

Topics in Likelihood Inference

STA4508H

Nancy Reid
University of Toronto

February 9, 2022

A4 | NEWS

G THE GLOBE AND MAIL | WEDNESDAY, FEBRUARY 9, 2022

Omicron less lethal than Delta, research finds

Although variant isn't as likely to lead to severe illness, transmissibility has caused daily death counts to eclipse those of previous wave

ANDREA WOO

People infected with the Omicron variant of SARS-CoV-2 are less likely to die or experience severe outcomes compared with those infected with the Delta variant, emerging research shows.

However, the variant's hyper-transmissibility has led to daily death numbers that have eclipsed those of the previous wave, and a recent study of hospitalizations, with Canada's older population still at highest risk.

Federal and provincial governments have begun to signal a turning point in the fight against COVID-19, with some leaders pointing to the variant's relatively less severe outcomes as one reason to lift remaining restrictions within weeks. Inexplicably, that has not yet been done, only when governments have done everything possible to protect the most vulnerable residents.

Samir Sinha, director of geriatric research at St. Michael's Hospital in Toronto, says the evidence is clear:

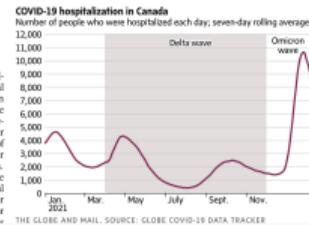
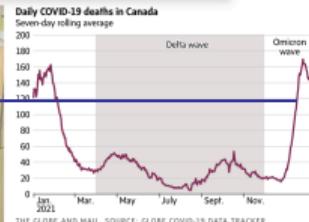


Nurses tend to a COVID-19 patient in the intensive-care unit at a Sarnia, Ont., hospital last month. Canada recorded almost 170 deaths a day during the peak of the Omicron wave in late January, based on a seven-day rolling average. CHRIS YOUNG/THE CANADIAN PRESS

evident among vaccinated and unvaccinated patients and those with and without prior COVID-19 infection, write the authors of a study that should have eased

concerns about the Omicron variant. Of those hospitalized with Omicron, 37.5 per cent needed critical care, while that number for Delta was 41.4 per cent. In comparison, 33.1 per cent of those hospitalized with Omicron required critical care, and 6.5 per cent died. The median length of hospital stay was seven days for Delta, and four days for Omicron.

Reviews of 550 cases of people across B.C. admitted to hospital with COVID-19 show that 44 per cent who were hospitalized for reasons unrelated to the virus



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. **likelihood inference in high dimensions**
6. simulated likelihood, indirect inference
7. bootstrap likelihood, h -likelihood, weighted likelihood, pseudo-likelihood, local likelihood, sieve likelihood

Presentations

- Feb 16 Angela: Cox (2013)
 Robert: Barndorff-Nielsen and Cox (1979)
 Shiki: Solomon and Cox (1992)
- Feb 23 Hengchao: Rotnitzky et al. (2000)
 Siyue: De Stavola and Cox (2008)
 Manuel: Battey and Cox (2018)
 Ziang: Cox (1975)

Feb 16 in SS 1087; Feb 23 online

typo in exercise 2 for Feb 2 $\exp\{\ell(\theta) - \ell(\hat{\theta})\}$

correct version now on web page

exercises Jan 26 has details about report structure

- y_1, \dots, y_n independent observations $\sim G(\cdot)$, density $g(\cdot)$
- we fit the **incorrect** model $f(y; \theta)$
- Kullback-Liebler (KL) divergence between $f(y; \theta)$ and $g(y)$ is defined as

$$KL(\theta) = \int \log \left\{ \frac{g(y)}{f(y; \theta)} \right\} g(y) dy$$

Wikipedia writes $D_{KL}(G||F)$ or more precisely $D_{KL}(P_G||P_F)$

- $KL(\theta) \geq 0$, and $KL(\theta) = 0 \iff f(y; \theta) \equiv g(y)$
- define $\theta^* = \arg \min KL(\theta)$
- $f(y; \theta^*)$ is **closest to $G(\cdot)$** in the family $\{f(\cdot; \theta), \theta \in \Theta\}$

... misspecified models

$$KL(\theta) = \int \log \left\{ \frac{g(y)}{f(y; \theta)} \right\} g(y) dy$$

- $\theta^* = \arg \min_{\theta} KL(\theta)$
- $\theta^* = \arg \max_{\theta} \int \log\{f(y; \theta)\} g(y) dy = \arg \max_{\theta} E_G\{\ell(\theta; y)\}$
- leads to a proof that the maximum likelihood estimator converges to θ^* under some smoothness conditions, etc.
- if $g(y) = f(y; \theta_0)$, then $\theta^* = \theta_0$ true density is in the model family
- otherwise θ^* is the ‘least false’ parameter value

... misspecified models

- Example: true model G is log-normal $\log y \sim N(\mu, \sigma^2)$ $g(y) = ?$
- fitted model has density $f(y; \theta) = \frac{1}{\theta} \exp(-\frac{y}{\theta})$
- $E_G\{\ell(\theta; y)\} = -\log \theta - E_G\left(\frac{y}{\theta}\right)$
- $\theta^* = E_G(y) = \exp(\mu + \sigma^2/2)$ $\arg \max_{\theta} E_G\{\ell(\theta; y)\}$
- If we fit $\{f(y; \theta) : \theta > 0\}$ to a sample y_1, \dots, y_n we get $\hat{\theta} = \bar{y}$
- WLLN under sampling from $G(\cdot)$, $\bar{y} \xrightarrow{P} E_G(y) = \theta^*$

