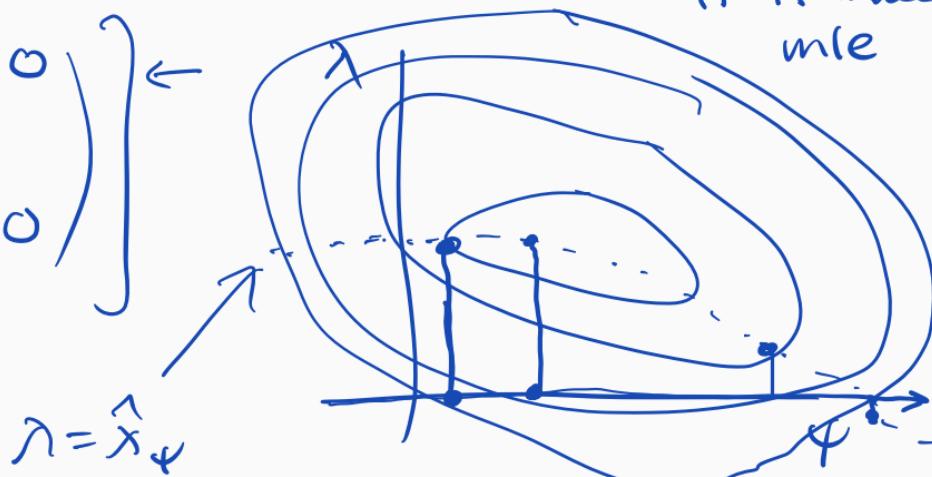
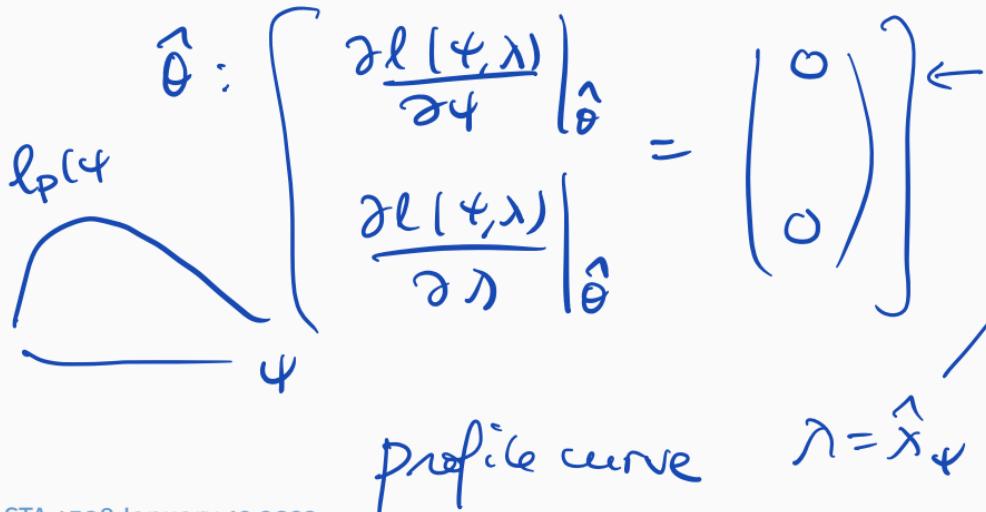
Nuisance parameters

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

~~fix λ~~

$$l_p'(\psi) = \frac{\partial l(\psi, \lambda)}{\partial \psi} + \frac{\partial l(\psi, \lambda)}{\partial \lambda} \frac{\partial \lambda}{\partial \psi}$$

$\underline{l_p'(\psi) = 0}$ defines an est. $\hat{\psi} \leftarrow$
it is indeed mle



