

Topics in Likelihood Inference

STA4508H

Nancy Reid
University of Toronto

February 2, 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

Presentations

Feb 9 Shiki: Solomon and Cox (1992)

Feb 16 Angela:
Robert: Barndorff-Nielsen and Cox (1979)

Feb 23 Hengchao:
Siyue: De Stavola and Cox (2008)
Manuel: Battey and Cox (2018)
Ziang: Cox (1975) ← Partial likelihood

Feb 9 and Feb 16 in SS 2120; ←
Feb 23 online

Recap: nuisance parameters

profile, marginal, conditional, saddlepoint, Laplace

$$\theta = (\psi, \lambda) \quad \psi \in \mathbb{R}$$

$$Y_1, \dots, Y_n \sim f(y_i; \psi, \lambda)$$

$$f(y_i | x_i; \psi, \lambda)$$

$$L(\psi, \lambda) \propto \prod_{i=1}^n f(y_i; \psi, \lambda)$$

$$L_p(\psi) = \prod_{i=1}^n f(y_i; \psi, \hat{\lambda}_\psi)$$

$$\frac{\partial \log L(\psi, \lambda)}{\partial \lambda} \Bigg|_{\lambda = \hat{\lambda}_\psi} = 0$$

$$l_\lambda(\psi, \hat{\lambda}_\psi) = 0$$

$$l_p(\psi) = \log l_p(\psi)$$

$$(\hat{\psi} - \psi) \{ -l_p''(\hat{\psi}) \}^{1/2} \checkmark$$

if dim λ is large relative to n ,

$$\rightarrow \pm \sqrt{2 \{ l_p(\hat{\psi}) - l_p(\psi) \}} \checkmark$$

$$l_p'(\psi) \{ j_p(\hat{\psi}) \}^{-1/2} \checkmark$$

$$E_\theta l_p'(\psi) \neq 0 = O\left(\frac{1}{n}\right) (?)$$

$$E_\theta l_p''(\psi) + \text{var}[l_p(\psi)]$$

Recap: nuisance parameters

Bartlett ident. $\neq 0$

profile, marginal, conditional, saddlepoint, Laplace

1) we might have $f(y; \psi, \lambda) \propto f_c(s(y) | t(y); \psi) \underbrace{f(t(y); \psi, \lambda)}_{\text{cond'l lik for } \psi}$

2) " " $f(y; \psi, \lambda) \propto f_c(s(y) | t(y); \psi, \lambda) \underbrace{f_m(t(y); \psi)}_{\text{marg'l lik}}$

$$\psi \rightarrow (\text{suff stat.}) \rightarrow (s, t)$$

Since $L_c(\psi) \propto f_c(s | t; \psi)$

$L_m(\psi) \propto f_m(t; \psi)$

"genuine" lik. $\int f_s$ b/c constructed from density for observable

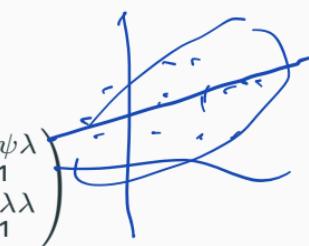
... Nuisance parameters

- partition score vector: $U(\theta) = \begin{pmatrix} U_\psi(\theta) \\ U_\lambda(\theta) \end{pmatrix}$; $\frac{1}{\sqrt{n}} U_\psi(\theta) \xrightarrow{d} N_q\{0, i_1^{\psi\psi}(\theta)\}$

$$\mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \star & \star \\ \star & \star \end{pmatrix}\right)$$

$$L(\psi, \lambda) = L_1(\psi) L_2(\lambda)$$

- partition information matrix: $i_1(\theta) = \begin{pmatrix} i_{\psi\psi} & i_{\psi\lambda} \\ i_{\lambda\psi} & i_{\lambda\lambda} \end{pmatrix}$ $i_1^{-1}(\theta) = \begin{pmatrix} i_1^{\psi\psi} & i_1^{\psi\lambda} \\ i_1^{\lambda\psi} & i_1^{\lambda\lambda} \end{pmatrix}$



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

$$\sqrt{n}(\hat{\psi} - \psi) \doteq \frac{1}{\sqrt{n}}(i_1^{\psi\psi})^{-1}(U_\psi - i_{\psi\lambda} i_{\lambda\lambda}^{-1} U_\lambda)$$

$$2\{\ell_p(\hat{\psi}) - \ell_p(\psi)\} \doteq (\hat{\psi} - \psi)^T i^{\psi\psi} (\hat{\psi} - \psi)$$

$$\boxed{\sqrt{n}(\hat{\theta} - \theta) \doteq \frac{1}{\sqrt{n}} i_1^{-1}(\theta) U(\theta)}$$

MLE
Score

projection of U_ψ on space spanned by U_λ
 $\sqrt{n}(\hat{\psi} - \psi) \xrightarrow{d} N_q\{0, i_1^{\psi\psi}(\theta)\}$

$$(X, Y)^\top \sim \mathcal{N}_2 \left[\begin{pmatrix} \mu_X \\ \mu_Y \end{pmatrix}, \begin{pmatrix} \Sigma_{XX} & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_{YY} \end{pmatrix} \right]$$

$\Sigma_{XX} \Sigma_{XX}^{-1} \triangleq \beta$
 $Y|X \sim N\left(\mu_Y - \Sigma_{YX} \Sigma_{XX}^{-1} (X - \mu_X), \Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}\right)$

Linear exponential families

$\hat{\beta}$

$$\sum_{i=1}^n y_i x_i / \sum_{i=1}^n (x_i^2)$$

- conditional density free of nuisance parameter

$$f(y_i; \psi, \lambda) = \exp\{\psi^T s(y_i) + \lambda^T t(y_i) - k(\psi, \lambda)\} h(y_i)$$

$$f(y; \psi, \lambda) = \exp\{\psi^T \sum s(y_i) + \lambda^T \sum t(y_i) - nk(\psi, \lambda)\} \prod h(y_i)$$

$$\psi^T = (\beta_1, \beta_2)$$

$$\lambda = (\beta_3, \dots, \beta_p)$$

Let $s = \sum s(y_i), t = \sum t(y_i)$] ← closed → suff' + st.