θ^* is a ‘meaningful’ parameter, regardless of the underlying model

... misspecified models

- viewing θ as a convenient summary of the data, we can consider properties of likelihood-based inference under the true model g Kent, White, 1982
- this can be cumbersome: studying robustness to local departures from an assumed model might be more relevant in practice
- composite likelihood is a special type of misspecification Lindsay, 1988
- another is the framework of generalized estimating equations, with dependence modelled by using a ‘working covariance’ Liang & Zeger, 1986
- indirect inference also uses a working (simplified) model that is adjusted using simulations from the true model Gouerieroux et al, 1993

Likelihood inference in misspecified models

- maximum likelihood estimate as usual: $(\partial/\partial\theta)\ell(\hat{\theta}; y) = \mathbf{0}$
- consistent for θ^* , the ‘least-false’ value
- $E_G U(\theta) = E_G(\partial/\partial\theta)\ell(\theta; y) = \int (\partial/\partial\theta)\ell(\theta; y)g(y)dy = \mathbf{0}$, only at θ^*
- $E_G(-\partial^2/\partial\theta^2)\ell(\theta; y) = \int (-\partial^2/\partial\theta^2)\ell(\theta; y)g(y)dy \equiv H(\theta)$
 $= H_G(\theta)$
- $E_G\{(\partial/\partial\theta)\ell(\theta; y)\}^2 = \int \{(\partial/\partial\theta)\ell(\theta; y)\}^2 g(y)dy \equiv J(\theta) \neq H(\theta)$
- $(\hat{\theta} - \theta^*) = H(\theta^*)^{-1}U(\theta^*)\{1 + o_p(1)\}$
 $U(\theta) = (\partial/\partial\theta)\ell(\theta)$
- $\hat{\theta} \stackrel{\text{d}}{\sim} N\{\theta^*, \mathcal{G}^{-1}(\theta^*)\}$
 $\mathcal{G}(\theta) = H(\theta)J^{-1}(\theta)H(\theta)$

Examples

- y_1, \dots, y_n i.i.d. $\sim G$; we assume $f(y; \theta)$ is $N(\mu, \sigma^2)$

- $\hat{\mu} = \bar{y}, \quad \hat{\sigma}^2 = \frac{1}{n} \sum (y_i - \bar{y})^2 \quad \mu^* = \mu_G, \sigma^{*2} = \sigma_G^2$

- $\partial_\theta \ell(\theta; y) = [\Sigma(y_i - \mu)/\sigma^2, \quad -(n/2\sigma^2) + \Sigma(y_i - \mu)^2/(2\sigma^4)]$

- $H(\theta_G) = n \begin{bmatrix} 1/(\sigma_G^2) & 0 \\ 0 & 1/(2\sigma_G^4) \end{bmatrix} \quad H(\theta) = -E_G \ell''(\theta; y)$

$$J(\theta) = \text{Var}_G \ell'(\theta; y)$$

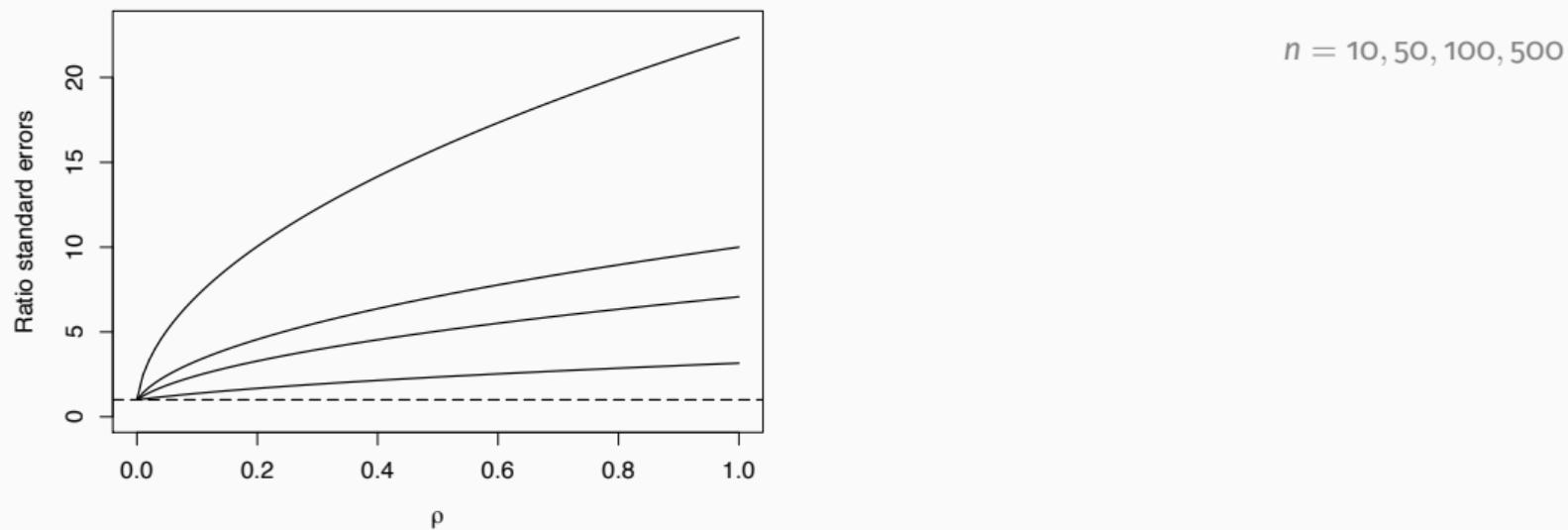
$$J(\theta_G) = n \begin{bmatrix} 1/(\sigma_G^2) & \mu_3/(2\sigma_G^6) \\ \mu_3/(2\sigma_G^6) & (\mu_4 - \sigma_G^4)/(4\sigma_G^8) \end{bmatrix} \quad \text{ntbc}$$

