

Approximate Likelihoods

Nancy Reid
University of Toronto

August 14, 2015



Models and likelihood

- **Model** for the probability distribution of y given x
- **Density** $f(y | x)$ with respect to, e.g., Lebesgue measure
- **Parameters** for the density $f(y | x; \theta)$, $\theta = (\theta_1, \dots, \theta_d)$
- **Data** $y = (y_1, \dots, y_n)$ often independent

- **Likelihood function** $L(\theta; y) \propto f(y; \theta)$ (y_1, \dots, y_n)
- **log-likelihood function** $\ell(\theta; y) = \log L(\theta; y)$

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

- θ could have very large dimension, $d > n$

- θ could have infinite dimension in principle
 $E(y | x) = \theta(x)$ 'smooth'

Why likelihood?

- makes probability modelling central $\ell(\theta; y) = \log f(y; \theta)$
- emphasizes the inverse problem of reasoning $y \rightarrow \theta$
- converts a 'prior' probability to a posterior $\pi(\theta) \rightarrow \pi(\theta | y)$
- provides a conventional set of summary quantities:
maximum likelihood estimator, score function, ...
- provides summary statistics with known limiting distribution
- these define approximate pivotal quantities, based on normal distribution
- basis for comparison of models, using AIC or BIC

Widely used



Cold Regions Science and Technology

Available online 4 October 2013

In Press, Accepted Manuscript — Note to users



A Generalized Probabilistic Model of Ice Load Peaks on Ship Hulls in Broken-Ice Fields

A. Suyuthi^a, B.J. Leira^a, K. Riska^{b, c}

^a Department of Marine Technology, NTNU, Trondheim, Norway

^b Centre of Ships and Offshore Structures (CeSOS), Trondheim, Norway

^c Ilmarinen Oy, Helsinki, Finland

... widely used

▷ PP-A09-12

A Semiparametric Empirical Likelihood on the Linear Models with Covariates Parametrically Transformed

Zhang, Jing Hua
Xue, Liugen

Beijing Univ. of Tech.
Beijing Univ. of Tech.

Screen Shot 2015-08-11 at 11.32.44 PM

▷ PP-A09-17

Empirical Likelihood in Generalized Linear Models for Longitudinal Data with Dropout

Screen Shot 2015-08-11 at 11.32.54 PM

Guo, Donglin
Xue, Liugen

Beijing Univ. of Tech.
Beijing Univ. of Tech.

Screen Shot 2015-08-11 at 11.32.58 PM

▷ PP-A09-8

Generalized Empirical Likelihood Inference for Longitudinal Data with Missing Response Variables and Error-Prone Covariates

Liu, Juanfang
Xue, Liugen

Beijing Univ. of Tech.
Beijing Univ. of Tech.

Screen Shot 2015-08-11 at 11.32.34 PM

▶ MS-Fr-D-48-3

Image Reconstruction and Interpretation in Positron Emission Tomography for Small Animals (micro-PET)

Garbarino, Sara

Department of Mathematics, Univ. of Genoa

14:30-15:00

Screen Shot 2015-08-11 at 11.32.19 PM

▶ MS-Fr-D-36-2

A Randomized Likelihood Method for Data Reduction in Large-scale Inverse Problems

14:00-14:30

Screen Shot 2015-08-11 at 11.32.06 PM

Mathematics (ICIAM) is
mathematicians held
Council for Industrial
mathematicians from
ICIAM to be held at
Olympic Green.



... why likelihood?

