3.1 Basic Results slide 93

Motivation

- □ Likelihood
 - provides a general paradigm for inference on parametric models, with many generalisations and variants:
 - uses only minimal sufficient statistics;
 - is a central concept in both frequentist and Bayesian statistics;
 - has a simple, general and widely-applicable 'large-sample' theory; but
 - is not a panacea!
- ☐ Plan below:
 - give (fairly) general setup;
 - prove main results for scalar parameter;
 - discussion of inference;
 - vector parameter, nuisance parameters, ...

stat.epfl.ch Autumn 2024 – slide 94

Basic setup

Let $Y, Y_1, \ldots, Y_n \stackrel{\text{iid}}{\sim} g$, and define the Kullback-Leibler divergence from the data-generating model g to a candidate density f,

$$\mathrm{KL}(g, f) = \mathrm{E}_g\{\log g(Y) - \log f(Y)\} = \mathrm{E}_g\left[-\log\left\{\frac{f(Y)}{g(Y)}\right\}\right] \ge 0,$$

where the inequality holds because $-\log x$ is convex and is strict unless $f \equiv g$ (Jensen).

□ In a parametric setting f belongs to a parametric family $\mathcal{F} = \{f_{\theta} : \theta \in \Theta\}$, so minimising $\mathrm{KL}(g,f)$ over f is equivalent to maximising $\mathrm{E}_g \log f(Y;\theta)$, which is estimated by

$$\overline{\ell}(\theta) = n^{-1} \sum_{j=1}^{n} \log f(Y_j; \theta) \xrightarrow{P} E_g \log f(Y; \theta), \quad n \to \infty.$$

- $\square \quad \theta_g = \operatorname{argmax}_{\theta} \mathbb{E}_g \log f(Y; \theta) \text{ gives the optimal large-sample fit of } f_{\theta} \text{ to } g.$
- \square In an ideal case $g \in \mathcal{F}$, so $g = f_{\theta_g}$, but the theory does not require this (yet).
- \square The natural estimator of θ_q is the maximum likelihood estimator

$$\widehat{\theta} = \operatorname{argmax}_{\theta} \overline{\ell}(\theta),$$

but we need conditions on $\overline{\ell}$ to ensure that $\widehat{\theta} \stackrel{P}{\longrightarrow} \theta_g$ or (better) $\widehat{\theta} \stackrel{\text{a.s.}}{\longrightarrow} \theta_g$ as $n \to \infty$.

stat.epfl.ch

Regular models

- $\square \quad \text{Notation: } \nabla h(\theta) = \partial h(\theta)/\partial \theta \text{ and } \nabla^2 h(\theta) = \nabla \nabla^{\mathrm{T}} h(\theta) = \partial^2 h(\theta)/\partial \theta \partial \theta^{\mathrm{T}}.$
- ☐ The asymptotic properties of the MLE rely on regularity conditions:
 - (C1) θ_q is unique and interior to $\Theta \subset \mathbb{R}^d$ for some finite d, and Θ is compact;
 - (C2) the densities f_{θ} defined by any two different values of $\theta \in \Theta$ are distinct;
 - (C3) there is a neighbourhood $\mathcal N$ of θ_g within which the first three derivatives of the log likelihood with respect to θ exist almost surely, and for $r,s,t=1,\ldots,d$ satisfy $|\partial^3 \log f(Y;\theta)/\partial \theta_r \partial \theta_s \partial \theta_t| < m(Y)$ with $\mathrm{E}_q\{m(Y)\} < \infty$; and
 - **(C4)** within \mathcal{N} , the $d \times d$ matrices

$$i_1(\theta) = \mathcal{E}_g \left\{ -\nabla^2 \log f(Y; \theta) \right\}, \quad \hbar_1(\theta) = \mathcal{E}_g \left\{ \nabla \log f(Y; \theta) \nabla^{\mathrm{T}} \log f(Y; \theta) \right\},$$

are finite and positive definite. When $g=f_{\theta_g}$ we shall see that $\hbar_1(\theta_g)=\imath_1(\theta_g)$.

stat.epfl.ch Autumn 2024 – slide 96

Regularity conditions

- \Box (C1) ensures that $\widehat{\theta}$ can be 'on all sides' of θ_g in the limit if it fails, then any limiting distribution cannot be normal;
- \square (C2) is essential for consistency, otherwise $\widehat{\theta}$ might not converge it often fails in mixture models, for which care is needed;
- □ (C3) is needed to bound terms of a Taylor series can be replaced by other conditions, see van der Vaart (1998, *Asymptotic Statistics*, Chapter 5); and
- \Box (C4) ensures that the asymptotic variance of $\widehat{ heta}$ is positive definite.

stat.epfl.ch Autumn 2024 – slide 97

Consistency of the MLE

Lemma 49 If $Y_1, \ldots, Y_n \sim g$ and $n \to \infty$, then under (C1) and (C2) a sequence of maximum likelihood estimators $\widehat{\theta}$ exists such that $\widehat{\theta} \stackrel{P}{\longrightarrow} \theta_q$.

This result:

- \Box does not require f_{θ} to be smooth, so it is quite general;
- \square guarantees that a consistent sequence exists, but not that we can find it;
- but if the log likelihood is concave (as in exponential families, for example), then there is (at most) one maximum for any n, and if it exists this must converge to θ_q ;
- \Box can be generalized to vector θ , but the argument is more delicate.

Note to Lemma 49

- \square We prove this for θ scalar.
- \square As the θ s correspond to different densities, precisely one θ_q minimises $\mathrm{KL}(g, f_\theta)$.
- $\Box \quad \text{Take any } \varepsilon > 0 \text{ and let } \theta_+, \theta_- = \theta_g \pm \varepsilon, \text{ write } D_n(\theta) = \overline{\ell}(\theta_g) \overline{\ell}(\theta), \text{ so } D_n(\theta_g) = 0, \text{ and note that as } n \to \infty,$

$$D_n(\theta_+) \xrightarrow{P} \mathrm{KL}(g, f_{\theta_+}) - \mathrm{KL}(g, f_{\theta_g}) = a_+ > 0, \quad D_n(\theta_-) \xrightarrow{P} \mathrm{KL}(g, f_{\theta_-}) - \mathrm{KL}(g, f_{\theta_g}) = a_- > 0.$$

 \square If A_n and B_n denote the events $D_n(\theta_+) > 0$ and $D_n(\theta_-) > 0$, Boole's inequality gives

$$P(A_n \cap B_n) = 1 - P(A_n^c \cup B_n^c) \ge 1 - P(A_n^c) - P(B_n^c).$$

Now

$$P(A_n^c) = P\{D_n(\theta_+) < 0\} = P\{a_+ - D_n(\theta_+) > a_+\} < P\{|D_n(\theta_+) - a_+| > a_+\} \to 0, \quad n \to \infty,$$

and likewise $P(B_n^c) \to 0$. Hence $P(A_n \cap B_n) \to 1$.

Hence there is a local minimum of $D_n(\theta)$, or equivalently a local maximum of $\overline{\ell}(\theta)$, inside the interval $(\theta_g - \epsilon, \theta_g + \epsilon)$ with probability one as $n \to \infty$, and as this is true for arbitrary ε , the corresponding sequence of maximisers $\widehat{\theta}$ satisfies $P(|\widehat{\theta} - \theta_g| > \varepsilon) \to 0$ and therefore is consistent.

stat.epfl.ch

Autumn 2024 - note 1 of slide 98

Asymptotic normality of the MLE

Theorem 50 If $Y_1, \ldots, Y_n \stackrel{\text{iid}}{\sim} g$, then under (C1)–(C4) the consistent sequence of maximum likelihood estimators $\widehat{\theta}$ satisfies

$$n^{1/2}(\widehat{\theta} - \theta_q) \xrightarrow{D} \mathcal{N}_d\{0, i_1^{-1}(\theta_q)\hbar_1(\theta_q)i_1^{-1}(\theta_q)\}, \quad n \to \infty,$$

where for a single observation Y we define

$$i_1(\theta) = \mathrm{E}_g \left\{ -\nabla^2 \log f(Y; \theta) \right\}, \quad \hbar_1(\theta) = \mathrm{E}_g \left\{ \nabla \log f(Y; \theta) \nabla^{\mathrm{T}} \log f(Y; \theta) \right\}.$$

 \square This implies that for large n we can use the approximation

$$\widehat{\theta} \stackrel{\cdot}{\sim} \mathcal{N}_d \{ \theta_g, \imath^{-1}(\theta_g) \hbar(\theta_g) \imath^{-1}(\theta_g) \},$$

where $i(\theta) = ni_1(\theta)$ and $\hbar(\theta) = n\hbar_1(\theta)$ correspond to a random sample of size n.

☐ This provides tests and confidence intervals based on the approximate pivots

$$v_{rr}^{-1/2}(\widehat{\theta}_r - \theta_{g,r}) \sim \mathcal{N}(0,1), \quad r = 1, \dots, d,$$

where v_{rr} are the diagonal elements of an estimate of $i^{-1}(\theta_q)\hbar(\theta_q)i^{-1}(\theta_q)$.