- $G(\theta_G) = n^{-1} \begin{bmatrix} \sigma_G^2 & \mu_3 \\ \mu_3 & \mu_4 - \sigma_G^4 \end{bmatrix}$ not uncorrelated

... examples

- linear regression model $G: y = X\beta_0 + \epsilon, \quad \epsilon \sim (0, \sigma^2 R)$ linear regression; correlated errors
- working model $f(y; \beta) = N(X\beta, \sigma^2 I)$ uncorrelated errors
- $\hat{\beta} = (X^T X)^{-1} X^T y$ least squares estimator; mle under normality
- $E_G(\hat{\beta}) = \beta_0$ LSE unbiased (consistent)
- $\text{var}_G(\hat{\beta}) = \sigma^2 (X^T X)^{-1} X^T R X (X^T X)^{-1}$ sandwich variance
- working model variance is incorrect $\sigma^2 (X^T X)^{-1}$
- single intercept model $\hat{\beta} = \bar{y}, R_{ij} = \rho, \text{var}(\bar{y}) = (\sigma^2/n)\{1 + \rho(n - 1)\}$ dependence is a killer

- $y_i \sim N(\beta, \sigma^2)$, $R_{ij} = \rho$
 - $\text{var}_G(\hat{\beta}) = (\sigma^2/n)\{1 + \rho(n - 1)\}$
- $\rho = 0.1, n = 100$ ratio of se's 3.3



Recap: Composite likelihood

- Vector observation: $Y \sim f(y; \theta)$, $Y \in \mathcal{Y} \subset \mathbb{R}^m$, $\theta \in \mathbb{R}^d$
- Set of events: $\{\mathcal{A}_k, k \in K\}$
- Composite Log-Likelihood:

Lindsay, 1988

$$cl(\theta; y) = \sum_{k \in K} w_k \ell_k(\theta; y)$$

- $\ell_k(\theta; y) = \log\{f(\{y \in \mathcal{A}_k\}; \theta)\}$ log-likelihood for an event
- $\{w_k, k \in K\}$ a set of weights
- Note that θ is assumed to have the same meaning in ℓ_k and in the full model
- This makes it different from a completely mis-specified model

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

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

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

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

- Sample: Y_1, \dots, Y_n , i.i.d., $CL(\theta; \underline{y}) = \prod_{i=1}^n CL(\theta; y_i)$
- $\hat{\theta}_{CL} - \theta \sim N\{\mathbf{o}, G^{-1}(\theta)\}$ $G_n(\theta) = H(\theta)J(\theta)^{-1}H(\theta)$
- $U(\hat{\theta}_{CL}) \doteq U(\theta) + (\hat{\theta}_{CL} - \theta)\partial_\theta U(\theta)$ $U = U_{CL}$
- $\hat{\theta}_{CL} - \theta \doteq -\partial_\theta U(\theta)^{-1}U(\theta) \doteq H^{-1}(\theta)U(\theta)$
- $U(\theta) \sim N\{\mathbf{o}, J(\theta)\}$
- $H^{-1}(\theta)U(\theta) \sim N\{\mathbf{o}, H^{-1}(\theta)J(\theta)H^{-T}(\theta)\}$
- conclude

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

... inference

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

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

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

- non-central χ^2 limit

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

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

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

H, J both scalars

Example: symmetric normal

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

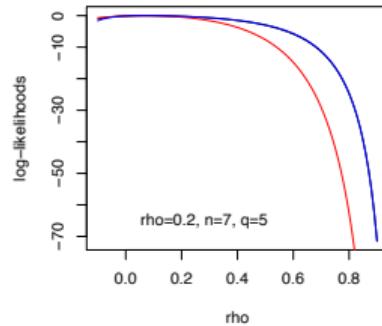
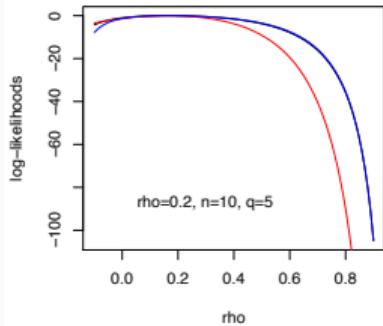
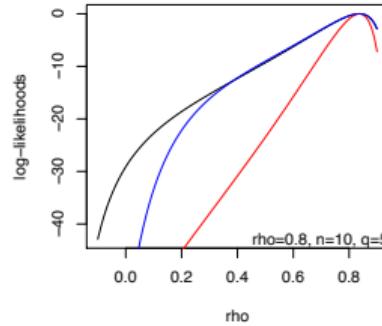
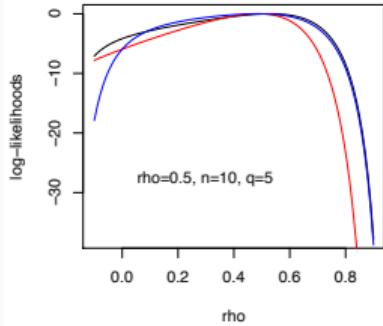
$$c\ell(\rho; y_1, \dots, y_n) = -\frac{nm(m-1)}{4} \log(1-\rho^2) - \frac{m-1+\rho}{2(1-\rho^2)} SS_w$$