$$f(s, t; \psi, \lambda) = \exp\{\psi^T s + \lambda^T t - nk(\psi, \lambda)\} \tilde{h}(s)$$

$$\begin{aligned} f(s | t; \psi) &= \frac{f(s, t; \psi, \lambda)}{\int f(s, t; \psi, \lambda) ds} \\ &= \frac{\exp\{\psi^T s + \lambda^T t - nk(\psi, \lambda)\} \tilde{h}(s)}{\int \exp\{\psi^T s + \lambda^T t - nk(\psi, \lambda)\} \tilde{h}(s) ds} \\ &= \frac{\exp\{\psi^T s\} \tilde{h}(s)}{\int \exp\{\psi^T s\} \tilde{h}(s) ds} \\ &= \exp\{\psi^T s - n \tilde{k}_t(\psi)\} \tilde{h}_t(s) \end{aligned}$$

$$e^{\psi s} \tilde{h}(s) / \int e^{\psi s} \tilde{h}(s) ds$$

$$= e^{\ell_t(\hat{\psi}) - \ell_t(\psi)}$$

suff.

✓ \tilde{k}_t, \tilde{h}_t just convenient notation for integral of denominator

$$\log f(s|t; \psi) = l_c(\psi)$$

Bayesian v.

$$\pi_m(\psi|y) = \int \pi(\psi, \lambda|\psi) d\lambda$$

↑
Lapl approx.
= ... =

$$\propto e^{\ell_p(\psi) - \ell_p(\hat{\psi})} \quad \downarrow \quad \{$$

$$\log \pi_m(\psi) = \ell_p(\psi) - \frac{1}{2} \underbrace{\sum_{\lambda} (\log |j_{\lambda\lambda}(\psi, \hat{\lambda}_{\psi})|)}_{(II)} + \underbrace{\log \pi(\psi)}_{(I)}$$

$$l_c(\psi) \doteq \ell(\psi, \hat{\lambda}_{\psi}) + \underbrace{\frac{1}{2} \log |j_{\lambda\lambda}(\psi, \hat{\lambda})|}_{\text{using similar analysis to Laplace}}$$

$$f_m(t; \psi) f_c(s|t; \psi, \lambda)$$

Approximate conditional and marginal inference

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

$$e^{\psi s + x^T t - c(\psi, x)} h(s, t)$$

~~$i_{\lambda\lambda}(\theta) = 0$~~

Bayesian π -aft prior

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

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

(ψ, η) , with $i_{\psi\eta}(\theta) \equiv 0$
orthog. param

adjusted profile log-likelihood

$$\ell_A(\psi) = \ell_p(\psi) + A(\psi)$$

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

- 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

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

SM §12.4.1

makes LHS inv.

Semi-parametric models

to review

- Recall: y_1, \dots, y_n jumps of a Poisson process

- rate function $\lambda(\cdot)$ observed on $(0, \tau)$
- events at $0 < y_1 < \dots < y_n < \tau$
- likelihood function

$$\underbrace{L\{\lambda(\cdot); y\}}_{f} = \left\{ \prod_{i=1}^n \lambda(y_i) \right\} \exp\left\{- \int_0^\tau \lambda(u) du\right\} = \prod \lambda(y_i) e^{-\lambda(\tau)}$$

- log-likelihood function

$$\ell\{\lambda(\cdot); y\} = \sum_{i=1}^n \log \lambda(y_i) - \int_0^\tau \lambda(u) du$$

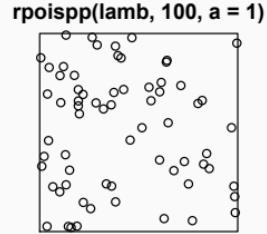
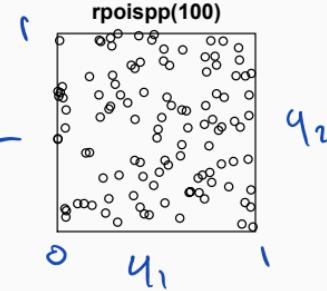
- in space:

$$\ell\{\lambda(\cdot); y\} = \sum_{i=1}^n \log \lambda(y_i) - \int_S \lambda(u) du$$

SM §6.5.1

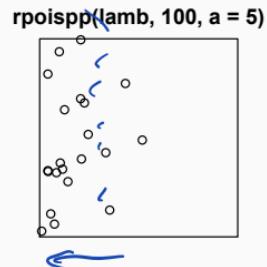
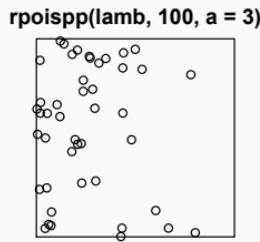


$$\alpha = 0$$



spatstat

$rpoispp$



e.g.

$$x(t) = e^{\beta_0 + \beta_1 t}$$

rate per unit area

$$\lambda(y_1, y_2) = 100 \exp(-ay_1)$$

Survival data

- Example: Survival data $(y_i, d_i), i = 1, \dots, n$

$$y_i = \min(y_j^0, c_i)$$

$$d_i = \mathbf{1}\{y_i = y_i^0\} \quad c_i > y_i^0$$

$$f(y_i, d_i; \theta) = [f(y_i; \theta)\{1 - G(y_i)\}]^{d_i} [\{1 - F(y_i; \theta)\}g(y_i)]^{1-d_i}$$