 \square When $g = f_{\theta_g}$, $i_1(\theta_g) = h_1(\theta_g)$ and the variance (matrix) becomes $i(\theta_g)^{-1}$.

stat.epfl.ch

Note to Theorem 50: A (fairly) simple argument

□ Write

$$0 = \nabla \overline{\ell}(\widehat{\theta}) = \nabla \overline{\ell}(\theta_g) + \int_0^1 \nabla^2 \overline{\ell} \{\theta_g + t(\widehat{\theta} - \theta_g)\} (\widehat{\theta} - \theta_g) dt,$$

and note that $U_n = n^{1/2} \nabla \overline{\ell}(\theta_g) \stackrel{D}{\longrightarrow} U \sim \mathcal{N}_d\{0, h_1(\theta_g)\}$, so writing $Z_n = n^{1/2} (\widehat{\theta} - \theta_g)$ we have

$$i_1(\theta_g)^{-1}U_n = i_1(\theta_g)^{-1} \left\{ -\int_0^1 \nabla^2 \overline{\ell}(\theta_g + tn^{-1/2}Z_n) dt \right\} Z_n = i_1(\theta_g)^{-1}J_n^* Z_n,$$

say, and as $n \to \infty$, $J_n^* \doteq -\int_0^1 \nabla^2 \overline{\ell}(\theta_g) \, \mathrm{d}t \stackrel{P}{\longrightarrow} \imath_1(\theta_g)$ and thus $\imath_1(\theta_g)^{-1} J_n^* \stackrel{P}{\longrightarrow} I_d$. Hence

$$i_1(\theta_g)^{-1}J_n^*Z_n = i_1(\theta_g)^{-1}U_n \xrightarrow{D} i_1(\theta_g)^{-1}U \sim \mathcal{N}_d\{0, i_1(\theta_g)^{-1}\hbar_1(\theta_g)i_1(\theta_g)^{-1}\}.$$

For a more careful treatment of the integral, we need a uniform law of large numbers (ULLN), which requires that $J_n(\theta) = -\nabla^2 \overline{\ell}(\theta)$ is measurable and continuous in θ within a compact subset \mathcal{N}' of \mathcal{N} , for almost all y, and that there exists a function d(Y) whose expectation is finite and for which $\|J_n(\theta)\| < d(Y)$ for all $\theta \in \mathcal{N}'$, where $\|\cdot\|$ is a matrix norm. Then $\mathrm{E}\{J_n(\theta)\} = \imath_1(\theta)$ is continuous in θ and

$$\sup_{\theta \in \mathcal{N}} \|J_n(\theta) - i_1(\theta)\| \stackrel{P}{\longrightarrow} 0, \quad n \to \infty.$$

 $\Box \quad \text{Let } \delta > 0 \text{ be small enough that } B_\delta = \{\theta: |\theta - \theta_g| \leq \delta\} \subset \mathcal{N}' \text{ and let } A_n = \{|n^{-1/2}Z_n| \leq \delta\} \text{ and } C_n = \|J_n^* - \iota_1(\theta_g)\|. \text{ Then for } \varepsilon > 0 \text{ we have}$

$$P(C_n > \varepsilon) = P(\{C_n > \varepsilon\} \cap A_n) + P(\{C_n > \varepsilon\} \cap A_n^c) \le P(\{C_n > \varepsilon\} \cap A_n) + P(A_n^c),$$

where the last term tends to zero because $n^{-1/2}Z_n=\widehat{\theta}-\theta_g\stackrel{P}{\longrightarrow} 0$. Now if A_n holds, then $\theta_g+tn^{-1/2}Z_n\in B_\delta$ when $0\leq t\leq 1$, so

$$C_{n} = \left\| \int_{0}^{1} \left\{ J_{n}(\theta_{g} + tn^{-1/2}Z_{n}) - \imath_{1}(\theta_{g} + tn^{-1/2}Z_{n}) + \imath_{1}(\theta_{g} + tn^{-1/2}Z_{n}) - \imath_{1}(\theta_{g}) \right\} dt \right\|$$

$$\leq \int_{0}^{1} \left\| J_{n}(\theta_{g} + tn^{-1/2}Z_{n}) - \imath_{1}(\theta_{g} + tn^{-1/2}Z_{n}) \right\| dt + \int_{0}^{1} \left\| \imath_{1}(\theta_{g} + tn^{-1/2}Z_{n}) - \imath_{1}(\theta_{g}) \right\| dt$$

$$\leq \sup_{\theta \in B_{\delta}} \|J_{n}(\theta) - \imath_{1}(\theta)\| + \sup_{\theta \in B_{\delta}} \|\imath_{1}(\theta) - \imath_{1}(\theta_{g})\|$$

$$= D_{n} + E_{n},$$

say. If $C_n > \varepsilon$ then at least one of D_n and E_n must exceed $\varepsilon/2$, so

$$\begin{split} \mathrm{P}(\{C_n > \varepsilon\} \cap A_n) & \leq & \mathrm{P}\left(\{\{D_n > \varepsilon/2\} \cup \{E_n \geq \varepsilon/2\}\} \cap A_n\right) \\ & \leq & \mathrm{P}(\{D_n > \varepsilon/2\} \cap A_n) + \mathrm{P}(\{E_n \geq \varepsilon/2\} \cap A_n) \\ & \leq & \mathrm{P}(D_n > \varepsilon/2) + \mathrm{P}(E_n > \varepsilon/2). \end{split}$$

Now $D_n \stackrel{P}{\longrightarrow} 0$ using the ULLN, and the continuity of $i_1(\theta)$ at θ_g implies that E_n can be made smaller than $\varepsilon/2$ by a suitable choice of $\delta > 0$, in which case

$$P(C_n > \varepsilon) \leq P(\{C_n > \varepsilon\} \cap A_n) + P(A_n^c)$$

$$\leq P(D_n > \varepsilon/2) + P(E_n > \varepsilon/2) + P(A_n^c)$$

$$\to 0, \quad n \to \infty,$$

which implies that $J_n^* \xrightarrow{P} \imath_1(\theta_g)$ and therefore that $\imath_1(\theta_g)^{-1} J_n^* \xrightarrow{P} I_d$, as required.

stat.epfl.ch

Autumn 2024 - note 1 of slide 99

Note to Theorem 50: Another approach

 \square We first note that under the given conditions, θ_g gives a stationary point of $\mathrm{KL}(g,f_\theta)$, and therefore

$$0 = \nabla \mathrm{KL}(g, f_{\theta})|_{\theta = \theta_g} = -\nabla \int \log f(y; \theta) g(y) \, \mathrm{d}y \bigg|_{\theta = \theta_g} = -\int \nabla \log f(y; \theta) \bigg|_{\theta = \theta_g} g(y) \, \mathrm{d}y,$$

so $E_g\{\nabla \log f(Y;\theta)\} = 0$.

 \square As $\widehat{\theta}$ gives a local maximum of the differentiable function $\overline{\ell}(\theta) = n^{-1} \sum_{j=1}^n \log f(Y_j; \theta)$,

$$0 = \nabla \overline{\ell}(\widehat{\theta}) = n^{-1} \sum_{j=1}^{n} \nabla \log f(Y_j; \widehat{\theta}),$$

and (supposing now that θ is scalar, to simplify the expressions), Taylor series expansion gives

$$0 = \nabla \overline{\ell}(\theta_g) + (\widehat{\theta} - \theta_g) \nabla^2 \overline{\ell}(\theta_g) + \frac{1}{2} (\widehat{\theta} - \theta_g)^2 \nabla^3 \overline{\ell}(\theta^*),$$

where θ^* lies between θ_g and $\widehat{\theta}$ (so $\theta^* \stackrel{P}{\longrightarrow} \theta_g$). Hence

$$n^{1/2}(\widehat{\theta} - \theta_g) = \frac{n^{1/2} \nabla \overline{\ell}(\theta_g)}{-\nabla^2 \overline{\ell}(\theta_g) - R_n/2}, \quad R_n = (\widehat{\theta} - \theta_g) \nabla^3 \overline{\ell}(\theta^*). \tag{3}$$

□ Now

$$n^{1/2}\nabla \overline{\ell}(\theta_g) = n^{-1/2} \sum_{j=1}^n \nabla \log f(Y_j; \theta_g)$$

has mean (vector) zero and variance (matrix)

$$\operatorname{var}\left\{n^{-1/2}\sum_{j=1}^{n}\nabla\log f(Y_{j};\theta_{g})\right\} = n^{-1}\sum_{j=1}^{n}\operatorname{E}_{g}\left\{\nabla\log f(Y_{j};\theta_{g})\nabla^{\mathrm{T}}\log f(Y_{j};\theta_{g})\right\} = \hbar_{1}(\theta_{g}).$$

so the numerator of (3) converges in distribution to $\mathcal{N}\{0, \hbar_1(\theta_q)\}$, using the CLT.

☐ Moreover the weak law of large numbers gives

$$-\nabla^2 \overline{\ell}(\theta_g) = -\frac{1}{n} \sum_{j=1}^n \nabla^2 \log f(Y_j; \theta_g) \xrightarrow{P} i_1(\theta_g).$$

- \square Lemma 51 shows that $R_n \stackrel{P}{\longrightarrow} 0$, so the denominator of (3) tends in probability to $i_1(\theta_g)$.
- \square Putting the pieces together, we find that

$$n^{1/2}(\widehat{\theta} - \theta_g) \xrightarrow{D} \mathcal{N}_d\{0, \imath_1(\theta_g)^{-1}\hbar_1(\theta_g)\imath_1(\theta_g)^{-1}\}, \quad n \to \infty,$$

where the variance formula is also valid when i_1 and \hbar_1 are $d \times d$ matrices.

 \Box The information quantities based on a random sample of size n are $\imath(\theta_g)=n\imath_1(\theta_g)$ and $\hbar(\theta_g)=n\hbar_1(\theta_g),$ giving

$$\widehat{\theta} \sim \mathcal{N}_d(\theta_g, \imath(\theta_g)^{-1} \hbar(\theta_g) \imath(\theta_g)^{-1} \},$$

in which the variance is of the usual order 1/n.

stat.epfl.ch

Note: A useful lemma

Lemma 51 Under the conditions of Theorem 50, $R_n = (\widehat{\theta} - \theta_g) \nabla^3 \overline{\ell}(\theta^*) \stackrel{P}{\longrightarrow} 0$ as $n \to \infty$.