$$-\frac{(m-1)(1-\rho)}{2(1-\rho^2)} \frac{SS_b}{m}$$

$$SS_w = \sum_{i=1}^n \sum_{s=1}^m (y_{is} - \bar{y}_{i.})^2, \quad SS_b = \sum_{i=1}^n y_{i.}^2$$

$$\ell(\rho; y_1, \dots, y_n) = -\frac{n(m-1)}{2} \log(1-\rho) - \frac{n}{2} \log\{1 + (m-1)\rho\}$$
$$-\frac{1}{2(1-\rho)} SS_w - \frac{1}{2\{1 + (m-1)\rho\}} \frac{SS_b}{m}$$

... symmetric normal LRT



Nuisance parameters $\theta = (\psi, \lambda)$

- constrained estimator: $\tilde{\theta}_\psi = \sup_{\theta=\theta(\psi)} c\ell(\theta; y)$
- $\sqrt{n}(\hat{\psi}_{CL} - \psi) \sim N\{0, G^{\psi\psi}(\theta)\}$ $G(\theta) = H(\theta)J(\theta)^{-1}H(\theta)$
- profile composite log-likelihood test $w(\psi) = 2\{c\ell(\hat{\theta}_{CL}) - c\ell(\tilde{\theta}_\psi)\} \sim \sum_{a=1}^{d_o} \mu_a Z_a^2$
- μ_1, \dots, μ_{d_o} are the eigenvalues of $(H^{\psi\psi})^{-1}G^{\psi\psi}$
- Godambe information needs to be estimated

Kent, 1982

$$\begin{aligned}\hat{H}(\theta) &= -\partial^2 c\ell(\hat{\theta}_{CL}) / \partial \theta \partial \theta^T \\ \hat{J}(\theta) &= n^{-1} \sum_{i=1}^n U_{CL}(\theta; y_i) U_{CL}^T(\theta; y_i)\end{aligned}$$

latter needs independent samples; can be biased and/or inefficient

- Akaike's information criterion

Varin and Vidoni, 2005

$$AIC = -2c\ell(\hat{\theta}_{CL}) + 2 \dim(\theta)$$

- derivation of AIC for misspecified likelihood leads to

$$TIC = -2c\ell(\hat{\theta}_{CL}) + 2 \text{tr}\{H(\hat{\theta}_{CL})G^{-1}(\hat{\theta}_{CL})\}$$

Takeuchi information criterion

- Bayesian information criterion

Gao and Song, 2009

$$BIC = -2c\ell(\hat{\theta}_{CL}) + \log n \dim(\theta)$$

used for selection of tuning parameters

Example: CL with dichotomized MV Normal

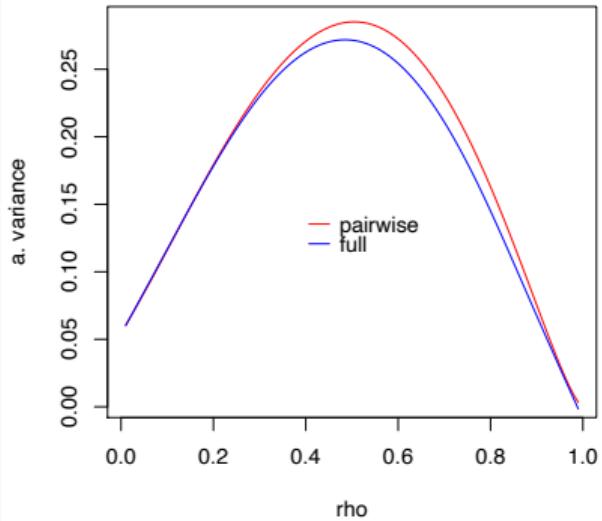
$$Y_{ir} = \mathbf{1}\{Z_{ir} > 0\} \quad Z \sim N(0, R) \quad r = 1, \dots, m; i = 1, \dots, n$$

$$\begin{aligned} \ell_2(\rho) = & \sum_{i=1}^n \sum_{s < r} \{ y_{ir} y_{is} \log P(y_r = 1, y_s = 1) + y_{ir} (1 - y_{is}) \log P_{10} \\ & + (1 - y_{ir}) y_{is} \log P_{01} + (1 - y_{ir}) (1 - y_{is}) \log P_{00} \} \end{aligned}$$

$$\text{a.var}(\hat{\rho}_{CL}) = \frac{1}{n} \frac{4\pi^2}{m^2} \frac{(1 - \rho^2)}{(m - 1)^2} \text{var}(T) \quad T = \sum_{s < r} (2y_{ir} y_{is} - y_{ir} - y_{is})$$

$$\begin{aligned} \text{var}(T) = & m^4(p_{1111} - 2p_{111} + 2p_{11} - p_{11}^2 + \frac{1}{4}) + \\ & m^3(-6p_{1111} \dots) + m^2(\dots) + m(\dots) \end{aligned}$$