$y_i^0 \sim F(\cdot; \theta)$, $c_i \sim G$; y_i^0 independent of c_i
uncensored observation

joint density

$$\ell(\theta) = \sum_{i=1}^n [d_i \log f(y_i; \theta) + (1 - d_i) \log \{1 - F(y_i; \theta)\}]$$

+ terms depending on G

$$= \sum \{d_i \log \lambda(y_i; \theta) - \Lambda(y_i; \theta)\}$$

$$\boxed{\Lambda(y; \theta)} = -\log \{1 - F(y; \theta)\}; \quad \lambda(y; \theta) = f(y; \theta) / \{1 - F(y; \theta)\}$$

hazard

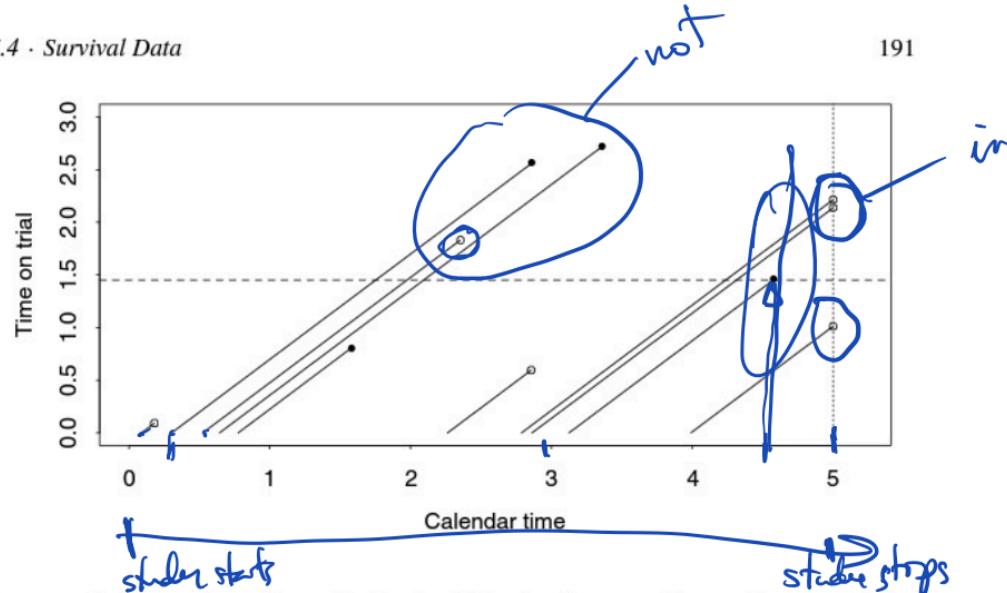
$\Lambda(t; \theta)$
 $= \int_0^t \pi(u; \theta) du$
Cumulative hazard

Survival data

$$d_i \log f(y_i; \theta) - \log\{-F(y_i; \theta)\}$$

$$d_i \log \frac{f}{1-F} \triangleq d_i \lambda(y_i; \theta)$$

Figure 5.8 Lexis diagram showing typical pattern of censoring in a medical study. Each individual is shown as a line whose x coordinates run from the calendar time of entry to the trial to the calendar time of failure (blob) or censoring (circle). Censoring occurs at the end of the trial, marked by the vertical dotted line, or earlier. The vertical axis shows time on trial, which starts when individuals enter the study. The risk set for the failure at calendar time 4.5 comprises those individuals whose lines touch the horizontal dashed line; see page 543.



thus we study events on the vertical axis. Calendar time may be used to account for changes in medical practice over the course of a trial.

In applications the assumption that C_j and Y_j^0 are independent is critical. There would be serious bias if the illlest patients drop out of a trial because the treatment makes them feel even worse, thereby inducing association between survival and censoring variables because patients die soon after they withdraw.

The examples above all involve *right-censoring*. Less common is left-censoring, where the time of origin is not known exactly, for example if time to death from a disease is observed, but the time of infection is unknown.

In practice a high proportion of the data may be censored, and there may be a serious loss of efficiency if this is ignored (Example 4.20). These will also be discussed.

the risk set
at t = 5

0+	1+	1+	3+	3+	7	10+	11+	12+	12+	15+	18+
20+	22+	22+	24+	25+	26+	31+	36+	36+	36	38	40
47+	47+	49+	53+	53+	55+	56+	57+	61+	67+	67+	70
73	75+	77+	83+	84+	88+	89+	99	121+	122+	123+	141+
0+	0+	2+	2+	2+	2+	3	3+	4+	5+	9+	10+
11	12+	13	13+	18+	22+	22+	24+	24+	24+	25+	26+
27	28	32+	35+	36	40+	43+	50+	54			

Table 5.3

Blalock-Taussig shunt data (Oakes, 1991). The table gives survival time of shunt (months after operation) for 48 infants aged over one month at time of operation, followed by times for 33 infants aged 30 or fewer days at operation. Infants whose shunt has not yet failed are marked +.

$$\ell(\theta; y_i, d_i) = \sum_{i=1}^n \{d_i \log \lambda(y_i; \theta) - \Lambda(y_i; \theta)\}$$


Proportional hazards regression

- semi-parametric model: $\lambda(y_i; x_i, \beta) = \lambda_0(y_i) \exp(x_i^T \beta)$
- log-likelihood function

$$\ell(\beta, \lambda_0; y, d) = \sum_{i=1}^n d_i \log\{\lambda(y_i; x_i, \beta)\} - \Lambda(y_i; x_i, \beta)$$

~~$\lambda_0(y_i)$~~

$$= \sum_{i=1}^n [d_i \{x_i^T \beta + \log \lambda_0(y_i)\} - \Lambda(y_i) \exp(x_i^T \beta)] = \ell(\beta, \lambda_0; y, d)$$

- partial log-likelihood function

$$\ell_{part}(\beta; y, d) = \sum_{i=1}^n d_i \{x_i^T \beta - \log \sum_{j \in \mathcal{R}_i} \exp(x_j^T \beta)\}$$