Nuisance parameters

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

Nuisance parameters

Limit distribution

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

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

$$\frac{1}{m} U(\theta) \xrightarrow{d} N_p(0, i(\theta))$$

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

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N_p(0, i^{-1}(\theta))$$

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

$$2\{\ell(\hat{\theta}) - \ell(\theta)\} \xrightarrow{d} \chi_p^2$$

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

- $\ell_P(\psi) = \ell(\psi, \hat{\lambda}_\psi), \quad j_P(\psi) = -\ell''_P(\psi)$

$$\begin{pmatrix} \hat{\psi} - \psi \\ \hat{\lambda} - \lambda \end{pmatrix} \sim N_p \left[\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} i^{\psi\psi}(\theta) & i^{\psi\lambda}(\theta) \\ i^{\lambda\psi}(\theta) & i^{\lambda\lambda}(\theta) \end{pmatrix} \right]$$

Approximations from limiting distributions, nuisance parameters

$$(\hat{\psi} - \psi) \sim N_q(0, i^{**}(\hat{\theta})) \quad \hat{\theta} = (\hat{\psi}, \hat{\lambda}) : (\psi, \hat{\lambda}_\psi) \stackrel{d}{=} \hat{\theta}_q$$

$$\left. \frac{\partial \ell(\psi, \lambda)}{\partial \psi} \right|_{\lambda = \hat{\lambda}_\psi}$$

$$q \geq 1$$

$$w_u(\psi) = U_\psi(\psi, \hat{\lambda}_\psi)^T \{ i^{\psi\psi}(\psi, \hat{\lambda}_\psi) \} U_\psi(\psi, \hat{\lambda}_\psi) \underset{q \geq 1}{\sim} \chi_q^2$$

$$\rightarrow w_e(\psi) = (\hat{\psi} - \psi) \{ i^{\psi\psi}(\hat{\psi}, \hat{\lambda}) \}^{-1} (\hat{\psi} - \psi) \underset{q \geq 1}{\sim} \chi_q^2$$

$$\rightarrow w(\psi) = 2\{\ell(\hat{\psi}, \hat{\lambda}) - \ell(\psi, \hat{\lambda}_\psi)\} = 2\{\ell_P(\hat{\psi}) - \ell_P(\psi)\} \underset{q \geq 1}{\sim} \chi_q^2;$$

$$q=1 \quad w_e \sim \chi_1^2 \quad \underbrace{(\hat{\psi} - \psi)^2 \{ i^{**}(\hat{\psi}, \hat{\lambda}) \}^{-1}}_{\text{approximate pivot}} \underset{q=1}{\sim} \mathcal{N}(0, 1)$$

Approximate Pivots, $q = 1$

$\left\{ \begin{array}{l} \pm \sqrt{w_u} = r_u(\psi) = \ell'_P(\psi) j_P(\hat{\psi})^{-1/2} \sim N(0, 1), \\ \pm \sqrt{w_e} = r_e(\psi) = (\hat{\psi} - \psi) j_P(\hat{\psi})^{1/2} \sim N(0, 1), \\ \pm \sqrt{w} = r(\psi) = \text{sign}(\hat{\psi} - \psi)[2\{\ell_P(\hat{\psi}) - \ell_P(\psi)\}]^{1/2} \sim N(0, 1) \end{array} \right.$

$$\pm \sqrt{w} \sim N(0, 1)$$

$$\ell'_P(\psi) = \frac{\partial \ell(\psi, \hat{\lambda}_\psi)}{\partial \psi} + \frac{\partial \ell(\psi, \hat{\lambda}_\psi)}{\partial \lambda}$$