- provides a conventional set of summary quantities:
maximum likelihood estimator, score function, ...
- provides summary statistics with known limiting distribution

Important summaries

- maximum likelihood estimator

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

- observed Fisher information

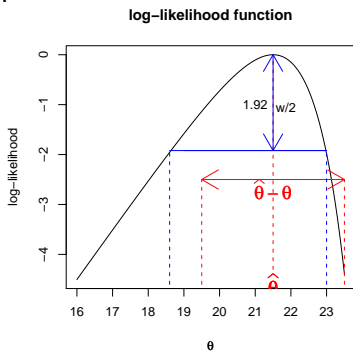
$$j(\hat{\theta}) = - \left. \frac{\partial^2 \ell(\theta)}{\partial \theta^2} \right|_{\hat{\theta}}$$

- efficient score function

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

$$\ell'(\hat{\theta}) = \mathbf{0} \text{ assuming enough regularity}$$

- $\ell'(\theta; \mathbf{y}) = \sum_{i=1}^n (\partial / \partial \theta) \log f_{Y_i}(y_i; \theta)$, y_1, \dots, y_n independent

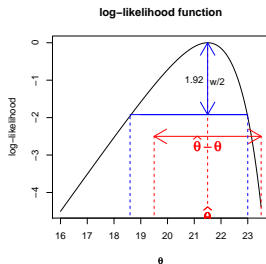


Limit theorems

- $\ell'(\theta)j^{-1/2}(\hat{\theta}) \xrightarrow{\mathcal{L}} N(0, 1)$
- $(\hat{\theta} - \theta)j^{1/2}(\hat{\theta}) \xrightarrow{\mathcal{L}} N(0, 1)$
- $2\{\ell(\hat{\theta}) - \ell(\theta)\} \xrightarrow{\mathcal{L}} \chi_1^2$
- under the model $f(y; \theta)$
regularity conditions
- approximate pivots

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

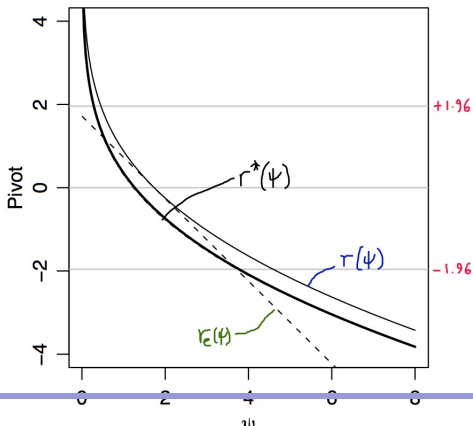
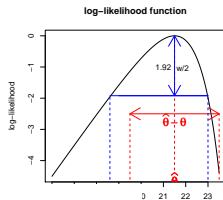
$$r(\theta) = \pm\sqrt{2\{\ell(\hat{\theta}) - \ell(\theta)\}}$$



... approximate pivots

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

$$r(\theta) = \pm\sqrt{2\{\ell(\hat{\theta}) - \ell(\theta)\}}$$



Complicated likelihoods

generalized linear mixed models

GLM: $y_{ij} \mid u_i \sim \exp\{y_{ij}\eta_{ij} - b(\eta_{ij}) + c(y_{ij})\}$

linear predictor: $\eta_{ij} = \mathbf{x}_{ij}^T \beta + \mathbf{z}_{ij}^T \mathbf{u}_i \quad j=1, \dots, n_i; \quad i=1, \dots, m$

random effects: $\mathbf{u}_i \sim N_k(\mathbf{0}, \Sigma)$

log-likelihood:

$$\begin{aligned} \ell(\beta, \Sigma) &= \sum_{i=1}^m \left(\mathbf{y}_i^T \mathbf{X}_i \beta - \frac{1}{2} \log |\Sigma| \right. \\ &\quad \left. + \log \int_{\mathbb{R}^k} \exp\{ \mathbf{y}_i^T \mathbf{Z}_i \mathbf{u}_i - \mathbf{1}_i^T b(\mathbf{X}_i \beta + \mathbf{Z}_i \mathbf{u}_i) - \frac{1}{2} \mathbf{u}_i^T \Sigma^{-1} \mathbf{u}_i \} d\mathbf{u}_i \right) \end{aligned}$$

Ormerod & Wand 2012

... complicated likelihoods

multivariate extremes: example, wind speed at d locations

vector observations: $(X_{1i}, \dots, X_{di}), i = 1, \dots, n$

component-wise maxima: $Z_1, \dots, Z_d; Z_j = \max(X_{j1}, \dots, X_{jn})$

Z_j are transformed (centered and scaled)

joint distribution function:

