

Data Uncertainty, MCML and Sampling Density

Graham Byrnes

International Agency for Research on Cancer

27 October 2015

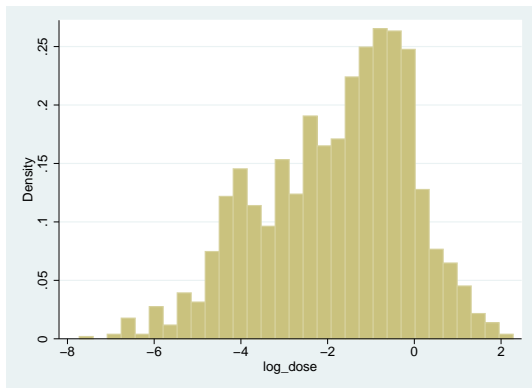
Outline...

- ▶ Correlated Measurement Error
- ▶ Maximal Marginal Likelihood
- ▶ Monte Carlo Maximum Likelihood
- ▶ Sampling Density

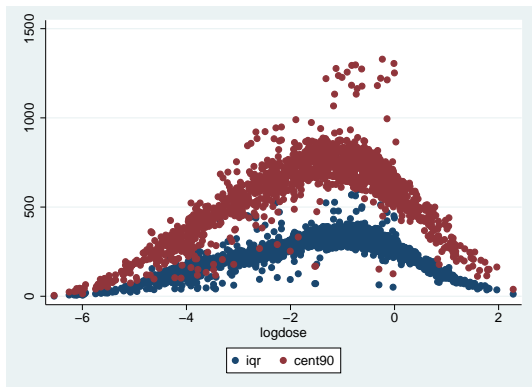
Example Data: Chernobyl Thyroid doses

- ▶ Range of doses is large;
- ▶ Dose rank varies dramatically by draw;
- ▶ Also rank within matched (risk) set.
- ▶ EPI-CT estimates still not ready

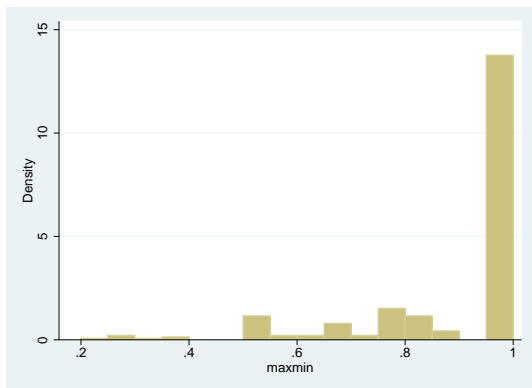
Log dose (central estimate)



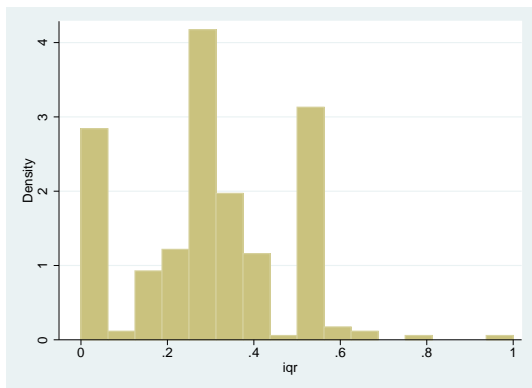
Rank variation



Within-set rank of case: max-min



Within-set rank of case: IQR



Berkson Error

Here $X = Z - \epsilon$ and $\langle Z, \epsilon \rangle = 0$, so $\langle X, \epsilon \rangle = -\tau^2$.
Also $\langle Y, \epsilon \rangle = \langle Y, X \rangle \langle X, \epsilon \rangle / \langle X, X \rangle$, so we have

$$\begin{aligned}\hat{\beta} &= \frac{\langle Y, Z \rangle}{\langle Z, Z \rangle} = \frac{\langle Y, X \rangle (1 - \tau^2 / \langle X, X \rangle)}{\langle X, X \rangle - \tau^2} \\ &= \frac{\langle Y, X \rangle}{\langle X, X \rangle}\end{aligned}$$

The point estimate is unbiased but the information matrix $\langle Z, Z \rangle = \langle X, X \rangle - \tau^2$ is biased downwards.

Marginal Likelihood

- ▶ Expected value of score is still zero
- ▶ Covariance matrix is still estimated by $-H^{-1}$
- ▶ Normality of estimators **does not** hold.

Approximation

After summing over observations and integrating, the first β derivative vanishes

$$\begin{aligned}\mathcal{L} &= \log \int \exp \left(\sum_i (A_i + (x_i - \tilde{x}_i)B_i + (x_i - \tilde{x}_i)^2 C_i) \right) dP(x|z) \\ &= \sum_i A_i + \frac{1}{2} B^T \text{Var}_P(X) B + \text{diag}(\text{Var}_P(X)) \cdot C \\ &\quad + \log \int O((x - \tilde{x})^3) dP(x|z).\end{aligned}$$