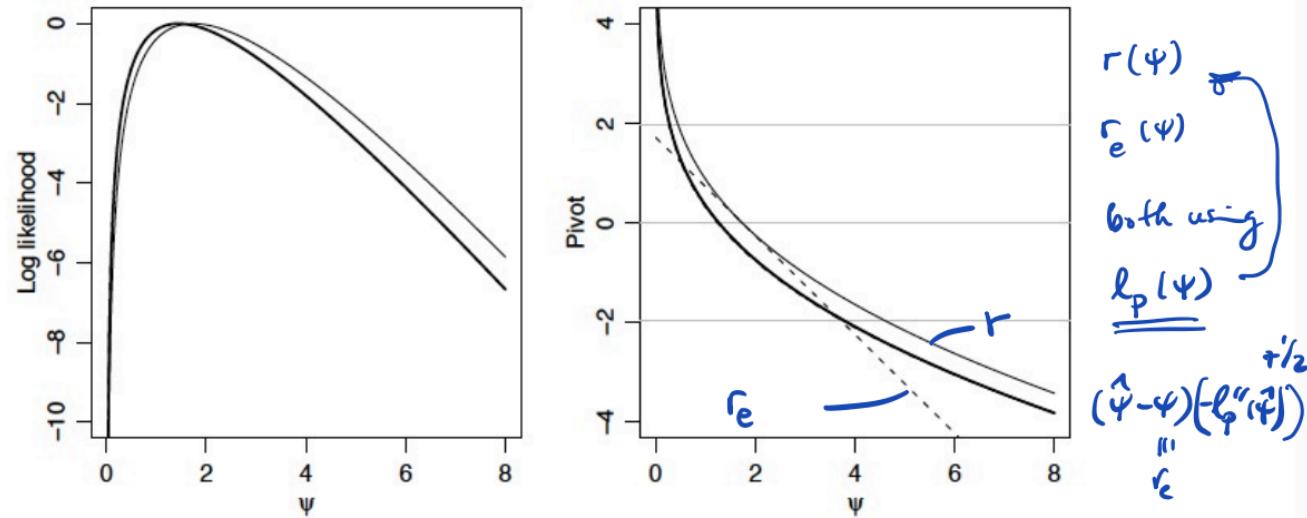


Figure 2.3: Inference for shape parameter ψ of gamma sample of size $n = 5$. Left: profile log likelihood ℓ_p (solid) and the log likelihood from the conditional density of u given v (heavy). Right: likelihood root $r(\psi)$ (solid), Wald pivot $t(\psi)$ (dashes), modified likelihood root $r^*(\psi)$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

$$\left\{ \begin{array}{l} r_u(\theta) = U(\theta) j^{-1/2}(\hat{\theta}) \sim N(0, 1) \\ r_e(\theta) = (\hat{\theta} - \theta) j^{1/2}(\hat{\theta}) \sim N(0, 1) \\ r(\theta) = \pm \sqrt{2 \{ \ell(\hat{\theta}) - \ell(\theta) \}} \sim N(0, 1) \end{array} \right.$$

$$\left\{ \begin{array}{l} r_u(\psi) = l_p'(\psi) j_p^{-1/2}(\hat{\psi}) \quad j_p(\psi) = -l_p''(\psi) \\ \qquad \qquad \qquad l_p(\psi) = \ell(\psi, \lambda_p) \\ r_e(\psi) = (\hat{\psi} - \psi) j_p^{1/2}(\hat{\psi}) \\ r(\psi) = \pm \sqrt{2 \{ l_p(\hat{\psi}) - l_p(\psi) \}} \\ \text{all } \sim N(0, 1) \quad (\text{tmt theorem}) \end{array} \right.$$

$$\begin{pmatrix} \hat{\psi} \\ \hat{\lambda} \end{pmatrix} \sim N \left(\begin{pmatrix} \psi \\ \lambda \end{pmatrix}, i^{-1}(\theta) \right)$$

$$\hat{\psi} \sim N(\psi, i^{\psi\psi}(\hat{\theta}))$$

pirotal q.: $(\hat{\psi} - \psi) \{ j^{\psi\psi}(\hat{\psi}, \lambda_p) \}^{-1} \sim N(0, 1)$