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

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

likelihood : need the joint derivatives of $V(\cdot)$

combinatorial explosion

Davison et al., 2012

... complicated likelihoods

Ising model:

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

observations: $y_i = \pm 1$; binary property of a node i
in a graph with K nodes

parameter: θ_{jk} measures strength of interaction between
nodes i and j

E is the set of edges between nodes

partition function:

$$Z(\theta) = \sum_{\mathbf{y}} \exp\left(\sum_{(j,k) \in E} \theta_{jk} y_j y_k\right)$$

Davison 2000 §6.2; Ravikumar et al. (2010); Xue et al. (2012)

... complicated likelihoods

$M/G/1$ queue: exponential arrival times, general service times, single server

observations y_i : times between departures from the queue

unobserved variables V_i : arrival time of customer i

model:

- $V_1 \sim \text{Exp}(\theta_3)$
- $V_i | V_{i-1} \sim V_{i-1} + \text{Exp}(\theta_3)$
- $Y_i | X_{i-1}, V_i \sim \text{Uniform}\{\theta_1 + \max(0, V_i - X_{i-1}), \theta_2 + \max(0, V_i - X_{i-1})\}$ $X_i = \sum_{j=1}^i Y_j$ $G = U(\theta_1, \theta_2)$

Likelihood

$$L(\theta; y) = \int \cdots \int f(v_1 | \theta) \prod_{i=1}^n f(v_i | v_{i-1}, \theta) \prod_{i=1}^n f(y_i | v_i, x_{i-1}, \theta) dv_1 \cdots dv_n$$

Heggland & Frigessi, 2004
Fearnhead & Prangle, 2012

What's a poor statistician to do?

- simplify the likelihood
 - composite likelihood
 - variational approximation
 - Laplace approximation to integrals
- change the mode of inference
 - quasi-likelihood
 - indirect inference
- simulate
 - approximate Bayesian computation
 - Markov chain Monte Carlo

Composite likelihood

- also called pseudo-likelihood Besag, 1975
- reduce high-dimensional dependencies by ignoring them
- for example, replace $f(y_{i1}, \dots, y_{ik}; \theta)$ by

pairwise marginal $\prod_{j < j'} f_2(y_{ij}, y_{ij'}; \theta),$ or

conditional $\prod_j f_c(y_{ij} \mid y_{\mathcal{N}(ij)}; \theta)$

- Composite likelihood function

$$CL(\theta; y) \propto \prod_{i=1}^n \prod_{j < j'} f_2(y_{ij}, y_{ij'}; \theta)$$

- Composite ML estimates are consistent, asymptotically normal, not fully efficient Lindsay, 1988; Varin R Firth, 2011

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

- pairwise composite likelihood used to compare the fits of several competing models
- model choice using “CLIC”, an analogue of AIC
$$-2 \log(\widehat{CL}) + \text{tr}(J^{-1}K)$$
- Davison et al. 2012 applied this to annual maximum rainfall at several stations near Zurich
- “fitting max-stable processes to spatial or spatio-temporal block maxima is awkward ... the use of composite likelihoods ... has become widely used” Davison & Huser

Example: Ising model

Ising model:

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

neighbourhood contributions

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

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

- in a Bayesian context, want $f(\beta | y)$
use an approximation $q(\beta)$
- dependence of q on y suppressed
- choose $q(\beta)$ to be
 - simple to calculate
 - close to posterior
- simple to calculate
 - $q(\beta) = \prod q_j(\beta_j)$
 - simple parametric family
- close to posterior: minimize Kullback-Leibler divergence between $q(\cdot)$ and $f(\cdot | y)$

- close to posterior: minimize Kullback-Leibler divergence

$$KL(q \parallel f_{post}) = \int q(\beta) \log\{q(\beta)/f(\beta | y)\} d\beta$$