 $\Box \quad \text{For } \varepsilon > 0, \ B_n = \{|R_n| > \varepsilon\}, \ A_n = \{|\widehat{\theta} - \theta_g| > \delta\} \ \text{and} \ \delta > 0 \ \text{small enough that} \ \mathcal{N} \ \text{contains a ball of radius} \ \delta \ \text{around} \ \theta_g, \ \text{we have}$

$$P(|R_n| > \varepsilon) = P(B_n \cap A_n) + P(B_n \cap A_n^c) \le P(A_n) + P(B_n \cap A_n^c),$$

where the first term tends to zero because the sequence $\widehat{\theta}$ is consistent.

 \Box If $|\widehat{\theta} - \theta_q| < \delta$, then (C3) implies that

$$|R_n| \le \delta n^{-1} \sum_{j=1}^n |\partial^3 \log f(Y_j; \theta^*) / \partial \theta^3| \le \delta n^{-1} \sum_{j=1}^n m(Y_j) = \delta \overline{M}_n,$$

say, and clearly $\overline{M}_n \stackrel{P}{\longrightarrow} M$, say. Therefore

$$P(B_n \cap A_n^c) = P(B_n \cap |\widehat{\theta} - \theta_g| > \delta) \le P(B_n \cap |R_n| \le \delta \overline{M}_n)$$

and for $\eta > 0$ this equals

$$P(B_n \cap |R_n| \le \delta \overline{M}_n \cap \overline{M}_n \le M + \eta) + P(B_n \cap |R_n| \le \delta \overline{M}_n \cap \overline{M}_n > M + \eta),$$

which is bounded by

$$P\{|R_n| > \varepsilon \cap |R_n| \le \delta(M+\eta)\} + P(|\overline{M}_n - M| > \eta).$$

The last term here tends to zero, because $\overline{M}_n \stackrel{P}{\longrightarrow} M$, and the first can be made equal to zero by choosing δ such that $\delta(M+\eta) < \varepsilon$. This proves the lemma.

stat.epfl.ch

Autumn 2024 - note 3 of slide 99

Classical asymptotics

- \square The true model is supposed to lie in the candidate family, i.e., $g \in \mathcal{F}$, so $\theta_g \in \Theta$.
- We saw on slide 38 that the moments of the $d \times 1$ score vector $U(\theta) = \nabla \ell(\theta)$ are given under mild conditions by the Bartlett identities, i.e.,

$$\mathrm{E}\{U(\theta)\} = 0, \quad \mathrm{var}\{U(\theta)\} = \mathrm{E}\left\{\nabla \ell(\theta) \nabla^{\mathrm{T}} \ell(\theta)\right\} = \mathrm{E}\left\{-\nabla^2 \ell(\theta)\right\}, \quad \dots$$

- $\Box \quad \text{Hence } \imath(\theta) = \hbar(\theta) \text{, and } \imath(\theta) = n \imath_1(\theta) = n \hbar_1(\theta) \text{ when } Y_1, \dots, Y_n \overset{\text{iid}}{\sim} f_{\theta_g}.$
- \square Mathematically speaking the assumption that $g \in \mathcal{F}$ is always false, but
 - the asymptotic results are supposed to provide guidelines on what to expect when fitting models checking the regularity conditions in practice would require knowledge of g, in which case there's no need for inference!
 - this is largely irrelevant if model-checking suggests that $f_{\widehat{\theta}_g}$ is 'close enough' to g.
- \square Crucially, the interest parameter ψ should have a stable interpretation for candidates likely to be close to g (i.e., within $n^{-1/2}$), so $\mathcal F$ is 'robustly specified' if the model is not quite right, then the interpretation of the crucial parameters will be unchanged.

stat.epfl.ch

Note: Stable interpretation of a parameter

To put some mathematical flesh on the discussion, suppose that $g(y) = f(y; \theta, \gamma)$ and the assumed model is $f(y; \theta, 0)$. Then for small γ , $\theta_q \equiv \theta_\gamma$ satisfies

$$0 = \int \nabla_{\theta} \log f(y; \theta_{\gamma}, 0) f(y; \theta, \gamma) dy$$

$$= \int \left\{ \nabla_{\theta} \log f(y; \theta, \gamma) + \nabla_{\theta}^{2} \log f(y; \theta, \gamma) (\theta_{\gamma} - \theta) + \nabla_{\gamma}^{T} \nabla_{\theta} \log f(y; \theta, \gamma) (0 - \gamma) + \cdots \right\} f(y; \theta, \gamma) dy$$

$$= 0 - i \theta_{\theta}(\theta, \gamma) (\theta_{\gamma} - \theta) + i \theta_{\gamma}(\theta, \gamma) \gamma + o(\gamma),$$

which implies that the effect of incorrectly assuming that $\gamma=0$ is that $\widehat{\theta}$ converges to

$$\theta_{\gamma} = \theta + i_{\theta\theta}^{-1}(\theta, \gamma) i_{\theta\gamma}(\theta, \gamma) \gamma + o(\gamma).$$

It is also easy to check that $\hbar_{\theta\theta}(\theta,0)=\imath_{\theta\theta}(\theta,0)+O(\gamma)$, so the two matrices become equal if $\gamma\to 0$, in which case $\imath_1(\theta,\gamma)^{-1}\hbar_1(\theta,\gamma)\imath_1(\theta,\gamma)^{-1}\to\imath_1(\theta,\gamma)^{-1}$, which implies that for small γ we have

$$n^{1/2}(\widehat{\theta} - \theta) = n^{1/2}(\widehat{\theta} - \theta_{\gamma}) + n^{1/2}(\theta_{\gamma} - \theta) \sim \mathcal{N}_{d}\{0, \imath_{1}(\theta, \gamma)^{-1}\} + n^{1/2}(\theta_{\gamma} - \theta).$$

 $\hfill \square$ Now if $\gamma=n^{-a}\delta$ for some a>0 , then

$$n^{1/2}(\theta_{\gamma} - \theta) = n^{1/2 - a} i_{\theta\theta}^{-1}(\theta, \gamma) i_{\theta\gamma}(\theta, \gamma) \delta,$$

which will tend to infinity if a<1/2 (should be obvious asymptotically), to zero if a>1/2 (can be ignored asymptotically) and to a constant if a=1/2. Hence there is an asymptotic bias for $\widehat{\theta}$ if there is misspecification, $\delta \neq 0$, unless $\imath_{\theta\gamma}(\theta,\gamma)=0$, i.e., the information matrix covariance for the scores for θ and γ is zero. This is known as orthogonality of θ and γ ; see later.

stat.epfl.ch

Autumn 2024 - note 1 of slide 100

In practice ...

 \square We usually assume classical asymptotics and replace the sandwich matrix $i(\theta_g)^{-1}\hbar(\theta_g)i(\theta_g)^{-1}$ by the inverse of the observed information matrix

$$\widehat{\jmath} = -\nabla^2 \ell(\widehat{\theta}),$$

which

- can be computed numerically without (possibly awkward) expectations,
- will (helpfully!) misbehave if the maximisation is questionable,
- has been found to give generally good results in applications,
- has the heuristic justification that $(\widehat{\theta},\widehat{\jmath})$ are approximately sufficient for θ_g , as

$$\ell(\theta_g) \doteq \ell(\widehat{\theta}) - \frac{1}{2}(\widehat{\theta} - \theta_g)^{\mathrm{T}} \widehat{\jmath}(\widehat{\theta} - \theta_g).$$

- $\hfill\Box$ Standard errors for $\widehat{\theta}$ are the square roots of the diagonal elements of $\widehat{\jmath}^{-1}.$
- \square If we <u>must</u> make the sandwich we can replace $\imath(\theta_g)$ by $\widehat{\jmath}$ and $\hbar(\theta_g)$ by (e.g.)

$$\widehat{h} = \sum_{j=1}^{n} \nabla \log f(Y_j; \widehat{\theta}) \nabla^{\mathrm{T}} \log f(Y_j; \widehat{\theta}),$$

though $\widehat{\jmath}^{-1}\widehat{h}\,\widehat{\jmath}^{-1}$ can be unstable because $\widehat{\hbar}$ misbehaves.

stat.epfl.ch

Autumn 2024 - slide 101

Related statistics

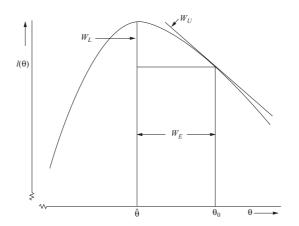


Figure 6.2. Three asymptotically equivalent ways, all based on the log likelihood function of testing null hypothesis $\theta = \theta_0$: W_E , horizontal distance; W_L vertical distance; W_U slope at null point.

From Cox (2006, Principles of Statistical Inference)

stat.epfl.ch

Related statistics

 \Box Classical asymptotics support inference for scalar θ based on any of the (approximate) pivots

$$\begin{split} T &= t(\theta_g) = \widehat{\jmath}^{1/2}(\widehat{\theta} - \theta_g) \stackrel{.}{\sim} \mathcal{N}(0,1), & \text{Wald statistic,} \\ S &= s(\theta_g) = \widehat{\jmath}^{-1/2}U(\theta_g) \stackrel{.}{\sim} \mathcal{N}(0,1), & \text{score statistic,} \\ W &= w(\theta_g) = 2\{\ell(\widehat{\theta}) - \ell(\theta_g)\} \stackrel{.}{\sim} \chi_1^2, & \text{likelihood ratio statistic,} \\ R &= r(\theta_g) = \mathrm{sign}(\widehat{\theta} - \theta_g)w(\theta_g)^{1/2} \stackrel{.}{\sim} \mathcal{N}(0,1), & \text{likelihood root.} \end{split}$$

The likelihood root has other names (e.g., directed likelihood ratio statistic).