$$p_{1111} = Pr(Z_r > 0, Z_s > 0, Z_t > 0, Z_u > 0)$$



Numbers incorrect in Cox & Reid 2004 Table 1

ρ	0.02	0.05	0.12	0.20	0.40	0.50
ARE	0.998	0.999	0.995	0.992	0.968	0.953

ρ	0.60	0.70	0.80	0.90	0.95	0.98
ARE	0.938	0.903	0.900	0.874	0.869	0.850

- latent variable: $z_{ir} = x'_{ir}\beta + b_i + \epsilon_{ir}, \quad \epsilon_{ir} \sim N(0, 1)$
- binary observations: $y_{ir} = 1(z_{ir} > 0); \quad r = 1, \dots, m_i; i = 1, \dots, n$
- probit model: $Pr(y_{ir} = 1 | b_i) = \Phi(x'_{ir}\beta + b_i); \quad b_i \sim N(0, \sigma_b^2)$
- likelihood

$$L(\beta, \sigma_b) = \prod_{i=1}^n \int_{-\infty}^{\infty} \prod_{r=1}^{m_i} \Phi(x'_{ir}\beta + b_i)^{y_{ir}} \{1 - \Phi(x'_{ir}\beta + b_i)\}^{1-y_{ir}} \phi(b_i, \sigma_b^2) db_i$$

- pairwise likelihood

$$CL(\beta, \sigma_b) = \prod_{i=1}^n \prod_{r < s} P_{11}^{y_{ir}y_{is}} P_{10}^{y_{ir}(1-y_{is})} P_{01}^{(1-y_{ir})y_{is}} P_{00}^{(1-y_{ir})(1-y_{is})}$$

- each $Pr(y_{ir} = j, y_{is} = k)$ evaluated using $\Phi_2(\cdot, \cdot; \rho_{irs})$

- computational effort doesn't increase with the number of random effects
- pairwise likelihood numerically stable
- efficiency losses, relative to maximum likelihood, of about 20% for estimation of β
- somewhat larger for estimation of σ_b^2

... Example

D. Renard et al. / Computational Statistics & Data Analysis 44 (2004) 649–667

663

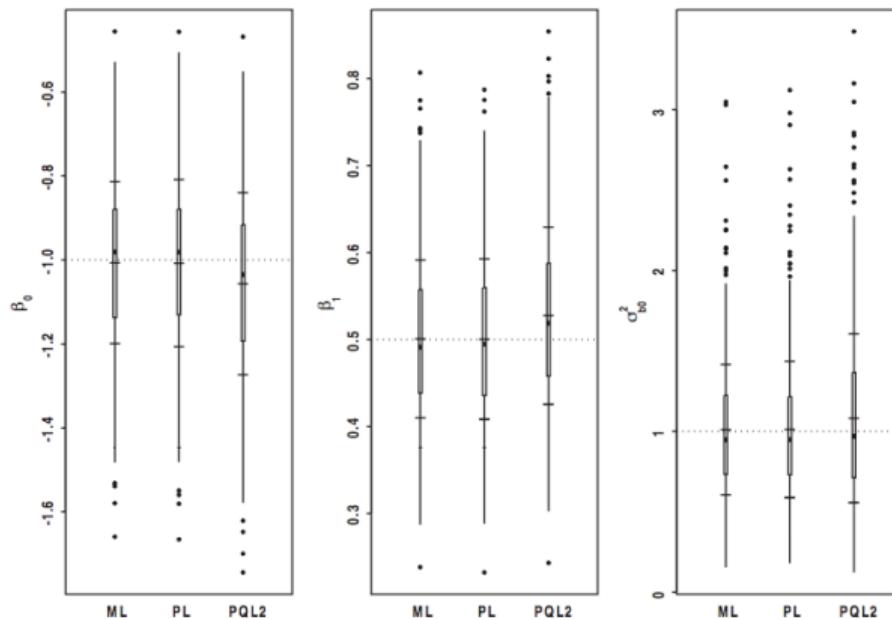


Fig. 5. Boxplots of ML, PL and PQL2 simulated parameter estimates under Model (10) with random intercept.

- subjects $i = 1, \dots, n$
- observations counts $y_{ir}, r = 1, \dots, m_i$
- model $y_{ir} \sim \text{Poisson}(u_{ir}x_{ir}^T\beta)$
- u_{i1}, \dots, u_{im_i} gamma-distributed random effects
- but correlated $\text{corr}(u_{ir}, u_{is}) = \rho^{|r-s|}$
- joint density has combinatorial number of terms in m_i ; impractical
- weighted pairwise composite likelihood

$$cL_{pair}(\beta) = \prod_{i=1}^n \frac{1}{m_i - 1} \prod_{r=1}^{m_i} \prod_{s=r+1}^{m_i} f(y_{ir}, y_{is}; \beta)$$

- weights chosen so that $\mathcal{L}_{pair} = \text{full likelihood if } \rho = 0$