... Laplace approximation

marginal posterior:

$$(\hat{\psi} - \psi) |_{\hat{\psi}}^{\text{h}} \sim N(0, 1)$$

Profile likelihood: examples

- regression

$$\underline{\underline{\psi = \beta_j}}$$

$$y = \underline{\underline{X\beta}} + \epsilon, \quad \epsilon \sim N(0, \sigma^2), \quad \underline{\underline{\psi = \sigma^2}}$$
$$\hat{\sigma}^2 = \frac{1}{n}(y - X\hat{\beta})^T(y - X\hat{\beta})$$

$$Y = \begin{matrix} y_1 \\ \vdots \\ y_n \end{matrix}$$

$$\hat{\beta} = (X^T X)^{-1} X^T y \quad l(\beta, \sigma^2)$$

$$l(\sigma^2, \beta) = -\frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} (y - X\beta)^T (y - X\beta)$$

$$\left(\begin{array}{l} \frac{\partial l}{\partial \beta} = 0 \\ \frac{\partial l}{\partial \sigma^2} = 0 \end{array} \right) \Rightarrow -\frac{1}{2\sigma^2} \cdot 2 \cdot (y - X\beta) X^T = 0 \quad \text{ntb}$$
$$\hat{\beta} = (X^T X)^{-1} X^T y = \hat{\beta}_{\sigma^2}$$

$$\frac{\partial \ell}{\partial \sigma^2} = \ell(\sigma^2, \hat{\beta}_{\sigma^2}) - \ell(\sigma^2, \hat{\beta})$$

$$\ell_P(\sigma^2) = -\frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} (y - x\hat{\beta})^T (y - x\hat{\beta})$$

$$\ell_P'(\sigma^2) = -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} (y - x\hat{\beta})^T (y - x\hat{\beta})$$

$$\left. \ell_P'(\sigma^2) \right|_{\sigma^2=0} = \frac{-n\sigma^2}{2\sigma^4} + \frac{1}{2\sigma^4} (y - x\hat{\beta})^T (y - x\hat{\beta})$$

$$\hat{\sigma}^2 = \frac{(y - x\hat{\beta})^T (y - x\hat{\beta})}{n}$$

$$E \hat{\sigma}^2 = \frac{(n-p)\sigma^2}{n} \quad \text{not } \sigma^2$$

$$E \hat{\sigma}^2 < \sigma^2 \quad \hat{\sigma}^2 = \frac{(y - x\hat{\beta})^T (y - x\hat{\beta})}{n-p}$$

- yes, we do have issues w/ profile lik
when p is large relative to n

- limit theory gives
poor approximations when n is fixed

+ p is large

Profile likelihood: examples

- regression

$$y = X\beta + \epsilon, \quad \epsilon \sim N(0, \sigma^2), \quad \psi = \sigma^2$$

$$\hat{\sigma}^2 = \frac{1}{n}(y - X\hat{\beta})^T(y - X\hat{\beta})$$

- Neyman-Scott

$$\dim(\underline{x}) = \dim(\underline{\mu}) \rightarrow \infty \text{ if } m \rightarrow \infty$$

$$\Psi = \sigma^2 \quad \lambda = (\mu_1, \mu_2, \dots, \mu_m)$$

$$y_{ij} \sim N(\hat{\mu}_i = \bar{y}_{i \cdot}, \sigma^2), \quad j = 1, \dots, k, \quad i = 1, \dots, m$$
$$\hat{\sigma}^2 = \frac{1}{mk} \sum_{j=1}^m (y_{ij} - \bar{y}_{i \cdot})^2 \quad \text{too small}$$

