

Likelihood based inference

1. Overview of classical asymptotics
2. Profile likelihood and nuisance parameters NR 2013; 2010
3. p growing with n Portnoy 1984, 5, 8
4. $p > n$: regularization Buhlmann 2013; Taylor et al. 2014
5. approximate likelihoods composite, quasi, empirical, ...

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 $f(y | x; \theta)$, $\theta = (\theta_1, \dots, \theta_p)$
- ▶ **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, $p > n$

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

Examples

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

- ▶ example: clustered binary data

- ▶ latent variable:

$$z_{ir} = \mathbf{x}'_{ir}\beta + b_i + \epsilon_{ir}, \quad b_i \sim N(0, \sigma_b^2), \quad \epsilon_{ir} \sim N(0, 1)$$

- ▶ $r = 1, \dots, n_i$: observations in a cluster/family/school...
 $i = 1, \dots, n$ clusters

- ▶ random effect b_i introduces correlation between observations in a cluster

- ▶ observations: $y_{ir} = 1$ if $z_{ir} > 0$, else 0

- ▶ $Pr(y_{ir} = 1 | b_i) = \Phi(\mathbf{x}'_{ir}\beta + b_i) = p_i$ $\Phi(z) = \int^z \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx$

- ▶ likelihood $\theta = (\beta, \sigma_b)$

$$L(\theta; y) = \prod_{i=1}^n \log \int_{-\infty}^{\infty} \prod_{r=1}^{n_i} p_i^{y_{ir}} (1 - p_i)^{1-y_{ir}} \phi(b_i, \sigma_b^2) db_i$$

- ▶ more general: $z_{ir} = \mathbf{x}'_{ir}\beta + \mathbf{w}'_{ir}b_i + \epsilon_{ir}$

Renard et al. (2004)

... complicated likelihoods

Poisson $f(y_t | \alpha_t; \theta) = \exp(y_t \log \mu_t - \mu_t) / y_t!$

$$\log \mu_t = \beta + \alpha_t$$

autoregression

$$\alpha_t = \phi \alpha_{t-1} + \epsilon_t, \quad \epsilon_t \sim N(0, \sigma^2), \quad |\phi| < 1, \quad \theta = (\beta, \phi, \sigma^2)$$

likelihood

$$L(\theta; y_1, \dots, y_n) = \int \left(\prod_{t=1}^n f(y_t | \alpha_t; \theta) \right) f(\alpha; \theta) d\alpha$$

$L_{approx}(\theta; y)$ via Laplace with some refinements

Davis & Yau, 2011

... 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

Restricted Boltzmann machine:

$$f(\mathbf{v}; \theta) = \frac{1}{Z(\theta)} \exp\left\{ \sum_h (\mathbf{h}^T \mathbf{W} \mathbf{v} + \alpha^T \mathbf{h} + \beta^T \mathbf{v}) \right\}, \quad \theta = (\mathbf{W}, \alpha, \beta)$$

observations: v_1, \dots, v_n , independent $\sim f(\mathbf{v}; \theta)$; hidden units h

complete data likelihood

$$f(\mathbf{v}, \mathbf{h}; \theta) = \frac{1}{Z(\theta)} \exp\{f(\mathbf{v}, \mathbf{h}; \theta)\}$$

partition function: $Z(\theta) = \sum_{\mathbf{v}, \mathbf{h}} \exp\{f(\mathbf{v}, \mathbf{h}; \theta)\}$

Jan 30 MZ slides; GL slides

Why likelihood?

- ▶ makes probability modelling central
- ▶ emphasizes the inverse problem of reasoning from y to θ or $f(\cdot)$
- ▶ suggested by Fisher as a measure of plausibility

Royall, 1997

$L(\hat{\theta})/L(\theta) \in (1, 3)$ very plausible;

$L(\hat{\theta})/L(\theta) \in (3, 10)$ implausible;

$L(\hat{\theta})/L(\theta) \in (10, \infty)$ very implausible

- ▶ converts a 'prior' probability $\pi(\theta)$ to a posterior $\pi(\theta | y)$ via Bayes' Theorem
- ▶ provides a conventional set of summary quantities: maximum likelihood estimator, score function, ...
- ▶ leading to approximate pivotal functions, based on normal distribution
- ▶ basis for comparison of models, using AIC or BIC