 \square The distribution of W follows from the expansion on the previous slide.

 \Box If $\widehat{ heta}^{
m o}$ and $\jmath(\widehat{ heta}^{
m o})$ have been obtained for observed data $y^{
m o}$, then the approximation

$$P_q\{T(\theta_q) \le t^{o}(\theta_q)\} \doteq \Phi\{t^{o}(\theta_q)\}$$

leads to $(1-\alpha)$ Wald confidence interval $\widehat{\theta}^{\rm o}\pm\jmath(\widehat{\theta}^{\rm o})^{-1/2}z_{1-\alpha/2}$ based on T, while that based on W is

$$\{\theta: W^{\circ}(\theta) \le \chi_1^2(1-\alpha)\} = \{\theta: \ell^{\circ}(\theta) \ge \ell^{\circ}(\widehat{\theta}^{\circ}) - \frac{1}{2}\chi_1^2(1-\alpha)\},$$

where z_p and $\chi^2_{\nu}(p)$ are respectively the p quantiles of the N(0,1) and χ^2_{ν} distributions.

stat.epfl.ch Autumn 2024 – slide 103

Comparative comments

 \square Confidence intervals based on T are symmetric, but those based on W or R take the shape of ℓ into account and are parametrisation-invariant;

in small samples the distributional approximations for W and R are better than that for T, and that for W can be improved by Bartlett correction, using $W_B = W/(1 + b/n)$;

 \square confidence sets based on W may not be connected (and if so T or R are unreliable);

 \square the main use of S is for testing in situations where maximisation of ℓ is awkward, and then $\widehat{\jmath}$ is often replaced by $\imath(\theta_q)$;

 \square a variant of R, the modified likelihood root

$$R^* = r^*(\theta_g) = r(\theta_g) + \frac{1}{r(\theta_g)} \log \frac{q(\theta_g)}{r(\theta_g)},$$

often gives almost perfect inferences even in small samples (more later ...).

Example 52 Compute the above statistics when $y_1, \ldots, y_n \stackrel{\text{iid}}{\sim} \exp(\theta)$ and compare the resulting inferences with those from an exact pivot.

stat.epfl.ch

- □ The log likelihood is $\ell(\theta) = n(\log \theta \theta \overline{y})$, for $\theta > 0$, which is clearly unimodal with $\widehat{\theta} = 1/\overline{y}$ and $\jmath(\theta) = n/\theta^2$.
- ☐ Hence

 $t(\theta) = n^{1/2}(1 - \theta \,\overline{y}),$

 $s(\theta) = n^{1/2} \{ 1/(\theta \, \overline{y}) - 1 \},$

 $w(\theta) = 2n \{\theta \overline{y} - \log(\theta \overline{y}) - 1\},$

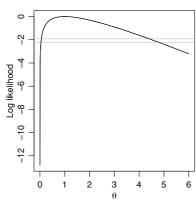
 $r(\theta) = \operatorname{sign}(1 - \theta \, \overline{y}) \left[2n \left\{ \theta \, \overline{y} - \log(\theta \, \overline{y}) - 1 \right\} \right]^{1/2}.$

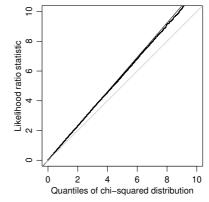
- \square The exact pivot is $\theta \sum Y_j$ whose distribution is gamma with unit scale and shape parameter n.
- Consider an exponential sample with n=1 and $\overline{y}=1$; then $\widehat{\jmath}=1$. The log likelihood $\ell(\theta)$, shown in the left-hand panel of the figure, is unimodal but strikingly asymmetric, suggesting that confidence intervals based on an approximating normal distribution for $\widehat{\theta}$ will be poor. The right-hand panel is a chi-squared probability plot in which the ordered values of simulated $w(\theta)$ are graphed against quantiles of the χ_1^2 distribution—if the simulations lay along the diagonal line x=y, then this distribution would be a perfect fit. The simulations do follow a straight line rather closely, but with slope $(1+b/n)\chi_1^2$, where b=0.1544. This indicates that the distribution of the Bartlett-adjusted likelihood ratio statistic $w(\theta)/(1+b/n)$ would be essentially χ_1^2 . The 95% confidence intervals for θ based on the unadjusted and adjusted likelihood ratio statistics are (0.058, 4.403) and (0.042, 4.782) respectively.

stat.epfl.ch

Autumn 2024 - note 1 of slide 104

Exponential example

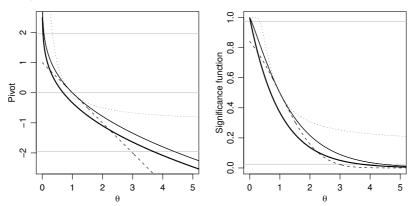




Likelihood inference for exponential sample of size n=1. Left: log likelihood $\ell(\theta)$. Intersection of the function with the two horizontal lines gives two 95% confidence intervals for θ : the upper line is based on the χ_1^2 approximation to the distribution of $w(\theta)$, and the lower line is based on the Bartlett-corrected statistic. Right: comparison of simulated values of likelihood ratio statistic $w(\theta)$ with χ_1^2 quantiles. The χ_1^2 approximation is shown by the line of unit slope, while the $(1+b/n)\chi_1^2$ approximation is shown by the upper straight line.

stat.epfl.ch

Exponential example



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 confidence functions, with horizontal lines at 0.025 and 0.975.

stat.epfl.ch Autumn 2024 – slide 106

Non-regular models

- \Box The regularity conditions (C1)–(C4) apply in many settings met in practice, but not universally. The most common failures arise when
 - some of the parameters are discrete (e.g., change point problems),
 - the model is not identifiable (distinct θ values give the same model),
 - θ_g is on the boundary of the parameter space (e.g., testing for a zero variance),
 - $d = \dim(\theta)$ grows (too fast) with n, or
 - the support of $f(y;\theta)$ depends on θ (so the Bartlett identities fail).
- □ Even when the conditions are satisfied there can be datasets for which maximum likelihood estimation fails, e.g.,
 - there is no unique maximum to the likelihood, or
 - the maximum is on the edge of the parameter space,
 and then penalisation (equivalent to using a prior) is often used.

Example 53 If $Y_1, \ldots, Y_n \stackrel{\text{iid}}{\sim} U(0, \theta)$, show that the limit distribution of $n(\theta - \widehat{\theta})/\theta$ when $n \to \infty$ is $\exp(1)$. Discuss.

 \Box In this case $1=\int f(y;\theta)\,\mathrm{d}y=\int_0^\theta\theta^{-1}\,\mathrm{d}y$, and differentiation with respect to θ gives

$$0 = 1/\theta + \int_0^\theta (-\theta^{-2}) \,\mathrm{d}y,$$

so the first Bartlett identity is not satisfied (because the support depends on θ , and $f(\theta;\theta) \neq 0$.

 \square Owing to the independence,

$$L(\theta) = \prod_{j=1}^{n} f_Y(y_j; \theta) = \prod_{j=1}^{n} \left\{ \theta^{-1} I(0 < y_j < \theta) \right\} = \theta^{-n} I(\max y_j < \theta), \quad \theta > 0,$$

and therefore $\widehat{\theta} = M = \max Y_i$, whose distribution is

$$P(M \le x) = (x/\theta)^n, \quad 0 < x < \theta.$$

Now

$$P\left\{n(\theta - \widehat{\theta})/\theta \le x\right\} = P(\widehat{\theta} \ge \theta - x\theta/n) = 1 - \{(\theta - x\theta/n)/\theta\}^n \to 1 - \exp(-x),$$

as required. Note that:

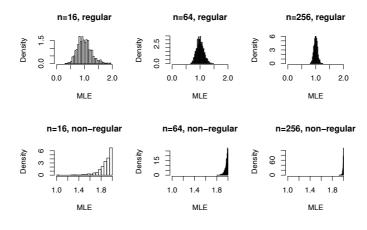
- the scaling needed to get a limiting distribution is much faster here than in the regular case (we have to multiply by n to get a non-degenerate limit);
- the limit is not normal.

stat.epfl.ch

Autumn 2024 - note 1 of slide 107

Uniform example

Comparison of the distributions of $\widehat{\theta}$ in a regular case (panels above, with standard deviation $\propto n^{-1/2}$) and in a nonregular case (Example 53, panels below, with standard deviation $\propto n^{-1}$). In other nonregular cases it might happen that the distribution is nasty (unlike here) and/or that the convergence is slower than in regular cases.