$\ell_p(\sigma^2)$ too

concentrated

$$E \hat{\sigma}^2 = \frac{1}{mk} \sigma^2 (k-1) m$$

$\hat{\sigma}^2 \rightarrow \sigma^2$ if $k \rightarrow \infty$ but not

Profile likelihood: examples

- regression

$$y = X\beta + \epsilon, \quad \epsilon \sim N(0, \sigma^2), \quad \psi = \sigma^2$$

$$\hat{\sigma}^2 = \frac{1}{n}(y - X\hat{\beta})^T(y - X\hat{\beta})$$

- Neyman-Scott

$$y_{ij} \sim N(\mu_i, \sigma^2), j = 1, \dots, k; i = 1, \dots, m$$

$$\hat{\sigma}^2 = \frac{1}{mk} \sum_{i=1}^m (y_{ij} - \bar{y}_{i.})^2 \quad \not\rightarrow \quad \frac{k-1}{k} \sigma^2 \neq \sigma^2$$

if $m \rightarrow \infty$

- 2×2 tables

$$y_{i1} \sim \text{Bern}(p_{i1}), y_{i2} \sim \text{Bern}(p_{i2}), i = 1, \dots, n, \quad \log \left\{ \frac{p_{i1}/(1-p_{i1})}{p_{i2}/(1-p_{i2})} \right\} = \psi + \lambda_i$$

↑ n.i.s. p.o.t.