Approximation

where

$$\sum_i A_i = \mathcal{L}(Y, \tilde{x}, \hat{\beta}_{\sim}) + \frac{1}{2}(\beta - \hat{\beta}_{\sim})^2 \mathcal{L}_{,\beta\beta}(Y, \tilde{x}, \hat{\beta}_{\sim});$$

$$\sum_i B_i = \mathcal{L}_{,x}(Y, \tilde{x}, \hat{\beta}_{\sim}) + (\beta - \hat{\beta}_{\sim}) \mathcal{L}_{,x\beta}(Y, \tilde{x}, \hat{\beta}_{\sim});$$

$$\sum_i C_i = \mathcal{L}_{,xx}(Y, \tilde{x}, \hat{\beta}_{\sim}).$$

Approximation

- ▶ A is the second order expansion of the imputed likelihood, the limiting where $P(X|z)$ is concentrated on a single exposure vector.
- ▶ If both $\mathcal{L}_{,x}(Y, \tilde{x}, \hat{\beta}_{\sim})$ and $\mathcal{L}_{,x\beta}(Y, \tilde{x}, \hat{\beta}_{\sim})$ are non-zero, then $\frac{\partial}{\partial \beta} \mathcal{L}_M(Y, z)|_{\beta=\hat{\beta}_{\sim}} \neq 0$.
- ▶ Consequently, MCML does not give the same point estimate as regression calibration even if the errors are uncorrelated: the estimate “pulls” x toward a better fitting value.

Monte-Carlo Marginal Likelihood

Implausible that we can evaluate the marginal likelihood analytically. Instead, make uniform draws from $P(X)$ and use a Monte-Carlo approximation:

$$\begin{aligned}\mathcal{L}_{MC}(Y, z; \beta) &= \frac{1}{s} \sum_k^s \mathcal{L}(Y, x^k; \beta) \\ &= \frac{1}{s} \sum_k^s e^{\sum_i L(Y_i, x_i^k; \beta)}.\end{aligned}$$

Shotgun?

- ▶ Cardis *et. al.* evaluated the likelihood on a grid of values of β , then did a second stage with finer spacing.
- ▶ Coded in Stata, extremely slow. Recoded in Fortran, still slow.
- ▶ The likelihood was skewed around the maximum. Almost symmetric using log-dose.
- ▶ Fearn *et. al.* were obliged to also use a grid, since maximum of sum is not the sum of the maxima...

Monte Carlo Likelihood: issues

- ▶ $E(D \ln \mathcal{L}_{MC}) = 0$, if it exists. Newton-Raphson code starting from the central estimate converges reliably on simulated data after removing extreme dose estimates ($> 10\text{Gy}$ to 98Gy);
- ▶ Log-Likelihood is not the sum of IID contributions, so CLT does not apply, so quadratic approximation is not guaranteed;
- ▶ LRT is “generally considered” more reliable than Wald test in such circumstances;
- ▶ But what is the null distribution? Not addressed in any of the above papers.

Sampling Density: test case

Try estimating the area under the standard normal density

$$I_s = \frac{2a}{s} \sum_{j=1}^s \phi(x_j), \quad x_j \sim U(-a, a)$$

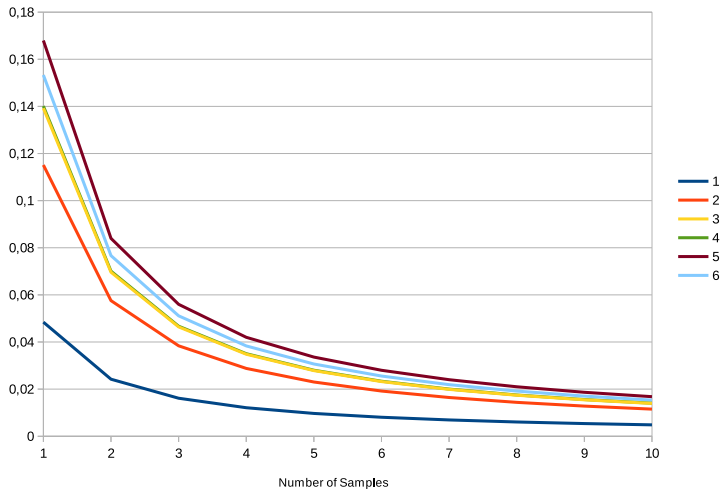
Unbiased, but variance depends on a and s

$$\begin{aligned} \text{Var}(I_s) &= \frac{2a}{s} \left(\int_{-a}^a \phi^2(x) dx - (EI)^2 \right) \\ &= \frac{a}{s} \left(\frac{2\Phi(a\sqrt{2}) - 1}{\sqrt{\pi}} - (2\Phi(a) - 1)^2 \right). \end{aligned}$$

Coeff of Variation

Determined by number of samples under the “fat part” of the density

CoV by domain width



Integral Transform View

- ▶ Can think of the integral over $P(x)$ as an iterated integral transform (approximately Laplace)
- ▶ Maps from density function on $x \in \mathbb{R}^n$ to posterior on $(\beta_1, \dots, \beta_n)$
- ▶ Then we restrict to the sub-domain $\beta_1 = \dots = \beta_n = \beta$.

Sample scaling

- ▶ In EPI-CT where $n \approx 10^6$ this could easily result in retaining finally only a single draw from the dose-set distribution.
- ▶ If this proportional loss remains constant, we would require the number of initial draws to increase exponentially with n .

Close our eyes?

- ▶ It may be tempting to say that worrying about measurement error is too difficult for large n .
- ▶ However for large n , $\text{Var}(\hat{\beta}) \rightarrow 0$ if we ignore measurement error. With error, it does not.
- ▶ Measurement error is the dominant source of uncertainty for sufficiently large n .

Thanks

- ▶ Ausra Kesminiene
- ▶ Deukwoo Kwon
- ▶ Elizabeth Cardis