- $y_1 < \dots < y_n$; $\mathcal{R}_i = \{j; y_j \geq y_i\}$

$\uparrow \uparrow \uparrow$

partial

1972 cond'l

... PH regression

$$\ell_{part}(\beta; y, d) = \sum_{i=1}^n d_i \{x_i^T \beta - \log \left(\sum_{j \in \mathcal{R}_i} \exp(x_j^T \beta) \right)\}$$

$$= \sum_{i=1}^n d_i \{x_i^T \beta - \underline{\log A_i(\beta)}\}$$

y_1, \dots, y_n

$U(\beta)$

$j(\beta)$

$\text{var } U \text{ est }$

$$\frac{\partial \ell_{part}(\beta)}{\partial \beta} = \sum_{i=1}^n d_i \left\{ x_i - \frac{A'_i(\beta)}{A_i(\beta)} \right\}$$

$$-\frac{\partial^2 \ell_{part}(\beta)}{\partial \beta \partial \beta^T} = \sum_{i=1}^n d_i \left\{ \frac{A''_i(\beta)}{A_i(\beta)} - \frac{A'_i(\beta) A'_i(\beta)^T}{A_i(\beta)^2} \right\}$$

$$A''_i = \frac{\partial^2 A_i}{\partial \beta \partial \beta^T}$$

$\approx 10^{-15} \text{ years}$

$$A'_i(\beta) = \frac{\partial A_i}{\partial \beta}$$

$Cov = \text{var}$

$U(\beta) j^{-1} h(\hat{\beta})$

notation is a bit careless
 $\sim N(0, I)$

... PH regression

- partial log-likelihood function

$$\ell_{part}(\beta; y, d) = \sum_{i=1}^n d_i \{x_i^T \beta - \log \sum_{j \in \mathcal{R}_i} \exp(x_j^T \beta)\}$$

- can be motivated as:

1. marginal log-likelihood of the **ranks** of the failure times

f_r

f_n

y_1, y_2, \dots
 x_1, x_2, \dots

Pf2
2. $\prod_{i=1}^n \Pr(\text{unit } i \text{ fails at } y_i \mid \text{history to } y_i^-, \text{ one failure from } \mathcal{R}_i)$

CL
not a cond'l
density

(discussi: an

Prentice)

2000
3. profile log-likelihood function Murphy - van der Vaart

- for inference, $\ell_{part}(\beta)$ has usual properties

1. $\hat{\beta}_{part} \sim N\{\beta, j_{part}^{-1}(\hat{\beta})\},$

$2\{\ell_{part}(\hat{\beta}_{part}) - \ell_{part}(\beta)\} \sim \chi_d^2$

Davison §10.8; Cox 1972, 1975

... semi-parametric models

(SM Ch. 10-8)

- partial log-likelihood function

$$\lambda(t; \beta) = \lambda_0(t) e^{\frac{x^T \beta}{\lambda_0(t)}}$$

$$\ell_{part}(\beta; y, d) = \sum_{i=1}^n d_i \left\{ x_i^T \beta - \log \sum_{j \in R_i} \exp(x_j^T \beta) \right\}$$

- is also, 3. profile log-likelihood function if $\lambda(\cdot)$ is represented by a vector of values

$$(\lambda_1, \dots, \lambda_n) = \{\lambda(y_1), \dots, \lambda(y_n)\}$$

\uparrow \uparrow $\dim(x) \uparrow \text{w.r.t. } ;$

- why does usual likelihood inference apply?

(y_i : failure only?
censoring?)

- can be connected to theory of empirical likelihood \leftarrow

Murphy & van der Waart, 2000; van der Waart 1998, Ch. 25

- $\ell(\beta, \lambda; y), \beta \in \mathbb{R}^d; \lambda = \lambda(\cdot)$

- $\ell_p(\beta; y) = \ell(\beta, \tilde{\lambda}_\beta; y); \quad \tilde{\lambda}_\beta = \arg \sup_\lambda \ell(\beta, \lambda; y)$

- example: failure times y with hazard $\lambda(y | x) = e^{x\beta} \lambda(y)$

β

$$f(y_i; \beta, \lambda) = e^{x_i \beta} \lambda(y_i) \exp\{-e^{x_i \beta} \Lambda(y_i)\}$$

- empirical likelihood:

$f(y_i; \beta, \Lambda(\cdot))$

$$EL(\beta, \Lambda; y) = \prod_{i=1}^n e^{x_i \beta} \Lambda\{y_i\} \exp\{-e^{x_i \beta} \Lambda(y_i)\}$$

- maximizing value of $\Lambda(\cdot)$ must have jumps at y_i only – replace $\Lambda(y_i)$ by sum

PH model, no censoring

$\Lambda = \int \lambda$

$\Lambda\{y_i\} = \lambda(y_i)$

$\lambda_1, \dots, \lambda_n$

$\sum_{j \geq i} \lambda_j$

... semi-parametric models

- empirical likelihood:

$$EL(\beta, \Lambda; y) = \prod_{i=1}^n e^{x_i \beta} \Lambda\{y_i\} \exp\{-e^{x_i \beta} \Lambda(y_i)\}$$

$$\hat{\Lambda}_\beta\{y_i\} = \left\{ \sum_{j:y_j \geq y_i} \exp(x_j \beta) \right\}^{-1}$$

constrained rule

- profile log-likelihood

$$L_p(\beta) = \prod_{i=1}^n \frac{e^{x_i \beta}}{\sum_{j:y_j \geq y_i} \exp(x_j \beta)}$$

]

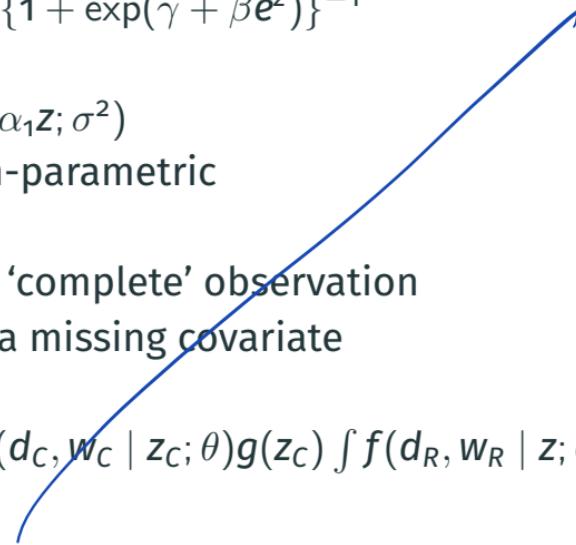
profile lik.