Derived quantities

- ▶ maximum likelihood estimator

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

- ▶ observed Fisher information

$$j(\hat{\theta}) = - \partial^2 \ell(\theta) / \partial \theta^2 |_{\hat{\theta}}$$

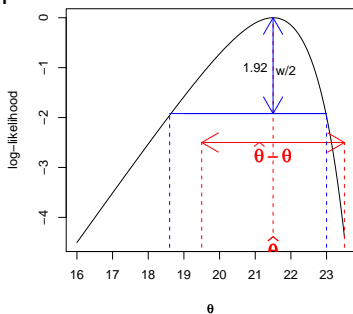
- ▶ efficient score function

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

$$\ell'(\hat{\theta}) = 0$$

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

log-likelihood function



assuming enough regularity

y_1, \dots, y_n independent

Limiting results

no nuisance parameters

$$\ell'(\theta)^\top j^{-1}(\hat{\theta}) \ell'(\theta) \xrightarrow{\mathcal{L}} \chi_p^2,$$

$$q(\theta) \equiv (\hat{\theta} - \theta)^\top j(\hat{\theta})(\hat{\theta} - \theta) \xrightarrow{\mathcal{L}} \chi_p^2,$$

$$w(\theta) \equiv 2\{\ell(\hat{\theta}) - \ell(\theta)\} \xrightarrow{\mathcal{L}} \chi_p^2$$

Approximate pivots $p = 1$

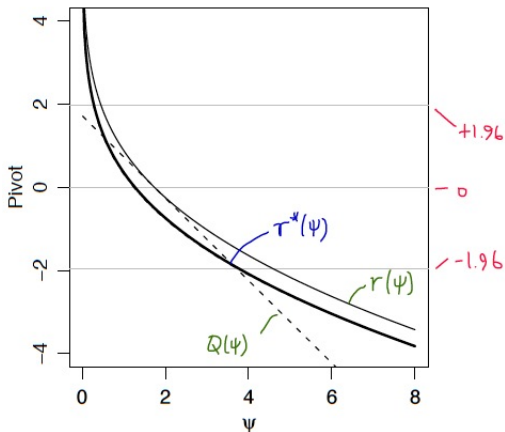
$$s(\theta) \equiv \ell'(\theta) j^{-1/2}(\hat{\theta}) \sim N(0, 1)$$

$$q(\theta) \equiv (\hat{\theta} - \theta) j^{1/2}(\hat{\theta}) \sim N(0, 1)$$

$$r(\theta) \equiv \pm \sqrt{2\{\ell(\hat{\theta}) - \ell(\theta)\}^{1/2}} \sim N(0, 1)$$

... approximate pivots

scalar parameter of interest



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

- ▶ $\hat{\lambda}_\psi$ constrained maximum likelihood estimator
- ▶ profile log-likelihood $\ell_p(\psi) = \ell(\psi, \hat{\lambda}_\psi)$

$$r_e(\psi; \mathbf{y}) = (\hat{\psi} - \psi) j_p^{1/2}(\hat{\psi}) \sim N(0, 1)$$

$$r(\psi; \mathbf{y}) = \pm \sqrt{2\{\ell_p(\hat{\psi}) - \ell_p(\psi)\}} \sim N(0, 1)$$

$$\pi_m(\psi | \mathbf{y}) \sim N\{\hat{\psi}, j_p^{-1/2}(\hat{\psi})\}$$

$$j_p(\psi) = -\ell_p''(\psi); \text{ profile information}$$

- ▶ treat profile log-likelihood as a one-parameter log-likelihood

The problem with profiling

- ▶ $\ell_p(\psi) = \ell(\psi, \hat{\lambda}_\psi)$ used as a ‘regular’ likelihood, with the usual asymptotics
- ▶ neglects errors in the estimation of the nuisance parameter
- ▶ can be very large when there are many nuisance parameters
- ▶ **example:** $Y \sim N(X\beta, \sigma^2 I)$, $\hat{\sigma}^2 = (y - X\hat{\beta})^T(y - X\hat{\beta})/n$
- ▶ badly biased if $\dim(\beta)$ large relative to n
- ▶ easy fix: $\tilde{\sigma}^2 = (y - X\hat{\beta})^T(y - X\hat{\beta})/(n - p)$
- ▶ **example:** $Y_{ij} \sim N(\mu_i, \sigma^2), j = 1, \dots, n; i = 1, \dots, p$
- ▶ $\hat{\sigma}^2 \xrightarrow{p} \frac{n-1}{n} \sigma^2$ as $p \rightarrow \infty, n$ fixed