stat.epfl.ch

Vector case

 \Box When θ is a vector and under classical asymptotics we base inference on the distributional approximations

$$\widehat{\theta} \stackrel{\cdot}{\sim} \mathcal{N}_d(\theta_g, \widehat{\jmath}^{-1}), \quad w(\theta_g) = 2 \left\{ \ell(\widehat{\theta}) - \ell(\theta_g) \right\} \stackrel{\cdot}{\sim} \chi_d^2, \quad s(\theta_g) = \widehat{\jmath}^{-1/2} U(\theta_g) \stackrel{\cdot}{\sim} \mathcal{N}_d(0, I_d),$$

with

- the first very commonly used for inferences on parameters;
- the second used to test whether $\theta = \theta_q$;
- the third much less used than the others, generally in the form $s(\theta_q)^{\mathrm{T}} s(\theta_q) \stackrel{\cdot}{\sim} \chi_d^2$.
- \square If θ divides into a $p \times 1$ interest parameter ψ and a $q \times 1$ nuisance parameter λ , then

$$\widehat{\theta} = \begin{pmatrix} \widehat{\psi} \\ \widehat{\lambda} \end{pmatrix} \stackrel{\cdot}{\sim} \mathcal{N}_{p+q} \left\{ \begin{pmatrix} \psi_g \\ \lambda_g \end{pmatrix}, \begin{pmatrix} \widehat{\jmath}_{\psi\psi} & \widehat{\jmath}_{\psi\lambda} \\ \widehat{\jmath}_{\lambda\psi} & \widehat{\jmath}_{\lambda\lambda} \end{pmatrix}^{-1} \right\},$$

where for brevity we now write $\widehat{\lambda}_{\psi} = \max_{\lambda} \ell(\psi, \lambda)$, $\widetilde{\theta} = \widehat{\theta}_{\psi} = (\psi, \widehat{\lambda}_{\psi})$,

$$\ell_{\psi} = \left. \frac{\partial \ell(\theta)}{\partial \psi} \right|_{\theta = \theta_g}, \quad \widehat{\jmath}_{\psi\psi} = -\widehat{\ell}_{\psi\psi} = -\left. \frac{\partial^2 \ell(\theta)}{\partial \psi \partial \psi^{\mathrm{T}}} \right|_{\theta = \widehat{\theta}}, \quad \widetilde{\ell}_{\psi\psi} = \left. \frac{\partial^2 \ell(\theta)}{\partial \psi \partial \psi^{\mathrm{T}}} \right|_{\theta = \widetilde{\theta}}, \quad \text{etc.}$$

stat.epfl.ch

Autumn 2024 - slide 110

Inference on ψ

 \Box Under classical asymptotics and setting $\hat{\jmath}^{\psi\psi} = (\hat{\jmath}_{\psi\psi} - \hat{\jmath}_{\psi\lambda}\hat{\jmath}_{\lambda\lambda}^{-1}\hat{\jmath}_{\lambda\psi})^{-1}$ we have

$$\begin{split} \widehat{\psi} \stackrel{.}{\sim} \mathcal{N}_p \left(\psi_g, \widehat{\jmath}^{\psi\psi} \right) & \text{maximum likelihood estimator,} \\ s(\psi_g) &= \widetilde{\ell}_\psi^{\scriptscriptstyle T} \widetilde{\jmath}^{\psi\psi} \widetilde{\ell}_\psi \stackrel{.}{\sim} \chi_p^2 & \text{score statistic,} \\ w_{\rm p}(\psi_g) &= 2 \left\{ \ell_{\rm p}(\widehat{\psi}) - \ell_{\rm p}(\psi_g) \right\} \stackrel{.}{\sim} \chi_p^2 & \text{(generalized) likelihood ratio statistic,} \end{split}$$

where we defined w_p using the profile log likelihood $\ell_p(\psi) = \ell(\psi, \widehat{\lambda}_{\psi}) = \max_{\lambda} \ell(\psi, \lambda)$.

 \Box If ψ is scalar (p=1, the usual situation), the likelihood root is defined as

$$r(\psi_g) = \operatorname{sign}\left(\widehat{\psi} - \psi_g\right) \sqrt{w(\psi_g)} \stackrel{\cdot}{\sim} \mathcal{N}(0, 1).$$

- ☐ Properties:
 - inferences using $w(\psi_g)$ and $r(\psi_g)$ are invariant to interest-respecting reparametrisation, so are preferable but more computationally burdensome;
 - $s(\psi_g)$ is mainly used for tests, since only λ must be estimated (as $\psi = \psi_g$ is known).
- \Box A $(1-\alpha)$ confidence set based on $w_{\rm p}(\psi_g)$ (or equivalently on $\ell_{\rm p}(\psi))$ is

$$\left\{\psi: w_{\mathbf{p}}(\psi) \leq \chi_p^2(1-\alpha)\right\} = \left\{\psi: \ell(\psi, \widehat{\lambda}_{\psi}) \geq \ell(\widehat{\psi}, \widehat{\lambda}) - \frac{1}{2}\chi_p^2(1-\alpha)\right\}.$$

stat.epfl.ch

Note: Large-sample distribution of the likelihood ratio statistic $w_{\rm p}(\psi_q)$

☐ We write

$$w_{\mathbf{p}}(\psi_q) = 2\{\ell(\widehat{\theta}) - \ell(\widehat{\theta}_{\psi})\} = 2\{\ell(\widehat{\theta}) - \ell(\theta_q)\} - 2\{\ell(\widehat{\theta}_{\psi}) - \ell(\theta_q)\}$$

and use Taylor series to approximate both terms by quadratic forms in $\widehat{\theta}-\theta_g$ and $\widehat{\lambda}_\psi-\lambda_g$.

- \square We shall need to express ℓ_{θ} , ℓ_{λ} and $\widehat{\lambda}_{\psi} \lambda_{q}$ in terms of $\widehat{\theta} \theta_{q}$. Taylor expansion gives

$$0 = \widehat{\ell}_{\theta} = \ell_{\theta} + \ell_{\theta\theta}(\widehat{\theta} - \theta_q) + \dots = \ell_{\theta} - \imath_{\theta\theta}(\widehat{\theta} - \theta_q) + \dots,$$

where \cdots denotes terms of smaller order containing third derivatives of ℓ . The λ component of this equation is

$$0 = \ell_{\lambda} - i_{\lambda\psi}(\widehat{\psi} - \psi_q) - i_{\lambda\lambda}(\widehat{\lambda} - \lambda_q) + \cdots$$

Likewise

$$0 = \tilde{\ell}_{\lambda} = \ell_{\lambda} + \ell_{\lambda\lambda}(\widehat{\lambda}_{\psi} - \lambda_g) + \dots = \ell_{\lambda} - \imath_{\lambda\lambda}(\widehat{\lambda}_{\psi} - \lambda_g) + \dots$$

Equating the expressions for ℓ_{λ} from the last two displays gives

$$\ell_{\lambda} \doteq \iota_{\lambda\psi}(\widehat{\psi} - \psi_g) + \iota_{\lambda\lambda}(\widehat{\lambda} - \lambda_g) = \iota_{\lambda\lambda}(\widehat{\lambda}_{\psi} - \lambda_g),$$

SO

$$\ell_{\theta} \doteq \imath_{\theta\theta}(\widehat{\theta} - \theta_q), \quad \ell_{\lambda} \doteq \imath_{\lambda\lambda}(\widehat{\lambda}_{\psi} - \lambda_q), \quad \widehat{\lambda}_{\psi} - \lambda_q \doteq \widehat{\lambda} - \lambda_q + \imath_{\lambda\lambda}^{-1} \imath_{\lambda\psi}(\widehat{\psi} - \psi_q).$$

☐ To obtain the quadratic forms we write

$$\ell(\widehat{\theta}) = \ell(\theta_g) + (\widehat{\theta} - \theta_g)^{\mathrm{T}} \ell_{\theta} + \frac{1}{2} (\widehat{\theta} - \theta_g)^{\mathrm{T}} \ell_{\theta\theta} (\widehat{\theta} - \theta_g) + \cdots$$

$$\stackrel{\cdot}{=} \ell(\theta_g) + (\widehat{\theta} - \theta_g)^{\mathrm{T}} \imath_{\theta\theta} (\widehat{\theta} - \theta_g) - \frac{1}{2} (\widehat{\theta} - \theta_g)^{\mathrm{T}} \imath_{\theta\theta} (\widehat{\theta} - \theta_g),$$

resulting in

$$2\{\ell(\widehat{\theta}) - \ell(\theta_g)\} \stackrel{:}{=} (\widehat{\theta} - \theta_g)^{\mathrm{T}} \imath_{\theta\theta} (\widehat{\theta} - \theta_g)$$

$$= (\widehat{\psi} - \psi_g)^{\mathrm{T}} \imath_{\psi\psi} (\widehat{\psi} - \psi_g) + 2(\widehat{\psi} - \psi_g)^{\mathrm{T}} \imath_{\psi\lambda} (\widehat{\lambda} - \lambda_g) + (\widehat{\lambda} - \lambda_g)^{\mathrm{T}} \imath_{\lambda\lambda} (\widehat{\lambda} - \lambda_g),$$

and likewise

$$2\{\ell(\widehat{\theta}_{\psi}) - \ell(\theta_{g})\} \stackrel{:}{=} (\widehat{\lambda}_{\psi} - \lambda_{g})^{\mathrm{T}} \imath_{\lambda\lambda} (\widehat{\lambda}_{\psi} - \lambda_{g})$$

$$\stackrel{:}{=} \{(\widehat{\lambda} - \lambda_{g}) + \imath_{\lambda\lambda}^{-1} \imath_{\lambda\psi} (\widehat{\psi} - \psi_{g})\}^{\mathrm{T}} \imath_{\lambda\lambda} \{(\widehat{\lambda} - \lambda_{g}) + \imath_{\lambda\lambda}^{-1} \imath_{\lambda\psi} (\widehat{\psi} - \psi_{g})\}$$

$$= (\widehat{\psi} - \psi_{g})^{\mathrm{T}} \imath_{\psi\lambda} \imath_{\lambda\lambda}^{-1} \imath_{\lambda\psi} (\widehat{\psi} - \psi_{g}) + 2(\widehat{\psi} - \psi_{g})^{\mathrm{T}} \imath_{\psi\lambda} (\widehat{\lambda} - \lambda_{g}) + (\widehat{\lambda} - \lambda_{g})^{\mathrm{T}} \imath_{\lambda\lambda} (\widehat{\lambda} - \lambda_{g}).$$

Subtracting the two quadratic forms gives

$$w_{\mathbf{p}}(\psi_g) = 2\{\ell(\widehat{\theta}) - \ell(\theta_g)\} - 2\{\ell(\widehat{\theta}_{\psi}) - \ell(\theta_g)\}$$

$$\doteq (\widehat{\psi} - \psi_g)^{\mathrm{T}} (\imath_{\psi\psi} - \imath_{\psi\lambda} \imath_{\lambda\lambda}^{-1} \imath_{\lambda\psi}) (\widehat{\psi} - \psi_g),$$

and as $\widehat{\psi} \sim \mathcal{N}\{\psi_g, (\imath_{\psi\psi} - \imath_{\psi\lambda}\imath_{\lambda\lambda}^{-1}\imath_{\lambda\psi})^{-1}\}$, we see that $w_p(\psi_g) \sim \chi_p^2$, as claimed.