- same as partial likelihood motivated by different arguments

↑ at moment no censoring

Murphy and ✓

- observation (D, W, Z) ; D and W are independent, given Z
- $\Pr(D = 0) = \{1 + \exp(\gamma + \beta e^z)\}^{-1}$
- $W \sim N(\alpha_0 + \alpha_1 z; \sigma^2)$
- $Z \sim g(\cdot)$, non-parametric
- (d_C, w_C, z_C) a ‘complete’ observation
- (d_R, w_R) has a missing covariate
- $f(x; \theta, g) = f(d_C, w_C | z_C; \theta)g(z_C) \int f(d_R, w_R | z; \theta)g(z)dz$



$$\begin{aligned}x &= (d_C, w_C, z_C, d_R, w_R) \\ \theta &= \gamma, \beta, \alpha_0, \alpha_1, \sigma^2\end{aligned}$$

$$EL(\theta, g) = f(d_C, w_C | z_C; \theta)g(z_C) \int f(d_R, w_R | z; \theta)g(z)dz$$

Profile likelihood

Murphy and vdW, 2000

$\mathbb{E}R^d$

$$1. \sqrt{n}(\hat{\theta} - \theta_0) = \frac{1}{\sqrt{n}}\tilde{i}^{-1}(\theta_0)\tilde{U}(\theta_0) + o_p(1)$$

$\beta \rightarrow \theta$

Semi-parametric
model $\lambda(\cdot)$, θ ERF

- $\tilde{U}(\theta_0) = \frac{\partial \ell(\theta, \lambda)}{\partial \theta} - \text{Proj}_g \frac{\partial \ell(\theta, \lambda)}{\partial \theta}$

- projection of $\partial \ell_\theta$ onto the closed linear span of the score functions for $\lambda(\cdot)$

- $\tilde{i}(\theta_0) = \text{var}\{\tilde{U}_j(\theta_0)\}$

$$2. \ell_p(\hat{\theta}) = \ell_p(\theta_0) + \frac{1}{2}n(\hat{\theta} - \theta_0)^T \tilde{i}(\theta_0)(\hat{\theta} - \theta_0) + o_p(1)$$

- for any random sequence $\tilde{\theta}_n \xrightarrow{P} \theta_0$, plus conditions on the model,

$$\ell_p(\tilde{\theta}_n) = \ell_p(\theta_0) + (\tilde{\theta}_n - \theta_0)^T \sum_{j=1}^n \tilde{U}_j(\theta_0) - \frac{1}{2}n(\tilde{\theta}_n - \theta_0)^T \tilde{i}^{-1}(\theta_0)(\tilde{\theta}_n - \theta_0)$$

$$+ o_p(\sqrt{n}||\tilde{\theta}_n - \theta_0|| + 1)^2$$

... inference

-

$$\begin{aligned}\ell_p(\tilde{\theta}_n) &= \ell_p(\theta_0) + (\tilde{\theta}_n - \theta_0)^T \sum_{j=1}^n \tilde{U}_j(\theta_0) - \frac{1}{2} n (\tilde{\theta}_n - \theta_0)^T \tilde{\iota}^{-1}(\theta_0) (\tilde{\theta}_n - \theta_0) \\ &\quad + o_p(\sqrt{n} \|\tilde{\theta}_n - \theta_0\| + 1)^2\end{aligned}$$

- this result (3.) gives (1.) and (2.)
- as in parametric models, lead to

$$(\hat{\theta} - \theta_0) \sim N\{\mathbf{0}, \tilde{\iota}^{-1}(\theta_0)\}$$

- and likelihood ratio test

$$2\{\ell_p(\hat{\theta}) - \ell_p(\theta_0)\} \sim \chi_d^2$$

- proof uses least favourable sub-models through the true model
- effectively turns infinite-dimensional parameter finite

-

$$\ell(\beta, \lambda(\cdot); y, d) = \sum_{i=1}^n [d_i \{x_i \beta + \log \lambda(y_i)\} - \Lambda(y_i) \exp(x_i \beta)]$$

- score function for β :

$$\partial \ell / \partial \beta = \sum_{i=1}^n \{d_i x_i - x_i e^{x_i \beta} \Lambda(y_i)\}$$

- score function for $\lambda(\cdot)$: in the ‘direction’ $h(\cdot)$

$$\sum_{i=1}^n d_i h(y_i) - e^{x_i \beta} \int_0^{y_i} h(t) d\Lambda(t)$$

- we need to project $\partial \ell / \partial \beta$ on the space spanned by the nuisance score functions
- result: $\sum_{i=1}^n d_i \left(x_i - \frac{M_1}{M_0}(y_i) \right) - e^{x_i \beta} \int_0^{y_i} \left(x_i - \frac{M_1}{M_0}(t) \right) d\Lambda(t)$

Semi-parametric models

- profile log-likelihood can (often) be defined using a least favorable sub-model finite dimensional
- standard likelihood asymptotics apply for inference based on the profile log-likelihood
- in other examples, we see that profiling out large numbers of nuisance parameters can lead to poor finite sample results
- ?does this happen in semi-parametric models?
- seems unlikely for proportional hazards regression complete separation of the parameters?
- other examples in vdW & M include current status data, gamma frailty models, partially missing data, ... *each one is analysed specially*