1 case $y_{i1} = 1$ if

case is ~~if~~, else 0
Samples are equal

$$\hat{\psi} \xrightarrow{P} \psi/2$$

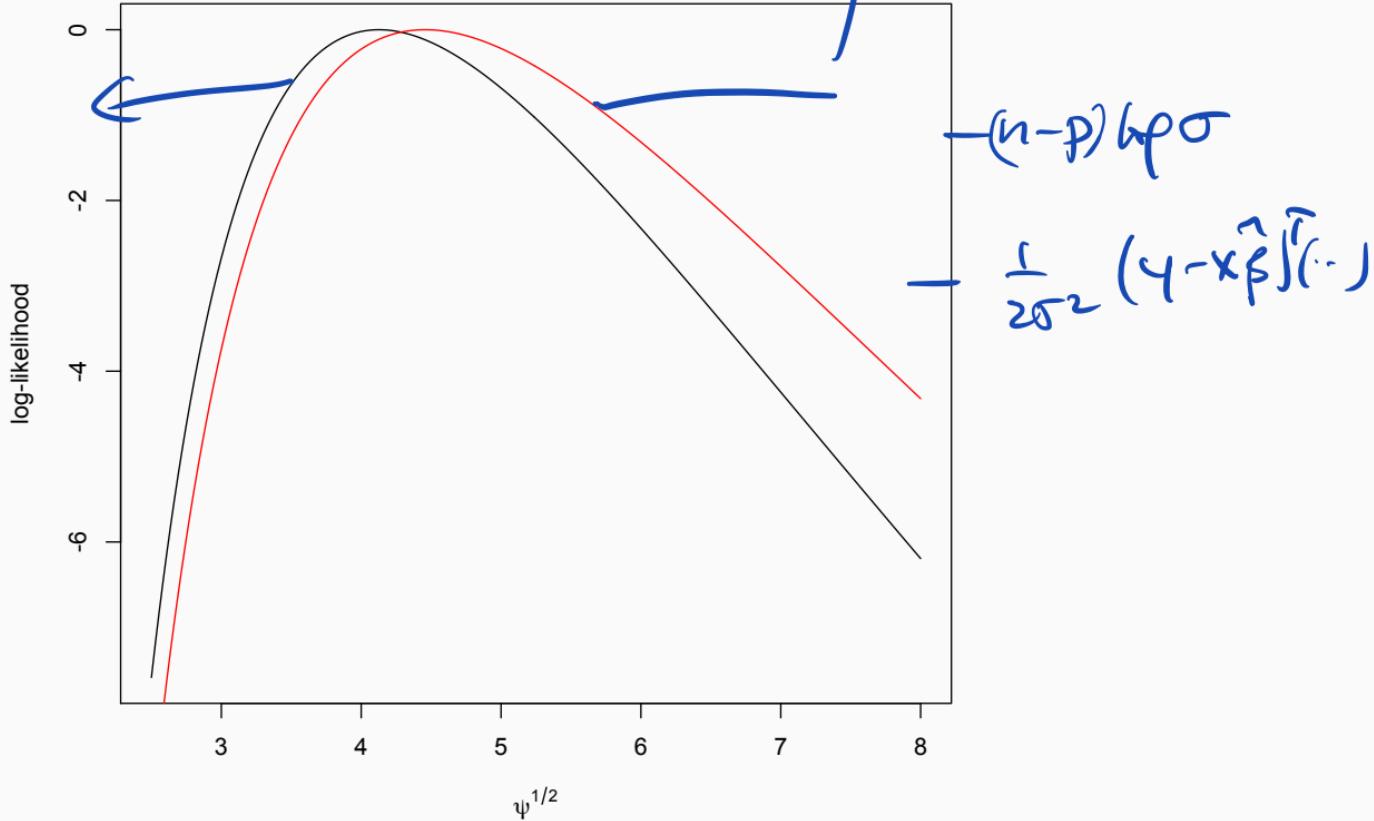
$$n \rightarrow \infty$$

2 controls $y_{i2} = 1$ if ~~case is 1~~

$l_a(\sigma^2)$

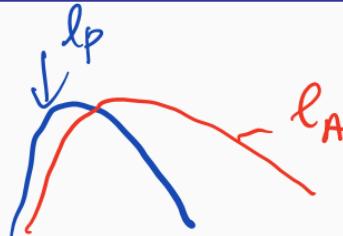
This is a plot of $-n \log \sigma - (y - X\hat{\beta})^T(y - X\hat{\beta})/2\sigma^2$ (black), and $-(n-p) \log \sigma - (y - X\hat{\beta})^T(y - X\hat{\beta})/2\sigma^2$ against σ (red) for given data

profile
 $l_p(\sigma^2)$



Approximate conditional inference

$$\cdot \ell_c(\psi) \doteq \ell_p(\psi) - \frac{1}{2} \log |j_{\lambda\lambda}(\psi, \hat{\lambda}_\psi)|$$



$$\cdot \ell_m(\psi) \doteq \ell_p(\psi) - \frac{1}{2} \log |j_{\lambda\lambda}(\psi, \hat{\lambda}_\psi)|$$

$$\cdot \ell_c(\psi) \doteq \ell_p(\psi) + \frac{1}{2} \log |j_{\eta\eta}(\psi, \hat{\eta}_\psi)|$$

$$i_{\psi\lambda}(\theta) = \exp\{\psi^T s + \eta^T t - c(\psi, \eta)\}$$

• adjusted profile log-likelihood

$$\boxed{\ell_A(\psi) = \ell_p(\psi) + A(\psi)}$$

$A(\psi)$ assumed to be $O_p(1)$

$$\cdot \text{generic form is } A_{FR}(\psi) = +\frac{1}{2} \log |j_{\lambda\lambda}(\psi, \hat{\lambda}_\psi)| - \log \left| \frac{d(\lambda)}{d\hat{\lambda}_\psi} \right|$$

Fraser 03

$$\cdot \text{closely related } A_{BN}(\psi) = -\frac{1}{2} \log |j_{\lambda\lambda}(\psi, \hat{\lambda}_\psi)| + \log \left| \frac{d\hat{\lambda}}{d\hat{\lambda}_\psi} \right|$$

SM §12.4.1

Cox & Reid adjusted likelihood

1987 showed that in general

it's a good 1st fix to

$$Q_{CR}^{(\psi)} = Q_a(\psi) = \underline{\underline{L_p(\psi)}} - \frac{1}{2} \log |\underline{\underline{j_m(\psi, \lambda_\psi)}}|$$

we need to work in the parametrization

(ψ, λ) in which $\psi \perp\!\!\!\perp \lambda$ (orthog)

meaning $i_{\psi\lambda}(\theta) = 0$

exp'd F. info.

$$\begin{bmatrix} i_{\psi\psi} & 0 \\ 0 & \text{Fix} \end{bmatrix}$$