Neyman & Scott, 1948

Reminder: deriving limit results

- ▶ $\ell'(\hat{\theta}; \mathbf{y}) = \mathbf{0} = \ell'(\theta; \mathbf{y}) + (\hat{\theta} - \theta)\ell''(\theta; \mathbf{y})$
- ▶ $\ell'(\theta; \mathbf{y})\{-\ell''(\theta; \mathbf{y})\}^{-1} \doteq \hat{\theta} - \theta$
- ▶ $\underbrace{\ell'(\theta; \mathbf{y})}_{\xrightarrow{\mathcal{L}} N(\mathbf{0}, i(\theta))} \underbrace{\{-\ell''(\theta; \mathbf{y})\}^{-1}}_{\xrightarrow{P} i^{-1}(\theta)} \implies (\hat{\theta} - \theta) \xrightarrow{\mathcal{L}} N(\mathbf{0}, i^{-1}(\theta))$
 $i(\theta) = E\{j(\theta)\} = E\{-\ell''(\theta)\} = \text{cov}\{\ell'(\theta)\}$
- ▶ **M estimator:** $\tilde{\theta}_\rho = \text{argmin}_\theta \Sigma \rho(\mathbf{y}_i; \theta)$
- ▶ Solution $\Sigma \psi(\mathbf{y}_i; \tilde{\theta}_\rho) = \mathbf{0}$, $\psi(\mathbf{y}; \theta) = \partial \rho(\mathbf{y}_i; \theta) / \partial \theta$
- ▶ $\tilde{\theta}_\rho \xrightarrow{\mathcal{L}} N\{\mathbf{0}, \mathbf{G}^{-1}(\theta)\}$ $E\{\psi(\mathbf{Y}; \theta)\} = \mathbf{0}$
- ▶ $\mathbf{G}(\theta) = E\{-\partial \psi(\mathbf{Y}; \theta) / \partial \theta\} [\text{cov}\{\psi(\mathbf{Y}; \theta)\}]^{-1} E\{-\partial \psi(\mathbf{Y}; \theta) / \partial \theta\}$

big data asymptotics

▶ Neyman-Scott problems: n fixed, $p \rightarrow \infty$

▶ Donoho $n, p \rightarrow \infty$, $p/n \rightarrow \beta < \infty$

▶ likelihood results $n, p \rightarrow \infty$, $p^2/n \rightarrow \beta < \infty$

Portnoy, 1984, 5, 8

▶ Laplace approx $n, p \rightarrow \infty$, $p = o(n^{1/3})$

Shun & McCullagh, 1995

▶ $p > n$: regularize

▶ lasso

$$\operatorname{argmin}_{\beta} (y - X\beta)^T (y - X\beta) - \lambda \sum_j |\beta_j|$$

no intercept

▶ ridge regression

$$\operatorname{argmin}_{\beta} (y - X\beta)^T (y - X\beta) - \lambda \sum_j \beta_j^2$$

- ▶ Model: $y_i = \mathbf{x}_i^T \beta + Z_i$, $i = 1, \dots, n$ independent
- ▶ M -estimation:

$$\sum_{i=1}^n \mathbf{x}_i \psi(y_i - \mathbf{x}_i^T \hat{\beta}) = 0 \quad (1)$$