Infinite-dimensional models

nonparametric like

- recall that $L(\theta; y) \propto f(y; \theta)$

$$L(f(\cdot); y) \propto f(y)$$

$f(y; \theta)$ a density w.r. to dominating measure

$f(\cdot) \in$ same class

- more abstract definition:

if a probability measure Q is absolutely continuous w.r. to a probability measure P , and both possess densities w.r. to a measure μ , then the likelihood of Q w.r. to P is the Radon-Nikodym derivative f

$$\frac{dQ}{dP} = \frac{q}{p}, \text{ a.e. } P$$

- some semi-parametric models have a dominating measure, and a family of densities
- some can be handled by the notion of empirical likelihood
- some may use mixtures of these

- Definition: Given a measure P , and a sample (y_1, \dots, y_n) , the **empirical likelihood function** is

$$EL(P; y) = \prod_{i=1}^n P(\{y_i\}),$$

where $P\{y\}$ is the measure of the one-point set $\{y\}$

- Definition: Given a model \mathcal{P} , a maximum likelihood estimator is the distribution \hat{P} that maximizes the empirical likelihood over \mathcal{P}
- may or may not exist

Example: the empirical distribution

vdW 25.68

- \mathcal{P} is the set of all probability distributions on a measurable space $(\mathcal{Y}, \mathcal{A})$
1-point sets are measurable

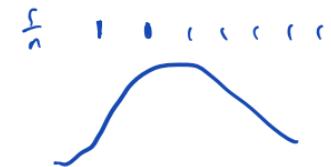
- suppose the observed values y_1, \dots, y_n are distinct

- $\{(P\{y_1\}, \dots, P\{y_n\}); P \in \mathcal{P}\} \iff (p_1, \dots, p_n), p_i \geq 0, \sum p_i = 1$

y_1, \dots, y_n is df.

- empirical likelihood maximized at

$$\left(\frac{1}{n}, \dots, \frac{1}{n}\right) \not\in \hat{\mathcal{P}}_{EL}$$



- empirical distribution function is the nonparametric MLE

$$F_n(\cdot) = n^{-1} \sum 1(Y_i \leq \cdot)$$

- EL is not the same as $\prod f(y_i)$, even if P has a density f

Compare Owen, Ch. 2

- for $y \in \mathbb{R}$, define $F(y) = \Pr(Y \leq y)$ and
 $F(y^-) = \Pr(Y < y)$

- for y_1, \dots, y_n the nonparametric likelihood function is

$$L(F) = \prod_{i=1}^n \{F(y_i) - F(y_i^-)\},$$

- hence 0 if F is continuous

- Theorem 2.1 of Owen:

$$\underline{L(F)} < \underline{L(F_n)}, \quad F_n(y) = \frac{1}{n} \sum \mathbf{1}\{y_i \leq y\}$$

- there is a likelihood function on the space of distribution functions for which the empirical c.d.f. is the maximum likelihood estimator

why does this fail for densities?

s.t. $T(F) = \theta$
↑
 f or cdf

Ex. $\int x dF(x) = \theta$

- profile version of empirical likelihood

$$\mathcal{R}(\theta) = \sup \left\{ \frac{L(F)}{L(F_n)} \mid F \in \mathcal{F}, T(F) = \theta \right\}$$

- example: $T(F) = \int x dF(x)$

$$\mathcal{R}(\theta) = \max \left\{ \prod_{i=1}^n np_i \mid \sum_{i=1}^n p_i y_i = \theta, p_i \geq 0, \sum p_i = 1 \right\}$$

- For y_1, \dots, y_n i.i.d. F_0 , $E(y_i) = \theta_0, \text{var}(y_i) < \infty$,

$$\pi(\theta|y) \sim N(\hat{\theta}, j_n^{-1}(\hat{\theta}))$$

with fl.

\mathcal{R} a relative likelihood, hence $n p_i$

n nuis. par.
 $F(\cdot)$

prof. led out

emp prof led (5k.)

$$-2 \log \mathcal{R}(\theta_0) \xrightarrow{d} \chi_1^2, \quad n \rightarrow \infty$$

Theorem 2.2 Owen

$$\hat{p}_i = \frac{1}{n} \frac{1}{\{1 + \alpha(y_i - \theta_0)\}}, \quad \frac{1}{n} \sum_{i=1}^n \frac{y_i - \theta_0}{1 + \alpha(y_i - \theta_0)} = 0$$

- $\Pr(Y = 1 | V, W) = \frac{e^{\theta V + \eta(W)}}{1 + e^{\theta V + \eta(W)}}$
- sample $(Y_i, V_i, W_i), i = 1, \dots, n$ independent

-

$$L(\theta, \eta; \underline{Y}) \propto \prod_{i=1}^n \left\{ \frac{e^{\theta V_i + \eta(W_i)}}{1 + e^{\theta V_i + \eta(W_i)}} \right\}^{y_i} \left\{ \frac{1}{1 + e^{\theta V_i + \eta(W_i)}} \right\}^{1-y_i}$$

- $\tilde{\eta}(w_i) = \infty$ when $y_i = 1$, $\tilde{\eta}(w_i) = -\infty$ when $y_i = 0$ gives

$$L(\theta, \tilde{\eta}) \rightarrow \infty$$

we can't maximize it

- suggestion: penalized log-likelihood

$$\log L(\theta, \eta; \underline{Y}) - \hat{\alpha}_n^2 \int \{\eta^{(k)}(w)\}^2 dw$$

Composite likelihood

(Likelihood when model is wrong)

- Vector observation: $Y \sim f(y; \theta)$, $Y \in \mathcal{Y} \subset \mathbb{R}^m$, $\theta \in \mathbb{R}^d$

y_1, \dots, y_n

- Set of events: $\{\mathcal{A}_k, k \in K\}$

=

ERan ...