- equivalent to

$$\max_q \int q(\beta) \log\{f(y, \beta)/q(\beta)\} d\beta$$

- because

$$\log f(y; \theta) \geq \int q(\beta) \log\{f(y, \beta; \theta)/q(\beta)\} d\beta$$

- in a likelihood context

$$\log f(y; \theta) = \log \int f(y | \beta; \theta) f(\beta) d\beta$$

here β represent random effects u , or b , or ...

log-likelihood:

$$\begin{aligned} \ell(\beta, \Sigma) &= \sum_{i=1}^m \left(y_i^T X_i \beta - \frac{1}{2} \log |\Sigma| \right. \\ &\quad \left. + \log \int_{\mathbb{R}^k} \exp\{y_i^T Z_i u_i - \mathbf{1}_i^T b(X_i \beta + Z_i u_i) - \frac{1}{2} u_i^T \Sigma^{-1} u_i\} du_i \right) \end{aligned}$$

variational approx:

$$\begin{aligned} \ell(\beta, \Sigma) &\geq \sum_{i=1}^m \left(y_i^T X_i \beta - \frac{1}{2} \log |\Sigma| \right) \\ &\quad + \sum_{i=1}^m E_{u \sim N(\mu_i, \Lambda_i)} \left(y_i^T Z_i u - \mathbf{1}_i^T b(X_i \beta + Z_i u) - \frac{1}{2} u^T \Sigma^{-1} u - \log\{\phi_{\Lambda_i}(u - \mu_i)\} \right) \end{aligned}$$

simplifies to k one-dim. integrals

-

$$\ell(\beta, \Sigma) \geq \ell(\beta, \Sigma, \mu, \Lambda)$$

- variational estimate:

$$\ell(\tilde{\beta}, \tilde{\Sigma}, \tilde{\mu}, \tilde{\Lambda}) = \arg \max_{\beta, \Sigma, \mu, \Lambda} \ell(\tilde{\beta}, \tilde{\Sigma}, \tilde{\mu}, \tilde{\Lambda})$$

- inference for $\tilde{\beta}, \tilde{\Sigma}$? consistency? asymptotic normality?

Hall, Ormerod, Wand, 2011; Hall et al. 2011

- emphasis on algorithms and model selection

e.g. Tan & Nott, 2013, 2014

- VL: approx $L(\theta; y)$ by a simpler function of θ , e.g. $\prod q_j(\theta)$

- CL: approx $f(y; \theta)$ by a simpler function of y , e.g. $\prod f(y_j; \theta)$

Some Links between Variational Approximation and Composite Likelihoods?

S. Robin

UMR 518 AgroParisTech / INRA Applied Math & Comput. Sc.



MSTGA, Paris, November 22-23, 2012

http://carlit.toulouse.inra.fr/AIGM/pub/Reunion_nov2012/MSTGA-1211-Robin.pdf

Zhang & Schneider 2012 JMLR V22; Grosse 2015 ICML

Indirect inference

- composite likelihood estimator solves $(\partial/\partial\theta) \log CL(\theta; y) = 0$
- solution converges to the true value under conditions ...
- because $E\{(\partial/\partial\theta) \log CL(\theta; y)\} = 0$

- what happens if an estimating equation $g(y; \theta)$ is **biased**?
- $g(y_1, \dots, y_n; \tilde{\theta}_n) = 0; \quad \tilde{\theta}_n \rightarrow \theta^* \quad E g(Y; \theta^*) = 0$

- $\theta^* = \tilde{k}(\theta)$; invertible? $\theta = k(\theta^*) \quad \tilde{k}^{-1} \equiv k$

- **new estimator** $\hat{\theta}_n = k(\tilde{\theta}_n)$
- $k(\cdot)$ is a **bridge** function, connecting wrong value of θ to the right one Yi & R, 2010; Jiang & Turnbull, 2004