Example: Varin & Czado 2010

- pain severity scores recorded at four time points morning, noon, evening, bed
- 119 patients; varying number of days per patient
- covariates: personal and weather
- response: pain score 0 1 2 3 4 5
- y_{ij} response at time t_{ij} for observation j on subject $i, j = 1, \dots, m_i$
- y_{ij}^* a **latent variable**, continuous $y_{ij}^* = x_{ij}^T \beta + u_i + \epsilon_{ij}$
- $y_{ij} = k \Leftrightarrow a_{k-1} < y_{ij}^* < a_k$
- if $u_i \sim N(0, \sigma^2)$ and $\epsilon_{ij} \sim N(0, 1)$

$$L(\theta; y) = \prod_{i=1}^n f(y_{i1}, \dots, y_{im_i}) = \\ \prod_{i=1}^n \int_{-\infty}^{\infty} \prod_{j=1}^{m_i} \{ \Phi(a_{y_{ij}} - x_{ij}^T \beta - u_i) - \Phi(a_{y_{ij}-1} - x_{ij}^T \beta - u_i) \} \phi\left(\frac{u_i}{\sigma}\right) du_i$$

$$\theta = (\underline{a}, \beta, \sigma^2)$$

... pain severity scores

- y_{ij}^* and $y_{ij'}^*$ have constant correlation $\sigma^2/(\sigma^2 + 1)$
- points nearer in time might be expected to have higher correlation
- change ϵ_{ij} i.i.d. $N(0, 1)$ to $\text{corr}(\epsilon_{ij}, \epsilon_{ij'}) = \exp(-\delta|t_{ij} - t_{ij'}|)$ $\tilde{a}_{ij} = a_{ij} - x_{ij}^\top \beta / \sqrt{\sigma^2 + 1}$

$$L(\theta; y) = \prod_{i=1}^n \int_{\tilde{a}_{y_{i1}-1}}^{\tilde{a}_{y_{i1}}} \cdots \int_{\tilde{a}_{y_{im_i}-1}}^{\tilde{a}_{y_{im_i}}} \phi_{n_i}(z_{i1}, \dots, z_{im_i}; R_i) dz_{i1} \dots dz_{im_i}$$

$$\bullet \quad R_{ijj'} = \frac{\sigma^2}{\sigma^2 + 1} + \frac{e^{-\delta|t_{ij} - t_{ij'}|}}{\sigma^2 + 1}$$

- pairwise log-likelihood:

$$c\ell(\theta; y) = \sum_{i=1}^n \sum_{j < j'}^{m_i} \log f_2(y_{ij}, y_{ij'}; \theta) \mathbf{1}_{[-q, q]}(t_{ij} - t_{ij'})$$

weights are 1 or 0, depending on distance between time points

Table 5: Migraine data. Estimates and standard errors from the pairwise likelihood with $q = 12$ for the base model (first two columns) and the best model (last two columns) accordingly to CLIC. The levels of the variable **change** are 1: change from low to high atmospheric pressure, 2: substantially unchanged atmospheric pressure, 3: change from high to low atmospheric pressure. The baseline is “no university degree, no intake of analgesics, change from low to high pressure”.

	est.	s.e.	est.	s.e.
α_2	0.588	0.046	0.588	0.046
α_3	1.136	0.069	1.136	0.069
α_4	1.786	0.079	1.787	0.080
α_5	2.505	0.109	2.506	0.111
intercept	-0.474	0.226	-0.522	0.223
university	-0.523	0.172	-0.523	0.174
analgesics	0.558	0.202	0.561	0.205
change2	—	—	0.031	0.051
change3	—	—	0.164	0.053
γ_F	0.415	0.094	0.424	0.094
γ_T	0.556	0.030	0.557	0.030
$\gamma_T - \gamma_F$	0.142	0.098	0.133	0.098
σ^2	0.566	0.110	0.564	0.111

- vector observations (X_{1i}, \dots, X_{mi}) , $i = 1, \dots, n$
- example rainfall at each of m locations
- component-wise maxima Z_1, \dots, Z_m ; $Z_j = \max(X_{j1}, \dots, X_{jn})$
- Z_j are transformed (centered and scaled)
- general theory says

$$\Pr(Z_1 \leq z_1, \dots, Z_m \leq z_m) = \exp\{-V(z_1, \dots, z_m)\}$$

- function $V(\cdot)$ can be parameterized via Gaussian process models
- example

$$\begin{aligned} V(z_1, z_2) &= z_1^{-1}\Phi\{(1/2)a(h) + a^{-1}(h)\log(z_2/z_1)\} + \\ &\quad z_2^{-1}\Phi\{(1/2)a(h) + a^{-1}(h)\log(z_1/z_2)\} \end{aligned}$$

$$Z(h) = (z_1, z_2), Z(0) = (0, 0), a(h) = h^T \Omega^{-1} h$$

$$\Pr(Z_1 \leq z_1, \dots, Z_d \leq z_m) = \exp\{-V(z_1, \dots, z_m)\}$$

- to compute log-likelihood function, need the density
- combinatorial explosion in computing joint derivatives of $V(\cdot)$

$D = 10$, one likelihood eval is a sum over 100,000 terms

- Davison et al. (2012, *Statistical Science*) used pairwise composite likelihood
- compared the fits of several competing models, using AIC analogue described above
- applied to annual maximum rainfall at several stations near Zurich