- Composite Log-Likelihood:

$$cl(\theta; y) = \sum_{k \in K} w_k \ell_k(\theta; y) \quad y \in \mathbb{R}^m$$

log of a density

Lindsay, 1988

- $\ell_k(\theta; y) = \log\{f(\underbrace{\{y \in \mathcal{A}_k\}}; \theta)\}$ log-likelihood for an event

- $\{w_k, k \in K\}$ a set of weights

- also called:

- pseudo-likelihood (spatial modelling)
- quasi-likelihood (econometrics)
- limited information method (psychometrics)

Examples of composite log-likelihood

$$\sum_{r=1}^m w_r \log \underline{f_1(y_r; \theta)}$$

An's

Independence

$$\sum_{r=1}^m \sum_{s>r} w_{rs} \log \underline{f_2(y_r, y_s; \theta)}$$

Pairwise

$$\sum_{r=1}^m w_r \log f(y_r | y_{(-r)}; \theta) \leftarrow \text{=}$$

Conditional

$$\sum_{r=1}^m \sum_{s>r} w_{rs} \log \underline{f(y_r | y_s; \theta)}$$

All pairs conditional

$$\sum_{r=1}^m w_r \log f(y_r | y_{r-1}; \theta) \quad y_1, \dots, y_m$$

Time series

$$\sum_{r=1}^m w_r \log f(y_r | \text{'neighbours' of } y_r; \theta)$$

Spatial

Small blocks of observations; pairwise differences; ...
your favourite combination...

Derived quantities

single response y with density $f(y; \theta)$, $y \in \mathbb{R}^m, \theta \in \mathbb{R}^d$

composite log-likelihood

$$cl(\theta; y) = \log cL(\theta; y) = \sum_k w_k \ell_k(\theta; y)$$

composite score function

$$U_{CL}(\theta) = \partial cl(\theta; y) / \partial \theta$$

sensitivity

and Bartlett
w/ Bartlett hole
in general

variability

but in earlier
work

Godambe information

$$H(\theta) = E_\theta \left\{ -\frac{\partial^2 cl(\theta; y)}{\partial \theta \partial \theta^T} \right\}$$

wrt true model $f(y; \theta)$

$$J(\theta) = E_\theta \left\{ U_{CL}(\theta) U_{CL}(\theta)^T \right\}$$



$$H(\theta) = J(\theta)$$

$$\text{var } U(\theta) = -E(\ell''(\theta))$$

in ~~the~~ usual th.

$$G(\theta) = H(\theta) J^{-1}(\theta) H(\theta)$$

... derived quantities

sample $y = (y_1, \dots, y_n)$ with joint density $f(y; \theta)$, $y \in \mathbb{R}^m, \theta \in \mathbb{R}^d$

score function

$$U_{CL}(\theta) = \frac{\partial}{\partial \theta} c\ell(\theta; y) = \sum_{i=1}^n \frac{\partial}{\partial \theta} c\ell(\theta; y_i)$$

y : vector
but we've used
some A_k

maximum composite
likelihood estimate

$$\hat{\theta}_{CL} = \hat{\theta}_{CL}(y) = \arg \sup_{\theta} c\ell(\theta; y)$$

score equation

$$U_{CL}(\hat{\theta}_{CL}) = c\ell'(\hat{\theta}_{CL}) = 0$$

composite LRT

$$w_{CL}(\theta) = 2\{c\ell(\hat{\theta}_{CL}) - c\ell(\theta)\}$$

Godambe information

$$G(\theta) = G_n(\theta) = H_n(\theta)J_n^{-1}(\theta)H_n(\theta) = O(n)$$

Inference

- Sample: Y_1, \dots, Y_n , i.i.d., $CL(\theta; \underline{y}) = \prod_{i=1}^n CL(\theta; y_i)$

$$\hat{\theta}_{CL} - \theta \sim N\{0, G^{-1}(\theta)\}$$

$$G_n(\theta) = H(\theta)J(\theta)^{-1}H(\theta)$$

$$U(\hat{\theta}_{CL}) \doteq U(\theta) + (\hat{\theta}_{CL} - \theta) \partial_\theta U(\theta)$$

$$\hat{\theta}_{CL} - \theta \doteq -\partial_\theta U(\theta)^{-1}U(\theta) \doteq H^{-1}(\theta)U(\theta)$$

$$U(\theta) \sim N\{0, J(\theta)\}$$

$$H^{-1}(\theta)U(\theta) \sim N\{0, H^{-1}(\theta)J(\theta)H^{-T}(\theta)\}$$

conclude

$$\sqrt{n}(\hat{\theta}_{CL} - \theta) \sim N\{0, G^{-1}(\theta)\}$$

$$U(\theta) = \mathcal{J}'(\theta)$$

~~$$= \cancel{\mathcal{J}'(\theta)}$$~~

$$\sum_{i=1}^n \mathcal{J}'_i(\theta)$$

$$U = U_{CL}$$

$$E_\theta U(\theta) = E_\theta \sum_{i=1}^n \{ \mathcal{J}'_i(\theta) \}$$

$$> \sum E_\theta \{ \mathcal{J}'_i(\theta) \} = 0$$

$$\text{var}_\theta U(\theta) = E_\theta \{ U^2(\theta) \}$$

$$= J(\theta)$$

$$d(\theta) = \sum_{i=1}^n \left\{ \sum_{r \leq s} h f_r(y_r, y_s; \theta) \right\}$$