- ▶ **result:** if ψ is monotone, and $p \log(p)/n \rightarrow 0$, and conditions on X , then

there is a solution of (1) satisfying $\|\hat{\beta} - \beta\|^2 = O(p/n)$

- ▶ “rows of X behave like a sample from a distribution in \mathbb{R}^p ”
- ▶ if $p^{3/2} \log n/n \rightarrow 0$, then

$$\max |\mathbf{x}_i^T (\hat{\beta} - \beta)| \xrightarrow{P} 0$$

- ▶ and

$$\mathbf{a}_n^T (\hat{\beta} - \beta) \xrightarrow{\mathcal{L}} N(0, \sigma^2)$$

$$\sigma^2 = \mathbf{a}_n^T (X^T X)^{-1} \mathbf{a}_n E \psi^2(Z) / \{E \psi'(Z)\}^2$$

- ▶ Model: $y_i \sim \exp\{\theta^T y - \psi(\theta)\}$, $i = 1, \dots, n$ independent; $p = p_n$
- ▶ maximum likelihood estimate $\psi'(\hat{\theta}_n) = \bar{y}_n$
- ▶ under conditions on the eigenvalues of $\psi''(\theta)$ and moment conditions on y ,
Fisher information matrix

$$\|\hat{\theta}_n - \theta_n\|^2 \leq c \frac{p}{n}, \text{ in probability,}$$

▶

$$\|\hat{\theta} - \theta - \bar{y}\| = O_p(p/n) \text{ if } p/n \rightarrow 0,$$

▶ $p^{3/2}/n \rightarrow 0$:

$$\sqrt{n} a_n^T (\hat{\theta} - \theta) \xrightarrow{\mathcal{L}} N(0, 1),$$

likelihood ratio test of simple hypothesis asymptotically χ_p^2

- ▶ “asymptotic approximations are trustworthy if $p^{3/2}/n$ is small, but may be very wrong if p^2/n is not small”
- ▶ MLE ‘will tend to be’ consistent if $p/n \rightarrow 0$

cf. also El Karoui et al., 2013, PNAS

$$\hat{\beta}_{\text{Lasso}} = \operatorname{argmin}_{\beta} (\|y - X\beta\|_2^2 + \lambda \|\beta\|_1)$$

- ▶ prediction: $\|X(\hat{\beta}_{\text{Lasso}} - \beta^0)\|_2^2/n$ 'small'
- ▶ estimation: $\|\hat{\beta}_{\text{Lasso}} - \beta^0\|_q$, $q \in 1, 2$ 'small'
- ▶ selection: $\mathbb{P}(\hat{S} = S_0)$ 'large' S_0 is the 'active set':
 $\{j : \beta_j^0 \neq 0\}$
- ▶ under restricted eigenvalue conditions on X , can get results like

$$\|\hat{\beta}_{\text{Lasso}} - \beta^0\|_1 = O_p(s_0 \sqrt{\log(p)/n}), \quad \lambda \approx \sqrt{\log(p)/n}$$

- ▶ what about estimated standard errors for $\hat{\beta}_{\text{Lasso}}$?
- ▶ Bühlmann, 2013: the **ridge** regression estimate

$$\hat{\beta}_R = \operatorname{argmin}_{\beta} (\|y - X\beta\|_2^2 + \lambda \|\beta\|_2),$$

$$\hat{\beta}_{\text{Lasso}} = \operatorname{argmin}_{\beta} (\|y - X\beta\|_2^2 + \lambda \|\beta\|_1)$$

First, we define some needed quantities. Let A be the active set just before λ_k , and suppose that predictor j enters at λ_k . Denote by $\hat{\beta}(\lambda_{k+1})$ the solution at the next knot in the path λ_{k+1} , using predictors $A \cup \{j\}$. Finally, let $\tilde{\beta}_A(\lambda_{k+1})$ be the solution of the lasso problem using only the active predictors X_A , at $\lambda = \lambda_{k+1}$. To be perfectly explicit,

$$(4) \quad \tilde{\beta}_A(\lambda_{k+1}) = \operatorname{argmin}_{\beta_A \in \mathbb{R}^{|A|}} \frac{1}{2} \|y - X_A \beta_A\|_2^2 + \lambda_{k+1} \|\beta_A\|_1.$$

We propose the *covariance test statistic* defined by

$$(5) \quad T_k = (\langle y, X \hat{\beta}(\lambda_{k+1}) \rangle - \langle y, X_A \tilde{\beta}_A(\lambda_{k+1}) \rangle) / \sigma^2.$$

$$T_k \xrightarrow{\mathcal{L}} \operatorname{Exp}(1)$$

Taylor et al. 2014

Likelihood in complex models

- ▶ 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

Composite likelihood

- ▶ also called pseudo-likelihood
- ▶ 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