162

A. C. DAVISON, S. A. PADOAN AND M. RIBATET

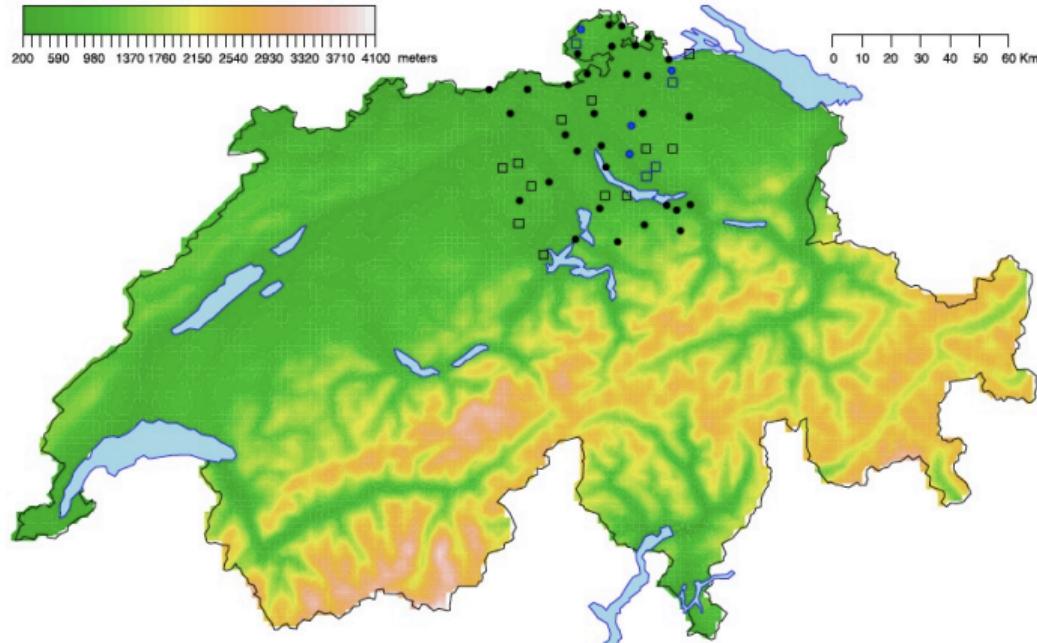


FIG. 1. Map of Switzerland showing the stations of the 51 rainfall gauges used for the analysis, with an insert showing the altitude. The 36 stations marked by circles were used to fit the models, and those marked with squares were used to validate the models. Data for the pairs of stations with blue symbols appear in Figure 2.

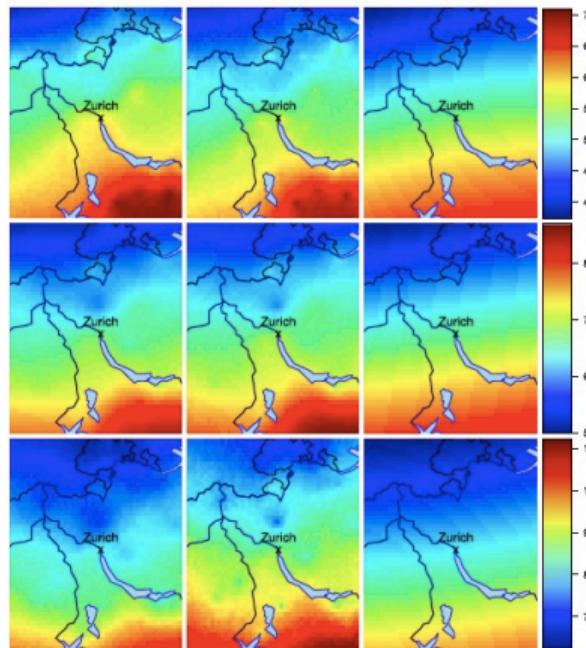


FIG. 3. Maps of the (predictive) pointwise 25-year return level estimates for rainfall (mm) obtained from the latent variable and max-stable models. The top and bottom rows show the lower and upper bounds of the 95% pointwise credible/confidence intervals. The middle row shows the predictive pointwise posterior mean and pointwise estimates. The left column corresponds to the latent variable model assuming $\text{Gamma}(5, 3)$ prior on λ . The middle column assumes the less informative priors $\lambda_\eta \sim \text{Gamma}(1, 100)$, $\lambda_\tau \sim \text{Gamma}(1, 10)$ and $\lambda_\xi \sim \text{Gamma}(1, 10)$. The right column corresponds to the extremal *t* copula model.

Example: Ising model

Ising model:

$$f(y; \theta) = \exp\left(\sum_{(j,k) \in E} \theta_{jk} y_j y_k\right) \frac{1}{Z(\theta)} \quad j, k = 1, \dots, K$$

neighbourhood contributions

$$f(y_j | y_{(-j)}; \theta) = \frac{\exp(2y_j \sum_{k \neq j} \theta_{jk} y_k)}{\exp(2y_j \sum_{k \neq j} \theta_{jk} y_k) + 1} = \exp \ell_j(\theta; y)$$

penalized CL estimation based on sample $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(n)}$

$$\max_{\theta} \left\{ \sum_{i=1}^n \sum_{j=1}^K \ell_j(\theta; \mathbf{y}^{(i)}) - \sum_{j < k} P_\lambda(|\theta_{jk}|) \right\}$$

Xue et al., 2012

Ravikumar et al., 2010

Some surprises

- Godambe information $G(\theta)$ can decrease as more component CLs are added
- pairwise CL can be less efficient than independence CL
- this can't always be fixed by weighting

Xu, 12

- parameter constraints can be important

- Example: binary vector Y ,

$$P(Y_j = y_j, Y_k = y_k) \propto \frac{\exp(\beta y_j + \beta y_k + \theta_{jk} y_j y_k)}{\{1 + \exp(\beta y_j + \beta y_k + \theta_{jk} y_j y_k)\}}$$

- this model is inconsistent

need $\theta_{jk} \equiv \theta$

- parameters may not be identifiable in the CL, even if they are in the full likelihood

Yi, 12

- Hammersley-Clifford theorem for conditionals; nothing similar (?) for marginals

when does a set of conditional densities determine a valid joint density

Quasi-likelihood

- Recall: generalized linear model y_1, \dots, y_n independent, with

$$f(y_i | x_i; \beta, \phi) = \exp[\{y_i\theta_i - c(\theta_i)\}/\phi + h(y_i, \phi)]$$