$$\begin{aligned} E_\theta\{d'(\theta)\} &= E\left(\sum_{i=1}^n \sum_{r \leq s} \frac{\partial}{\partial \theta} f_r(y_r, y_s; \theta)\right) \\ &= \sum_{i=1}^n \sum_{r \leq s} \left\{ \underbrace{\sum_{r \leq s} \left(\frac{\partial}{\partial \theta} f_r(y_r, y_s; \theta) \right)}_{\text{f}(y; \theta)} \underbrace{f_r(y_r; \theta)}_{dy_r dy_s} dy_r dy_s \right\} \\ &= \sum_{i=1}^n \sum_{r \leq s} \left\{ \underbrace{\frac{\partial}{\partial \theta} f_r(y_r, y_s; \theta)}_{\text{usual avg}} f_r(y_r, y_s; \theta) \right\} \\ \frac{\partial}{\partial \theta} \int f_r = 0 & \quad \text{usual avg} = 0 \end{aligned}$$

$$\hat{\theta}_{CL}'(\hat{\theta}_{CL}) = 0$$

$$= U_{CL}'(\theta) + (\hat{\theta}_{CL} - \theta) U_{CL}''(\theta) + R$$

$$\hat{\theta}_{CL} - \theta \underset{\substack{U_{CL}'(\theta) \leftarrow \text{CLT var } J \\ -U_{CL}''(\theta) \leftarrow \text{WLLN E H}}}{\sim}$$

$$\sqrt{n} \{ \hat{\theta}_{CL} - \theta \} \xrightarrow{d} N(0, G_1^{-1}(\theta))$$

$G_1(\theta)$ God. info $\frac{1}{HJH}$ $\frac{\partial s^2}{\partial \theta}$ $\hat{\theta}_{CL}$

... inference

- $w(\theta) = 2\{\text{cl}(\hat{\theta}_{CL}) - \text{cl}(\theta)\} \sim \sum_{a=1}^d \mu_a Z_a^2$ $Z_a \sim N(0, 1)$ $Z_a^2 \sim \chi_1^2$ χ_d^2 $d = d - \theta$

- μ_1, \dots, μ_d eigenvalues of $J(\theta)H(\theta)^{-1}$
- \cdot

$$\text{cl}(\hat{\theta}_{CL}) - \text{cl}(\theta) \doteq \frac{1}{2}(\hat{\theta}_{CL} - \theta)^T \{-\text{cl}''(\hat{\theta}_{CL})\}(\hat{\theta}_{CL} - \theta)$$

- non-central χ^2 limit

$(HJ^{-1}H)^{-1}$

- $J(\theta) = \text{var}_{\text{cl}} U(\theta), \quad H(\theta) = -E \partial_\theta U(\theta)$

- if $J(\theta) = H(\theta)$, $w(\theta) \sim \chi_d^2$

- if $d = 1$, $w(\theta) \sim \mu_1 \chi_1^2 = J(\theta) H^{-1}(\theta) \chi_1^2$

y_1, y_2, \dots, y_n H, J both scalars
 (\quad) dependence

Example: symmetric normal

$$(Y_1, \dots, Y_m)^T \sim N\left(\mu\left(\begin{array}{c} 0 \\ \vdots \\ 0 \end{array}\right), \Sigma\left(\begin{array}{cc} 1 & \rho & \dots & \rho \\ \rho & 1 & \dots & \rho \\ \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \dots & 1 \end{array}\right)\right)$$

- $Y_i \sim N(0, R)$, $\text{var}(Y_{ir}) = 1$, $\text{corr}(Y_{ir}, Y_{is}) = \rho$
- compound bivariate normal densities to form pairwise likelihood

$$\prod_{i=1}^n \prod_{s=1}^{m-1} f_2(y_{ir}, y_{is}; \rho) \quad N_2\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}\right)$$

$$\begin{aligned} \text{cl}(\rho; y_1, \dots, y_n) &= -\frac{nm(m-1)}{2} \log(1-\rho^2) - \frac{m-1+\rho}{2(1-\rho^2)} SS_w \\ &\quad - \frac{(m-1)(1-\rho)}{2(1-\rho^2)} \frac{SS_b}{m} \\ SS_w &= \sum_{i=1}^n \sum_{s=1}^{m-1} (y_{is} - \bar{y}_{i.})^2, \quad SS_b = \sum_{i=1}^n y_{i.}^2 \end{aligned}$$

$$\begin{aligned} |\quad | = 1 - \rho^2 \quad \ell(\rho; y_1, \dots, y_n) &= -\frac{n(m-1)}{2} \log(1-\rho) - \frac{n}{2} \log\{1 + (m-1)\rho\} \quad |\underline{R}|^{-1/2} \\ &\quad - \frac{1}{2(1-\rho)} SS_w - \frac{1}{2\{1 + (m-1)\rho\}} \frac{SS_b}{m} \end{aligned}$$

... symmetric normal

$$H J^{-1} H$$

$$\text{• a. } \text{var}(\hat{\rho}) = \frac{2}{nm(m-1)} \frac{(1 + (m-1)\rho)^2(1-\rho)^2}{1 + (m-1)\rho^2} \quad 0 < \rho < 1$$

$$\text{• a. } \text{var}(\hat{\rho}_{CL}) = \frac{2}{nm(m-1)} \frac{(1-\rho)^2 c(m, \rho)}{(1+\rho^2)^2}$$

$$\text{• } c(m, \rho) = (1-\rho)^2(3\rho^2+1) + m\rho(-3\rho^3+8\rho^2-3\rho+2) + m^2\rho^2(1-\rho)^2$$

$$\text{a.var}(\hat{\rho}_{CL}) = \frac{2}{nm(m-1)} \frac{(1-\rho)^2}{(1+\rho^2)^2} c(m, \rho)$$

$$O\left(\frac{1}{n}\right)$$
$$n \longrightarrow \infty$$

$$O(1)$$
$$\underbrace{m \longrightarrow \infty}_{n \text{ fixed}}$$

... symmetric normal

$$\frac{a.var(\hat{\rho})}{a.var(\hat{\rho}_{CL})}, \quad m = 3, 5, 8, 10$$

(Cox & Reid, 2004)