Besag, 1975; Lindsay, 1988

- ▶ Likelihood

$$L(\theta; y_1, \dots, y_n) = \int \left(\prod_{t=1}^n f(y_t | \alpha_t; \theta) \right) f(\alpha; \theta) d\alpha$$

- ▶ Composite likelihood

$$CL(\theta; y_1, \dots, y_n) = \prod_{t=1}^{n-1} \int \int f(y_t | \alpha_t; \theta) f(y_{t+1} | \alpha_{t+1}; \theta) f(\alpha_t, \alpha_{t+1}; \theta) d\alpha_t d\alpha_{t+1}$$

- ▶ consecutive pairs
- ▶ Time-series asymptotic regime one vector y of increasing length
- ▶ Composite ML estimator still consistent, asymptotically normal, estimable asymptotic variance
- ▶ Efficient, relative to a Laplace-type approximation
- ▶ Surprises: AR(1), fully efficient; MA(1), poor; ARFIMA(0,d,0), ok

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

- ▶ 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:

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

neighbourhood contributions

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

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

$$\max_{\theta} \left\{ \sum_{i=1}^n \ell_j(\theta; y^{(i)}) - \sum_j \sum_k P_{\lambda}(|\theta_{jk}|) \right\}$$

Xue et al., 2012

Ravikumar et al., 2010

Quasi-likelihood

- ▶ simplify the model



$$E(y_i; \theta) = \mu_i(\theta); \quad \text{Var}(y_i; \theta) = \phi \nu_i(\theta)$$

- ▶ consistent with generalized linear models
- ▶ example: over-dispersed Poisson responses
- ▶ PQL uses this construction, but with random effects

Molenberghs & Verbeke, Ch. 14

- ▶ why does it work?
- ▶ score equations are the same as for a 'real' likelihood

hence unbiased

- ▶ derivative of score function equal to variance function

special to GLMs

Indirect inference

- ▶ composite likelihood estimators are consistent
under conditions ...
- ▶ because $\log CL(\theta; y) = \sum_{i=1}^n \sum_{j < j'} \log f(y_j, y_{j'}; \theta)$
- ▶ derivative w.r.t. θ has expected value 0

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

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

- ▶ **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 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
- ▶ use simulated samples from model of interest to find the 'best' β
- ▶ 'best' θ gives data that reproduces $\hat{\theta}^*$

Shalizi, 2013

- ▶ 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 $p = d$ simply invert the ‘bridge function’
- ▶ usually $p > d$
 - ▶ $\hat{\theta}_1 = \arg \min_{\theta} \{ \hat{\theta}^*(\theta) - \hat{\theta} \}^T W \{ \hat{\theta}^*(\theta) - \hat{\theta} \}$
 - ▶ $\hat{\beta}_2 = \arg \min_{\theta} (\sum_t \log f(y_t | y_{t-1}, x_t, \hat{\theta}^*(\theta)) - \sum_t \log f(y_t | y_{t-1}, x_t, \hat{\theta}))$
- ▶ estimates of θ are consistent, asymptotically normal, but not efficient

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

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
- ▶ MCMC estimation of log-likelihood function

Geyer & Thompson, 1992

cond. comp. likelihood poor for Ising model

Okabayashi et al., 2011

- ▶ 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

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

- ▶ close to posterior:

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

- ▶ equivalent to

best LB for marginal $f(y)$

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

- ▶ in a likelihood context $\log f(y; \theta) = \log \int f(y | \beta; \theta) f(\beta) d\beta$

$$= \int q(\beta) \log\{f(y, \beta; \theta)/q(\beta)\} d\beta + KL(q || f_{post})$$

- ▶

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

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

log-likelihood:

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

variational approx:

$$\begin{aligned}
 \ell(\beta, \Sigma) &\geq \sum_{i=1}^m \left(y_i^T X_i \beta - \frac{1}{2} \log |\Sigma| \right) \\
 &\quad + \sum_{i=1}^m E_{u_i \sim \mathcal{N}(\mu_i, \Lambda_i)} \left(y_i^T Z_i u_i - 1_i^T b(X_i \beta + Z_i u_i) - \frac{1}{2} u_i^T \Sigma^{-1} u_i - \log\{\phi_{\Lambda_i}(u - \mu_i)\} \right) \\
 &\equiv \ell(\beta, \Sigma, \mu, \Lambda) \quad \text{simplifies to } k \text{ one-dim. integrals}
 \end{aligned}$$



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

- ▶ **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. $\prod f(y_j; \theta)$

Laplace approximation

$$\ell(\theta; \mathbf{y}) = \log \int f(\mathbf{y} | \mathbf{b}; \theta) g(\mathbf{b}) d\mathbf{b} = \log \int \exp\{Q(\mathbf{b}, \mathbf{y}, \theta)\} d\mathbf{b}, \text{ say}$$

$$\ell_{Lap}(\theta; \mathbf{y}) = Q(\tilde{\mathbf{b}}, \mathbf{y}, \theta) - \frac{1}{2} \log |Q''(\tilde{\mathbf{b}}, \mathbf{y}, \theta)| + c$$

using Taylor series expansion of $Q(\cdot, \mathbf{y}, \theta)$ about $\tilde{\mathbf{b}}$

simplification of the Laplace approximation leads to PQL:

$$\ell_{PQL}(\theta, \mathbf{b}; \mathbf{y}) = \log f(\mathbf{y} | \mathbf{b}; \theta) - \frac{1}{2} \mathbf{b}^T \Sigma^{-1} \mathbf{b}$$

Breslow & Clayton, 1993

to be jointly maximized over \mathbf{b} and θ

and parameters in Σ

PQL can be viewed as linearizing $E(\mathbf{y})$ and then using results for linear mixed models

Molenberghs & Verbeke, 2006

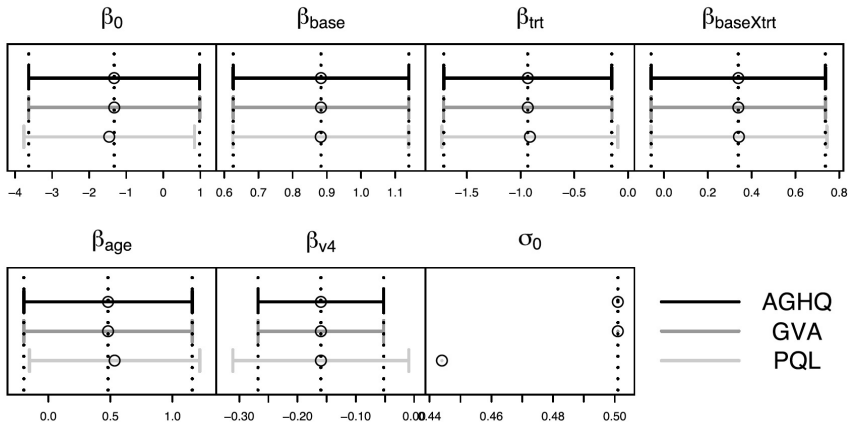


Figure 2. Point estimates and approximate 95% confidence intervals based on each of AGHQ, GVA, and PQL for the *Epilepsy* data random intercept model. The vertical dotted lines correspond to the AGHQ values.