- model of interest

$$y_t = G_t(y_{t-1}, x_t, \epsilon_t; \theta), \quad \theta \in \mathbb{R}^d$$

- likelihood is not computable, but
we can simulate from the model
- simple (wrong) model

$$y_t \sim f(y_t | y_{t-1}, x_t; \theta^*), \quad \theta^* \in \mathbb{R}^p$$

- find the MLE in the simple model, $\hat{\theta}^* = \hat{\theta}^*(y_1, \dots, y_n)$, say
- **simulate** from model of interest for some value θ , compute a new MLE in simple model
- ‘good’ values of θ give data that reproduces $\hat{\theta}^*$

- **simulate** samples y_t^m , $m = 1, \dots, M$ at some value θ
- compute $\hat{\theta}^*(\theta)$ from the simulated data

$$\hat{\theta}^*(\theta) = \arg \max_{\theta^*} \sum_m \sum_t \log f(y_t^m | y_{t-1}^m, x_t; \theta^*)$$

- choose θ so that $\hat{\theta}^*(\theta)$ is as close as possible to $\hat{\theta}^*$
- if both model parameters have the same dimension simply invert the ‘bridge function’
- usually not, so minimize some measure of distance between $\hat{\theta}(\beta)$ and $\hat{\theta}$
- estimates of θ are consistent, asymptotically normal, but not efficient

- simulate θ from prior density $\pi(\cdot)$
- simulate data y' from $f(\cdot; \theta)$
- if $y' = y$ then θ is an observation from posterior $\pi(\cdot | y)$
- actually $s(y') = s(y)$ for some set of statistics
- actually $\rho\{s(y'), s(y)\} < \epsilon$ for some distance function $\rho(\cdot)$

Fearnhead & Prangle, 2011

- many variations, using different MCMC methods to select candidate values θ

... approximate Bayesian computation

M/G/1 queue: exponential arrival times, general service times, single server

observations y_i : times between departures from the queue

unobserved variables V_i : arrival time of customer i

model:

- $V_1 \sim \text{Exp}(\theta_3)$
- $V_i | V_{i-1} \sim V_{i-1} + \text{Exp}(\theta_3)$
- $Y_i | X_{i-1}, V_i \sim \text{Uniform}\{\theta_1 + \max(0, V_i - X_{i-1}), \theta_2 + \max(0, V_i - X_{i-1})\}$ $X_i = \sum_{j=1}^i Y_j$
- service time $\sim U(\theta_1, \theta_2)$

ABC: use quantiles of departure times as summary statistics

Indirect Inference: use \bar{y} , $y_{(1)}$, $\hat{\theta}_2$ from steady-state model

Heggland & Frigessi, 2004

Table 7. Mean quadratic losses for various analyses of 50 $M/G/1$ data sets[†]

<i>Method</i>	θ_1	θ_2	θ_3
Comparison	1.1	2.2	0.0013
Comparison + regression	<i>0.020</i>	1.1	<i>0.0013</i>
Semi-automatic ABC	<i>0.022</i>	1.0	<i>0.0013</i>
Semi-automatic predictors	0.024	1.2	0.0017
Indirect inference	0.18	<i>0.42</i>	0.0033

[†]Losses within 10% of the smallest values for that parameter are italicized.

- both methods need a set of parameter values from which to simulate: θ' or θ
- both methods need a set of auxiliary functions of the data $s(y)$ or $\hat{\theta}^*(y)$
- in indirect inference, $\hat{\theta}^*$ is the 'bridge' to the parameters of real interest, θ
- C & K use orthogonal designs based on Hadamard matrices to chose θ'
- and calculate summary statistics focussed on individual components of θ

What's a poor statistician to do?

- simplify the likelihood
 - composite likelihood
 - variational approximation
 - Laplace approximation to integrals
- change the mode of inference
 - quasi-likelihood
 - indirect inference
- simulate
 - approximate Bayesian computation
 - MCMC

Summary

so much to do, so little time!

Summary

- empirical likelihood, weighted likelihood, local likelihood, sieve likelihood, simulated likelihood, ...
- likelihood provides a common set of tools:
 - summary statistics
 - e.g. point estimates and estimates of precision
 - comparison of models
- likelihood puts modelling first
- likelihood puts inference first
- contrast with 'black-box' predictions

Thank You!