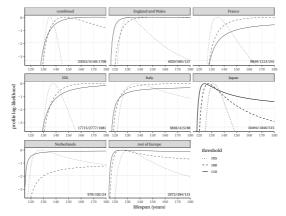
 \square Here we are under classical asymptotics, whereby the dimensions of ψ and λ are fixed and $n \to \infty$, and arguments along the lines of Theorem 50 show that the terms \cdots all tend in probability to zero, and thus do not affect the limiting distribution.

stat.epfl.ch

Autumn 2024 - note 1 of slide 111

Example: Human lifespan

Example 54 Profile log likelihoods for the endpoint ψ of a generalized Pareto model fitted to data on lifetimes of persons aged over 105 from different databases, with thresholds at 105, 108, 110 years. Here λ is scalar, so p=q=1, and the horizontal line at $-\frac{1}{2}\chi_1^2(0.95)=-1.92$ indicates 95% confidence regions.



From Belzile et al. (2022, Annual Review of Statistics and its Application).

stat.epfl.ch Autumn 2024 – slide 112

Model selection

☐ The fact that

$$KL(g, f) = E_g\{\log g(Y) - \log f(Y)\} = E_g\left[-\log\left\{\frac{f(Y)}{g(Y)}\right\}\right] \ge 0$$

is minimised when f=g suggested comparing competing models $\mathcal{F}_1,\ldots,\mathcal{F}_M$ by their maximised log likelihoods $\log f_m(y;\widehat{\theta}_m)=\widehat{\ell}_m$.

- \square But $\widehat{\ell}_m$ should be penalized, because
 - $\widehat{\ell}_m \geq \log f_m(y; \theta_m)$ even if \mathcal{F}_m is the true model class, and
 - enlarging θ_m will increase $\widehat{\ell}_m$ even if further parameters are unnecessary.
- Akaike proposed minimising $2E_g E_g^+ \left[-\log\{f(Y^+; \widehat{\theta})/g(Y^+)\} \right]$, where $Y^+, Y \stackrel{\text{iid}}{\sim} g$ are independent datasets. The idea is that if $\widehat{\theta} = \widehat{\theta}(Y)$ is estimated separately from Y^+ , there will be a penalty due to 'missing θ_g ' which will grow with $\dim(\theta)$ (picture . . .)
- \square This leads to choosing m to minimise the Akaike or the network information criteria

$$AIC_m = 2\left(d_m - \widehat{\ell}_m\right), \quad NIC_m = 2\left\{tr(\widehat{h}_m\widehat{\jmath}_m^{-1}) - \widehat{\ell}_m\right\},$$

where the first takes $\operatorname{tr}(\widehat{h}_m \widehat{\jmath}_m^{-1}) \approx d_m = \dim(\theta_m)$.

Note: Derivation of AIC/NIC

☐ As

$$2\mathbf{E}_g\mathbf{E}_g^+ \left[-\log\{f(Y^+;\widehat{\theta})/g(Y^+)\} \right] = 2\mathbf{E}_g^+ \left\{ \log g(Y^+) \right\} - 2\mathbf{E}_g\mathbf{E}_g^+ \left\{ \log f(Y^+;\widehat{\theta}) \right\},$$

we can ignore the first term in the minimisation over f. An unbiased estimator of the second term would be $-2\ell^+(\widehat{\theta})$, where ℓ^+ is the log likelihood based on Y^+ and $\widehat{\theta}$ is based on Y, but the estimator we have available is $-2\ell(\widehat{\theta})$, in which the log likelihood and $\widehat{\theta}$ are both based on Y. Clearly $\ell(\widehat{\theta})$ is upwardly biased, but by how much?

☐ To find out we consider the Taylor expansion

$$2\ell^{+}(\widehat{\theta}) = 2\ell^{+}(\widehat{\theta}^{+}) + 2(\widehat{\theta} - \widehat{\theta}^{+})^{\mathrm{T}}\ell_{\theta}^{+}(\widehat{\theta}^{+}) + (\widehat{\theta} - \widehat{\theta}^{+})^{\mathrm{T}}\ell_{\theta\theta}^{+}(\widehat{\theta}^{+})(\widehat{\theta} - \widehat{\theta}^{+}) + \cdots$$

$$= 2\ell^{+}(\widehat{\theta}^{+}) - \mathrm{tr}\left\{(\widehat{\theta} - \widehat{\theta}^{+})^{\mathrm{T}}\imath_{\theta\theta}(\theta_{g})(\widehat{\theta} - \widehat{\theta}^{+})\right\} + \cdots$$

$$= 2\ell^{+}(\widehat{\theta}^{+}) - \mathrm{tr}\left\{(\widehat{\theta} - \widehat{\theta}^{+})(\widehat{\theta} - \widehat{\theta}^{+})^{\mathrm{T}}\imath_{\theta\theta}(\theta_{g})\right\} + \cdots$$

where $\widehat{\theta}^+$ maximises $\ell^+(\theta)$, $\widehat{\theta}$ maximises $\ell(\theta)$, we have replaced $-\ell_{\theta\theta}^+(\widehat{\theta}^+)$ by its large-sample limit $\imath_{\theta\theta}(\theta_g)$ and neglected terms that are $o_p(1)$. Recall that θ_g is the large-sample limit of $\widehat{\theta}$ when data are sampled from g.

 $\square \quad \text{Now } \widehat{\theta}^+ \text{ and } \widehat{\theta} \text{ are independent and approximately } \mathcal{N}_d(\theta_g, V) \text{, where } V = \imath_{\theta\theta}^{-1}(\theta_g) \hbar(\theta_g) \imath_{\theta\theta}^{-1}(\theta_g) \text{, so } \widehat{\theta}^+ - \widehat{\theta} \stackrel{.}{\sim} \mathcal{N}_d(0, 2V) \text{, giving }$

$$-2\mathbf{E}_{g}\mathbf{E}_{g}^{+}\left\{\ell^{+}(\widehat{\theta})\right\} \stackrel{:}{=} -2\mathbf{E}_{g}\mathbf{E}_{g}^{+}\left\{\ell(\widehat{\theta})\right\} + \operatorname{tr}\left\{2V\imath_{\theta\theta}(\theta_{g})\right\} + o(1)$$

$$= 2\left[\operatorname{tr}\left\{\hbar(\theta_{g})\imath_{\theta\theta}^{-1}(\theta_{g})\right\} - \mathbf{E}_{g}\mathbf{E}_{g}^{+}\left\{\ell(\widehat{\theta})\right\}\right] + o(1).$$

- $\Box \quad \text{If } \hbar(\theta_g) \doteq \imath_{\theta\theta}(\theta_g) \text{, then this final expression can be estimated by } \mathrm{AIC} = 2\{d \ell(\widehat{\theta})\}, \text{ where } d = \dim(\theta), \text{ or by the } network \textit{ information criterion } \mathrm{NIC} = 2\{\mathrm{tr}(\widehat{h}\widehat{\jmath}^{-1}) \ell(\widehat{\theta})\}.$
- \square Neither AIC or NIC gives consistent selection of the true model, which would require the penalty to grow with n.
- ☐ The calculations above use generic large-sample likelihood approximations, and can be improved in specific cases (e.g., with normal errors).

stat.epfl.ch

Autumn 2024 - note 1 of slide 113

3.3 Nuisance Parameters

slide 114

Effect of nuisance parameters

Example 55 (Neyman–Scott) Find the profile log likelihood for σ^2 when $(y_{j1}, y_{j2}) \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu_j, \sigma^2)$, for $j = 1, \ldots, n$. Comment.

- Profiling over many nuisance parameters can lead to completely wrong inferences, as the previous example shows.
- \square Even when the number of nuisance parameters is o(n) we may run into trouble: in general

Bias
$$(\widehat{\psi}; \psi) = O(d^3/n)$$
,

so for the bias to tend to zero in large samples we require $d=o(n^{1/3})$ for consistency of $\widehat{\psi}$. Hence bias increases with $\dim(\lambda)$, at least in general.

☐ How can we rescue 'ordinary' likelihood inference when there are many nuisance parameters?