implemented in `lme4` as `glmer`, in `MASS` as `glmmPQL`

Ormerod & Wand, 2012

References

- Besag, J. (1975). Statistical analysis of non-lattice data. *The Statistician* **24**, 179–195.
- Breslow, N.E. & Clayton, D. G. (1993). Approximate inference in generalised linear mixed models. *J. Am. Statist. Assoc.* **88**, 9–25.
- Bühlmann, P. (2013). Statistical significance in high-dimensional linear models. *Bernoulli* **19**, 1212 – 1242.
- Bühlmann, P., Kalisch, M. and Meier, L. (2014). High-dimensional statistics with a view toward applications in biology. *Annual Review of Statistics and its Applications* **1**, 255–278.
- Cox, D.R. & Kartsonaki, C. (2012). The fitting of complex parametric models. *Biometrika* **99**, 741–747.
- Davis, R. & Yau, C.Y. (2011). Comments on pairwise likelihood in time series. *Statistica Sinica* **21**, 255–277.
- Davison, A.C., (2012). Statistical modeling of spatial extremes. *Statistical Science* **27**, 161–186.
- Davison, A.C. & Huser, R.(2015). Statistics of Extremes *Annual Reviews* **2**, to appear.
- El Karoui, N., Bean, D., Bickel, P.J., Lim, C. and Yu, B. (2013). On robust regression with hig-dimensional predictors. *PNAS* **110**, 14557 – 14562.
- Fearnhead, P. & Prangle, (2012). Approximate liekelihood methods for estimating local recombination rates *J. R. Statist. Soc. B* **64**, 657–680.
- Geyer, C. & Thompson, E.A. (1992). Constrained MC maximum likelihood... *J. R. Statist. Soc. B* **54**, 657–699.
- Jiang, W. & Turnbull, B. (2004). The indirect methods ... *Statistical Science* **19**, 239–263 .
- Lindsay, B. (1988). Composite likelihood methods. *Contemp. Math.* **80**, 220–239.
- Lockhart, R., Taylor, J., Tibshirani, R.J. and Tibshirani, R. (2014). A significance test for the lasso. *Ann. Statist.* **42**, 413 – 468.
- Marin, J.-M. et al. (2010). Approximate Bayesian computational methods. *Stat. & Computing* **22**, 1167–1180.
- Molenberghs, G. & Verbeke, G. (2006). *Discrete Longitudinal Data* Springer, New York.

... references

- Okabayashi, X. Johnson, L. & Geyer, C.J. (2011). Extending pseudo-likelihood *Statistica Sinica* **21**, 331–347.
- Ormerod, & Wand, M. (2012). Gaussian variational approximate inference... *J Comp Graph Statist* **21**, 2–17.
- Ormerod, & Wand, M. (2010). Explaining variational approximations. *Am. Stat.* **64**, 140–153.
- Portnoy, S. (1984). Asymptotic behaviour of M -estimators of p regression parameters when p^2/n is large. I. Consistency. *Ann. Statist.* **12**, 1298 – 1309.
- Portnoy, S. (1985). Asymptotic behaviour of M -estimators of p regression parameters when p^2/n is large. II. Normal approximation. *Ann. Statist.* **13**, 1403 – 1417.
- Portnoy, S. (1988). Asymptotic behaviour of likelihood methods for exponential families when the number of parameters tends to infinity. *Ann. Statist.* **16**, 356–366.
- Ravikumar et al. (2010). High-dimensional Ising model selection... *Ann. Statist.* **38**, 1287–1319.
- Reid, N. (2013). Aspects of likelihood inference. *Bernoulli* **19**, 1404–1418.
- Reid, N. (2010). Likelihood inference. *Wiley Interdisciplinary Reviews in Computational Statistics*, **5**, 517–525.
- Renard, D. Molenberghs, G. and Geys, H. (2004). A pairwise likelihood approach to estimation in multilevel probit models. *Comp. Stat. Data. Anal.* **44**, 649–667.
- Royall, R.J. (1997). *Statistical Evidence*.... Chapman & Hall, London.
- Shalizi, C. (2013). Notebooks. [indirect inference](#)
- Shun, Z. & McCullagh, P. (1995). Laplace approximation ... *J. R. Statist. Soc. B* **57**, 749–760.
- Smith, A.A. (2008). Indirect inference. in *New Palgrave Dictionary of Economics* 2nd ed.
- Taylor, J., Lockhart, R., Tibshirani, R.J. and Tibshirani, R. (2014). Exact post-selection inference for forward stepwise and least angle regression. <http://arxiv.org/pdf/1401.3889v4.pdf>
- Titterton, D.M. (2006). Bayesian methods for neural networks ... *Statistical Science* **19**, 128–139.
- Xue, L., Zou, H. & Cai, T. (2012). Nonconcave penalized composite conditional likelihood... *Ann. Statist.* **40**, 1403–1429.
- Yi, G. & Reid, N. (2010). A note on misspecified estimating equations. *Statistica Sinica* **20**, 1749–1769.