- ϕ a scale parameter in this exponential family

- $E(y_i) = \mu_i = c'(\theta_i)$

- $\text{var}(y_i) = \phi V(\mu_i) = \phi c''(\theta_i)$

variance function

- $g(\mu_i) = x_i^T \beta$

link function

- link function converts $\theta_{n \times 1}$ to $\beta_{p \times 1}$

- Standard $V(\mu)$: Normal- 1 ; Gamma- μ^2 ; Poisson- μ ; Bernoulli- $\mu(1 - \mu)$

$$\ell(\beta, \phi; \mathbf{y}) = \sum_{i=1}^n \left\{ \frac{y_i\theta_i - c(\theta_i)}{\phi} + h(y_i, \phi) \right\}$$

... quasi-likelihood

- log-likelihood

$$\ell(\beta, \phi; \mathbf{y}) = \sum_{i=1}^n \left\{ \frac{y_i \theta_i - c(\theta_i)}{\phi} + h(y_i, \phi) \right\}$$

- score function

$$\frac{\partial \ell}{\partial \beta_r} = \sum_{i=1}^n \frac{\partial \ell_i}{\partial \mu_i} \frac{\partial \mu_i}{\partial \beta_r} = \sum_{i=1}^n \frac{y_i - \mu_i}{\phi V(\mu_i)} \frac{\partial \mu_i}{\partial \beta_r}$$

- MLE

$$\sum_{i=1}^n \frac{y_i - \mu_i}{V(\mu_i)} \frac{x_{ir}}{g'(\mu_i)} = 0$$

- Bartlett identity:

$$E \left\{ \frac{\partial \ell^2}{\partial \beta_r \partial \beta_s} + \frac{\partial \ell}{\partial \beta_r} \frac{\partial \ell}{\partial \beta_s} \right\} = 0$$

... quasi-likelihood

- Suppose instead of a generalized linear model, we had only a partially specified model:

$$E(y_i) = \mu_i, \text{var}(y_i) = \phi V(\mu_i), g(\mu_i) = x_i^T \beta$$

- $g(\cdot), V(\cdot)$ known
- unbiased estimating equation

$$g(y; \beta) = \sum_{i=1}^n \frac{y_i - \mu_i}{V(\mu_i)} \frac{x_{ir}}{g'(\mu_i)}$$

- if $g(y; \tilde{\beta}) = 0$, then asymptotic variance of $\tilde{\beta}$ is

$$E \left\{ -\frac{\partial g(y; \beta)}{\partial \beta^T} \right\}^{-1} \text{var}\{g(y; \beta)\} E \left\{ -\frac{\partial g(y; \beta)}{\partial \beta} \right\}^{-1}$$

as with composite likelihood

... quasi-likelihood

- With $g(y; \beta) = \sum g_i(y_i; \beta) = \sum \frac{y_i - \mu_i}{\phi V(\mu_i)} \frac{x_{ir}}{g'(\mu_i)}$
- $E \left\{ -\frac{\partial g(y; \beta)}{\partial \beta^T} \right\} = \sum_{i=1}^n x_i x_i^T \frac{1}{g'(\mu_i)^2 \phi V(\mu_i)} = \phi^{-1} X^T W X = \text{var}\{g(y; \beta)\}$
- $W = \text{diag}(w_j), \quad w_j = \{g'(\mu_j)^2 V(\mu_j)\}^{-1}, j = 1, \dots, n$
- quasi-likelihood function

$$Q(\beta; y) = \sum_{i=1}^n \int_{y_i}^{\mu_i} \frac{y_i - u}{\phi V(u)} du$$

- this only works for models of this form
- called quasi-likelihood because $\partial Q / \partial \beta$ gives estimating equation with expected value 0, and 2nd Bartlett identity holds

Longitudinal data

- suppose now our observations come in groups: $y_{ij}, j = 1, \dots, m_i; i = 1, \dots, n$
- could be repeated measurements on subjects
- or measurements of members of the same cluster/family/group
- assume GLM-type structure $E(y_{ij}) = \mu_{ij}$, $g(\mu_{ij}) = x_{ij}^T \beta + z_{ij}^T b_i$
- random effects b_i induce correlation among observations in the same group; e.g. assume $b_i \sim N(0, \Omega_b)$
- GLM variance structure $\text{var}(y_{ij}) = V_i(\beta, \alpha)$ for example
- α are extra parameters in the variance-covariance matrix
- QL-type estimating equations

$$\sum_{i=1}^n \left(\frac{\partial \mu_i}{\partial \beta^T} \right)^T V_i^{-1}(\alpha, \beta) (y_i - \mu_i) = 0$$

$$\sum_{i=1}^n \left(\frac{\partial \mu_i}{\partial \beta^T} \right)^T V_i^{-1}(\alpha, \beta) (y_i - \mu_i) = 0$$

- parameter α in variance function doesn't divide out, as in univariate case
- we will need an estimate $\hat{\alpha}$ from somewhere
 - many suggestions in the literature
- Liang & Zeger suggested using a “working covariance matrix” to get an estimate of β
- e.g. could assume independence, or $AR(1)$, or ...
- estimates of β will still be consistent, but the asymptotic variance will be of the sandwich form as the model is misspecified
- there is no integrated function that serves as a quasi-likelihood in this setting