☐ The overall log likelihood is

$$\ell(\sigma^2, \mu_1, \dots, \mu_n) \equiv -\frac{1}{2} \left[(2n) \log \sigma^2 + \frac{1}{\sigma^2} \sum_{j=1}^n \left\{ (y_{j1} - \mu_j)^2 + (y_{j2} - \mu_j)^2 \right\} \right],$$

and differentiation with respect to μ_i gives that $\widehat{\mu}_i = (y_{i1} + y_{i2})/2$, so as

$${a - (a + b)/2}^2 + {b - (a + b)/2}^2 = (a - b)^2/2$$

we obtain

$$\ell_{\rm p}(\sigma^2) = -n\log\sigma^2 - \frac{1}{4\sigma^2} \sum_{j=1}^n (y_{j1} - y_{j2})^2, \quad \sigma^2 > 0.$$

- This is maximised at $\widehat{\sigma}_p^2 = (4n)^{-1} \sum_{j=1}^n (y_{j1} y_{j2})^2$, but as $Y_{j1} Y_{j2} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 2\sigma^2)$, we see that $\sigma_p^2 \stackrel{P}{\longrightarrow} \sigma^2/2$ as $n \to \infty$; this is a completely inconsistent estimator. Hence the profile log likelihood has its asymptotic maximum in completely the wrong place.
- \square In this example there are d=n+1 parameters of which n are nuisance parameters.

stat.epfl.ch

Autumn 2024 - note 1 of slide 115

Dealing with nuisance parameters

- \square Approaches to dealing with high-dimensional λ include:
 - basing inference on a marginal likelihood or a conditional likelihood,

$$f(y; \psi, \lambda) = f(w; \psi) \times f(y \mid w; \psi, \lambda) = f(y \mid w_{\psi}; \psi) \times f(w_{\psi}; \psi, \lambda),$$

where w_{ψ} may not depend on ψ (recall Lemmas 39 and 40) — OK for any configuration of λ s, but may lose information on ψ ;

- constructing a partial likelihood (like the above, but harder to build);
- higher-order inference, via, e.g., a modified profile likelihood or a modified likelihood root, which can approximate both conditional and marginal likelihoods;
- using orthogonal parameters, i.e., mapping $\lambda \mapsto \zeta(\lambda, \psi)$ which is orthogonal to ψ ;
- using a composite likelihood in which λ does not appear; or
- taking $\lambda \sim h(\cdot)$ and using the integrated likelihood $\int f(y;\psi,\lambda)h(\lambda)\,\mathrm{d}\lambda$ depends on h, like Bayesian inference.
- □ We have already seen examples of marginal and conditional likelihoods.
- ☐ Below we sketch some of the other approaches.

stat.epfl.ch

Modified profile likelihood

 \square Replace profile log likelihood $\ell_p(\psi)$ by the modified profile log likelihood

$$\ell_{\rm mp}(\psi) = \ell_{\rm p}(\psi) + m(\psi),$$

with $m(\psi)$ chosen to make $\ell_{\rm p}$ closer to a marginal or conditional log likelihood.

□ Taking

$$m(\psi) = -\frac{1}{2} \log \left| \jmath_{\lambda\lambda}(\psi, \widehat{\lambda}_{\psi}) \right| + \log \left| \frac{\partial \widehat{\lambda}}{\partial \widehat{\lambda}_{\psi}^{\mathrm{T}}} \right|$$

does this in some generality.

- □ The
 - first term of $m(\psi)$ can be obtained numerically if need be, but
 - the second term, a Jacobian needed to make $\ell_{\rm mp}$ invariant to interest-preserving reparametrisation, is hard to compute in general.
- \square Simpler to base a likelihood on the normal distribution of the modified likelihood root $r^*(\psi)$ (next).

stat.epfl.ch Autumn 2024 – slide 117

Higher-order inference . . .

 \square Classical theory gives first-order accuracy, i.e., with ψ scalar

$$P\{r(\psi_q) \le r^{o}(\psi_q)\} = \Phi\{r^{o}(\psi)\} + O(n^{-1/2}),$$

so tests and one-sided confidence sets

$$\{\psi: r^{\mathrm{o}}(\psi) \le z_{1-\alpha}\}$$

based on the observed data y^{o} have error $n^{-1/2}$.

 \square If we replace $r(\psi)$ by the modified likelihood root,

$$r^*(\psi) = r(\psi) + \frac{1}{r(\psi)} \log \left\{ \frac{q(\psi)}{r(\psi)} \right\},$$

where $q(\psi)$ depends on the model, then for continuous responses the error drops to $O(n^{-3/2})$, so

$$P\{r^*(\psi_a) \le r^{*o}(\psi_a)\} = \Phi\{r^{*o}(\psi_a)\} + O(n^{-3/2}),$$

so a one-sided confidence set

$$\{\psi: r^{*o}(\psi) \le z_{1-\alpha}\}$$

has error of order $n^{-3/2}$; often this almost exact even for tiny n (recall Example 52).

stat.epfl.ch

... with nuisance parameters

 \square With nuisance parameters, $r(\psi) = \operatorname{sign}(\widehat{\psi} - \psi) \sqrt{w_{p}(\psi)}$, and

$$q(\psi) = \frac{|\varphi(\widehat{\theta}) - \varphi(\widehat{\theta}_{\psi}) | \varphi_{\lambda}(\widehat{\theta}_{\psi})|}{|\varphi_{\theta}(\widehat{\theta})|} \left\{ \frac{|\widehat{\jmath}|}{|\jmath_{\lambda\lambda}(\widehat{\theta}_{\psi})|} \right\}^{1/2}$$

where φ is the $d \times 1$ canonical parameter of a local exponential family approximation to the model at the observed data y^{o} , with $\varphi_{\theta}(\theta) = \partial \varphi(\theta)/\partial \theta^{\mathrm{T}}$, etc.

 \Box In a general exponential family $\varphi(\theta)$ is the canonical parameter, and in a linear exponential family,

$$q(\psi) = (\widehat{\psi} - \psi) \left\{ \frac{|\widehat{\jmath}|}{|\jmath_{\lambda\lambda}(\widehat{\theta}_{\psi})|} \right\}^{1/2}.$$

☐ In general for independent continuous observations we write

$$\varphi(\theta)_{d \times 1} = V_{d \times n}^{\mathrm{T}} \left. \frac{\partial \ell(\theta; y)}{\partial y} \right|_{y = y^{\mathrm{o}}} = \sum_{j=1}^{n} V_{j}^{\mathrm{T}} \left. \frac{\partial \log f(y_{j}; \psi, \lambda)}{\partial y_{j}} \right|_{y = y^{\mathrm{o}}},$$

where the $1 \times d$ vectors $V_j = \partial y_j/\partial \theta^{\mathrm{T}}$ are evaluated at y^{o} and $\widehat{\theta}^{\mathrm{o}}$.

stat.epfl.ch Autumn 2024 – slide 119

Properties of higher order approximations

- ☐ Invariant to interest-respecting reparameterization.
- ☐ Computation almost as easy as first order versions.
- \square Error $O(n^{-3/2})$ in continuous response models, $O(n^{-1})$ in discrete response models.
- ☐ Relative (not absolute) error, so highly accurate in tails.
- ☐ Bayesian version is also available (and easier to derive).

Example 56 (Location-scale model) Compute $\varphi(\theta)$ for a location-scale model, in which independent observations Y_j have density $\tau^{-1}h\{(y-\eta)/\tau\}$. What about the normal density?

☐ In this case the overall log likelihood is

$$\ell(\eta, \tau) = -n \log \tau + \sum_{j=1}^{n} \log h\{(y_j - \eta)/\tau\},$$

so the vector $\partial \ell(\eta, \tau)/\partial y$ has components $\tau^{-1}(\log h)'\{(y_j - \eta)/\tau\}$, evaluated at the parameters η and τ and observed data vector $y_1^{\rm o}, \ldots, y_n^{\rm o}$.

- To compute the V_j we use the structural expression $y=\eta+\tau\varepsilon$, where $\varepsilon\sim h$. This represents y as a function of $\theta^{\rm T}=(\eta,\tau)$, and yields $\partial y_j/\partial\theta^{\rm T}=(1,\varepsilon_j)$. This has to be evaluated at the observed data point $y^{\rm o}$, and at that point the parameters are replaced by their maximum likelihood estimates, giving $V_i^{\rm T}=(1,(y_j^{\rm o}-\widehat{\eta}^{\rm o})/\widehat{\tau}^{\rm o})$.
- ☐ This yields

$$\varphi(\theta) = \sum_{j=1}^{n} \tau^{-1} (\log h)' \{ (y_j^{o} - \eta)/\tau \} (1, \varepsilon_j^{o})^{\mathrm{T}},$$

where we have set $\varepsilon_{j}^{\mathrm{o}}=(y_{j}^{\mathrm{o}}-\widehat{\eta}^{\mathrm{o}})/\widehat{\tau}^{\mathrm{o}}.$

 $\Box \quad \text{If h is normal, then } \log h(u) \equiv -u^2/2 \text{, so } (\log h)'\{(y_j^{\rm o} - \eta)/\tau\} = -(y_j^{\rm o} - \eta)/\tau^2 \text{, leading to }$

$$\varphi(\theta)^{\mathrm{T}} = \left(\sum_{j=1}^{n} (\eta - y_{j}^{\mathrm{o}}) / \tau^{2}, \sum_{j=1}^{n} (\eta - y_{j}^{\mathrm{o}}) / \tau^{2} \times e_{j}\right) \equiv (\eta / \tau^{2}, 1 / \tau^{2}),$$

because it turns out that inferences are invariant under non-singular affine transformations of $\varphi(\theta)$ (exercise).

stat.epfl.ch

Autumn 2024 - note 1 of slide 120

Orthogonal parameters

- If the expected information matrix is block diagonal, with $\imath_{\psi,\lambda}(\theta)=0$ for all θ , then $\widehat{\psi}$ is asymptotically independent of $\widehat{\lambda}$, and we can hope that the effect on $\widehat{\psi}$ of estimating λ will be limited. If so, we say that ψ and λ are orthogonal.
- \Box To see the effect of this, we expand the equation defining $\widehat{\lambda}_{\psi}$ around $\widehat{\theta}$, giving

$$0 = \frac{\partial \ell(\widehat{\theta}_{\psi})}{\partial \lambda} = \frac{\partial \ell(\widehat{\theta})}{\partial \lambda} + \frac{\partial^{2} \ell(\widehat{\theta})}{\partial \lambda \partial \theta^{T}} (\widehat{\theta}_{\psi} - \widehat{\theta}) + \cdots$$
$$= \frac{\partial^{2} \ell(\widehat{\theta})}{\partial \lambda \partial \lambda^{T}} (\widehat{\lambda}_{\psi} - \widehat{\lambda}) + \frac{\partial^{2} \ell(\widehat{\theta})}{\partial \lambda \partial \psi^{T}} (\psi - \widehat{\psi}) + \cdots$$
$$= \widehat{\jmath}_{\lambda \lambda} (\widehat{\lambda}_{\psi} - \widehat{\lambda}) + \widehat{\jmath}_{\lambda \psi} (\psi - \widehat{\psi}) + \cdots$$

which implies that

$$\widehat{\lambda}_{\psi} = \widehat{\lambda} + \widehat{\jmath}_{\lambda\lambda}^{-1} \widehat{\jmath}_{\lambda\psi} (\widehat{\psi} - \psi) + \cdots$$

- \square Hence if we can arrange the model so that $\widehat{\jmath}_{\lambda\psi}\approx 0$, for example by parametrising it so that $\imath_{\lambda\psi}(\theta)\equiv 0$, then $\widehat{\lambda}_{\psi}$ will depend only weakly on ψ , and we can ignore the Jacobian term in the modified profile likelihood.
- \square This suggests mapping an original parametrisation (ψ, γ) to (ψ, λ) , where $\lambda = \lambda(\psi, \gamma)$ is orthogonal to ψ .

stat.epfl.ch

Orthogonalisation

 \square Writing $\gamma = \gamma(\psi, \lambda)$ gives

$$\ell(\psi, \lambda) = \ell^* \{ \psi, \gamma(\psi, \lambda) \},$$

and differentiation with respect to ψ and λ leads to

$$\frac{\partial^2 \ell}{\partial \lambda \partial \psi} = \frac{\partial \gamma^{\mathrm{T}}}{\partial \lambda} \frac{\partial^2 \ell^*}{\partial \gamma \partial \psi} + \frac{\partial \gamma^{\mathrm{T}}}{\partial \lambda} \frac{\partial^2 \ell^*}{\partial \gamma \partial \gamma^{\mathrm{T}}} \frac{\partial \gamma}{\partial \psi} + \frac{\partial^2 \gamma^{\mathrm{T}}}{\partial \lambda \partial \psi} \frac{\partial \ell^*}{\partial \gamma}.$$

☐ For orthogonality this must have expectation zero, so

$$0 = \frac{\partial \gamma^{\mathrm{T}}}{\partial \lambda} i_{\gamma\psi}^* + \frac{\partial \gamma^{\mathrm{T}}}{\partial \lambda} i_{\gamma\gamma}^* \frac{\partial \gamma}{\partial \psi},$$

where $\imath_{\gamma\psi}^*$ and $\imath_{\gamma\gamma}^*$ are components of the expected information matrix in the non-orthogonal parametrization, so λ solves the system of q PDEs

$$\frac{\partial \gamma}{\partial \psi} = -i_{\gamma\gamma}^{*-1}(\psi, \gamma)i_{\gamma\psi}^*(\psi, \gamma).$$

 \square In fact an explicit expression for λ in terms of ψ and γ is not needed to compute ℓ_{mp} in the new parametrisation.

stat.epfl.ch

Autumn 2024 - slide 122

Orthogonal parametrisation

 \square A solution (possibly numerical) always exists when $\dim(\psi)=1$, but need not exist when ψ is vector, because then we must simultaneously solve

$$\frac{\partial \gamma}{\partial \psi_1} = -i_{\gamma\gamma}^{*-1}(\psi, \gamma) i_{\gamma\psi_1}^*(\psi, \gamma), \quad \frac{\partial \gamma}{\partial \psi_2} = -i_{\gamma\gamma}^{*-1}(\psi, \gamma) i_{\gamma\psi_2}^*(\psi, \gamma),$$

for all γ , ψ_1 and ψ_2 , but the compatibility condition

$$\frac{\partial^2 \gamma}{\partial \psi_1 \partial \psi_2} = \frac{\partial^2 \gamma}{\partial \psi_2 \partial \psi_1}$$

may fail.

Example 57 (Linear exponential family) What parameter is orthogonal to ψ in the linear exponential family with log likelihood

$$\ell^*(\psi,\gamma) \equiv s_1^{\mathrm{T}}\psi + s_2^{\mathrm{T}}\gamma - k(\psi,\gamma)?$$

Consider normal and Poisson likelihoods in particular.

stat.epfl.ch

 \square The parameters $\lambda = \lambda(\psi, \gamma)$ orthogonal to ψ are determined by

$$\frac{\partial \gamma}{\partial \psi^{\mathrm{T}}} = -k_{\gamma\gamma}^{-1}(\psi, \gamma)k_{\gamma\psi}(\psi, \gamma). \tag{4}$$

If we reparametrize in terms of ψ and $\lambda = k_{\gamma}(\psi, \gamma) = \partial k(\psi, \gamma)/\partial \gamma$, then in this new parametrization, γ is a function of ψ and λ , and

$$0 = \frac{\partial \lambda^{\mathrm{T}}}{\partial \psi} = \frac{\partial \gamma^{\mathrm{T}}}{\partial \psi} k_{\gamma\gamma}(\psi, \gamma) + k_{\psi\gamma}(\psi, \gamma),$$

so $\lambda=k_{\gamma}(\psi,\gamma)$ is a solution to (4). That is, the parameter orthogonal to ψ is the so-called complementary mean parameter $\lambda(\psi,\gamma)=\mathrm{E}(S_2;\psi,\gamma)$. By symmetry, $\mathrm{E}(S_1;\psi,\gamma)$ is orthogonal to γ .

- □ The normal distribution with mean μ and variance σ^2 has canonical parameter $(\mu/\sigma^2, -1/(2\sigma^2))$. The canonical statistic (Y, Y^2) has expectation $(\mu, \mu^2 + \sigma^2)$, so μ is orthogonal to $-1/(2\sigma^2)$, and hence to σ^2 , while μ/σ^2 is orthogonal to $\mu^2 + \sigma^2$.
- \square Independent Poisson variables Y_1 and Y_2 with means $\exp(\gamma)$ and $\exp(\gamma + \psi)$ have log likelihood

$$\ell^*(\psi,\gamma) \equiv (y_1 + y_2)\gamma + y_2\psi - e^{\gamma} - e^{\gamma + \psi}.$$

The discussion above suggests that

$$\lambda = E(Y_1 + Y_2) = \exp(\gamma) + \exp(\gamma + \psi) = e^{\gamma}(1 + e^{\psi})$$

is orthogonal to ψ , so $\gamma = \log \lambda - \log(1 + e^{\psi})$ and

$$\ell(\psi, \lambda) \equiv y_2 \psi - (y_1 + y_2) \log(1 + e^{\psi}) + (y_1 + y_2) \log \lambda - \lambda.$$

The separation of ψ and λ implies that the profile and modified profile likelihoods for ψ are proportional. They correspond to the conditional likelihood obtained from the density of Y_2 given Y_1+Y_2 .

stat.epfl.ch

Autumn 2024 - note 1 of slide 123

Composite likelihood

Used when full likelihood can't be computed but densities for distinct subsets of the observations, y_{S_1}, \ldots, y_{S_C} , are available, can use a composite (log) likelihood

$$\ell_{\mathrm{C}}(\theta) = \sum_{c=1}^{C} \log f(y_{\mathcal{S}_c}; \theta).$$

 \square The choice of subsets S_1, \ldots, S_C determines what parameters can be estimated.

☐ Special cases:

- independence likelihood takes $S_j = \{y_j\}$ and treats (possibly dependent) y_j as independent;
- pairwise likelihood uses subsets of distinct pairs $\{y_j, y_{j'}\}$.
- ☐ May be useful with spatial data, and then contributions from distant pairs may be downweighted or dropped entirely.
- $\[\] \] \ell_C(\theta)$ satisfies the first Bartlett identity, so can give consistent estimators $\tilde{\theta}$, but requires a sandwich variance matrix (or some other approach) to estimate $var(\tilde{\theta})$.
- ☐ Model comparisons use the composite likelihood information criterion

$$\mathrm{CLIC} = 2 \left[\mathrm{tr} \{ \hbar(\tilde{\theta}) \jmath(\tilde{\theta})^{-1} \} - \ell_{\mathrm{C}}(\tilde{\theta}) \right].$$

stat.epfl.ch Autumn 2024 – slide 124

Comments

☐ Other likelihoods and/or likelihood-like functions are widely used, especially

- partial likelihood, used to eliminate nuisance functions for inference (survival data),
- quasi-likelihood, used to model over-dispersion in exponential family models,
- pseudo-likelihood, treats data as Gaussian even when they are not (econometrics), and
- empirical likelihood, an extension of nonparametric modelling (econometrics).

☐ Strengths of likelihood approach:

- heuristic as plausibility of a model as explanation of data;
- we 'just' have to write down the density of the observed data;
- invariance to data and parameter transformations;
- general (and 'optimal') approximate theory for inference in regular models;
- close links to Bayesian inference.

☐ Weaknesses of likelihood approach:

- requires 'parametric' model for data;
- can fail in high-dimensional settings;
- not all models are regular.