

# On the Behavior of Marginal and Conditional Akaike Information Criteria in Linear Mixed Models

Thomas Kneib

Department of Mathematics  
Carl von Ossietzky University Oldenburg

Sonja Greven

Department of Biostatistics  
Johns Hopkins University

# Outline

- Akaike Information Criterion
- Linear Mixed Models
- Marginal Akaike Information Criterion
- Conditional Akaike Information Criterion
- Application: Childhood Malnutrition in Nigeria

# Akaike Information Criterion

- Most commonly used model choice criterion for **comparing parametric models**.

- Definition:

$$\text{AIC} = -2l(\hat{\boldsymbol{\psi}}) + 2k.$$

where  $l(\hat{\boldsymbol{\psi}})$  is the log-likelihood evaluated at the maximum likelihood estimate  $\hat{\boldsymbol{\psi}}$  for the unknown parameter vector  $\boldsymbol{\psi}$  and  $k = \dim(\boldsymbol{\psi})$  is the number of parameters.

- Properties:

- Compromise between **model fit** and **model complexity**.
- Allows to compare non-nested models.
- Selects rather too many than too few variables in variable selection problems.

- Data  $\mathbf{y}$  generated from a **true underlying model** described in terms of density  $g(\cdot)$ .
- Approximate the true model by a parametric class of models  $f_\psi(\cdot) = f(\cdot; \psi)$ .
- Measure the discrepancy between a model  $f_\psi(\cdot)$  and the truth  $g(\cdot)$  by the **Kullback-Leibler distance**

$$\begin{aligned} K(f_\psi, g) &= \int [\log(g(\mathbf{z})) - \log(f_\psi(\mathbf{z}))] g(\mathbf{z}) d\mathbf{z} \\ &= \mathbf{E}_z [\log(g(\mathbf{z})) - \log(f_\psi(\mathbf{z}))]. \end{aligned}$$

where  $\mathbf{z}$  is an independent replicate following the same distribution as  $\mathbf{y}$ .

- Decision rule: Out of a sequence of models, **choose the one that minimises  $K(f_\psi, g)$** .

- In practice, the parameter  $\psi$  will have to be **estimated as  $\hat{\psi}(\mathbf{y})$**  for the different models.
- To focus on average properties not depending on a specific data realisation, minimise the **expected Kullback-Leibler distance**

$$\mathbf{E}_{\mathbf{y}}[K(f_{\hat{\psi}(\mathbf{y})}, g)] = \mathbf{E}_{\mathbf{y}}[\mathbf{E}_{\mathbf{z}} [\log(g(\mathbf{z})) - \log(f_{\hat{\psi}(\mathbf{y})}(\mathbf{z}))]]]$$

- Since  $g(\cdot)$  does not depend on the data, this is equivalent to minimising

$$-2 \mathbf{E}_{\mathbf{y}}[\mathbf{E}_{\mathbf{z}}[\log(f_{\hat{\psi}(\mathbf{y})}(\mathbf{z}))]] \quad (1)$$

(the expected **relative** Kullback-Leibler distance).

- The best available estimate for (1) is given by

$$-2 \log(f_{\hat{\psi}(\mathbf{y})}(\mathbf{y})).$$

- While (1) is a **predictive quantity** depending on both the data  $\mathbf{y}$  and an independent replication  $\mathbf{z}$ , the density and the parameter estimate are **evaluated for the same data**.  
  
 $\Rightarrow$  **Introduce a correction term.**

- Let  $\tilde{\psi}$  denote the parameter vector minimising the Kullback-Leibler distance.
- Then

$$\begin{aligned} \text{AIC} = & -2 \log(f_{\hat{\psi}(\mathbf{y})}(\mathbf{y})) + 2 \mathbb{E}_{\mathbf{y}}[\log(f_{\hat{\psi}(\mathbf{y})}(\mathbf{y})) - \log(f_{\tilde{\psi}}(\mathbf{y}))] \\ & + 2 \mathbb{E}_{\mathbf{y}}[\mathbb{E}_{\mathbf{z}}[\log(f_{\tilde{\psi}}(\mathbf{z})) - \log(f_{\hat{\psi}(\mathbf{y})}(\mathbf{z}))]] \end{aligned}$$

is unbiased for (1).

- Consider the **regularity conditions**
  - $\psi$  is a  $k$ -dimensional parameter with parameter space  $\Psi = \mathbb{R}^k$  (possibly achieved by a change of coordinates).
  - $\mathbf{y}$  consists of independent and identically distributed replications  $y_1, \dots, y_n$ .
- In this case, the AIC simplifies since

$$2 \left[ \log(f_{\hat{\psi}(\mathbf{y})}(\mathbf{y})) - \log(f_{\tilde{\psi}}(\mathbf{y})) \right] \stackrel{a}{\sim} \chi_k^2,$$

$$2 \mathbb{E}_{\mathbf{z}} \left[ \log(f_{\tilde{\psi}}(\mathbf{z})) - \log(f_{\hat{\psi}(\mathbf{y})}(\mathbf{z})) \right] \stackrel{a}{\sim} \chi_k^2$$

and therefore an (asymptotically) unbiased estimate for (1) is given by

$$\text{AIC} = -2 \log(f_{\hat{\psi}(\mathbf{y})}(\mathbf{y})) + 2k.$$



# Linear Mixed Models

- Mixed models form a very useful class of regression models with general form

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon}$$

where  $\boldsymbol{\beta}$  are usual regression coefficients while  $\mathbf{b}$  are **random effects** with distributional assumption

$$\begin{bmatrix} \boldsymbol{\varepsilon} \\ \mathbf{b} \end{bmatrix} \sim \text{N} \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma^2 \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix} \right).$$

- In the following, we will concentrate on mixed models with **only one variance component** where

$$\mathbf{b} \sim \text{N}(\mathbf{0}, \tau^2 \mathbf{I}) \quad \text{or} \quad \mathbf{b} \sim \text{N}(\mathbf{0}, \tau^2 \boldsymbol{\Sigma})$$

with  $\boldsymbol{\Sigma}$  known.

- Special case I: Random intercept model for **longitudinal data**

$$y_{ij} = \mathbf{x}'_{ij}\boldsymbol{\beta} + b_i + \varepsilon_{ij}, \quad j = 1, \dots, J_i, \quad i = 1, \dots, I,$$

where  $i$  indexes individuals while  $j$  indexes **repeated observations** on the same individual.

- The random intercept  $b_i$  accounts for shifts in the individual level of response trajectories and therefore also for **intra-subject correlations**.

- Special case II: **Penalised spline smoothing** for nonparametric function estimation

$$y_i = m(x_i) + \varepsilon_i, \quad i = 1, \dots, n,$$

where  $m(x)$  is a **smooth, unspecified function**.

- Approximating  $m(x)$  in terms of a **spline basis of degree  $d$**  leads (for example) to the truncated power series representation

$$m(x) = \sum_{j=0}^d \beta_j x^j + \sum_{j=1}^K b_j (x - \kappa_j)_+^d$$

where  $\kappa_1, \dots, \kappa_K$  denotes a sequence of knots.

- Assume random effects distribution  $\mathbf{b} \sim \mathbf{N}(\mathbf{0}, \tau^2 \mathbf{I})$  for the basis coefficients of truncated polynomials to **enforce smoothness**.
- Works also for other basis choices (e.g. B-splines) and other types of flexible modelling components (varying coefficients, surfaces, spatial effects, etc.).

- Additive mixed models consist of a **combination of random effects and flexible modelling components** such as penalised splines.
- Example: Childhood malnutrition in Zambia.
- Determine the nutritional status of a child in terms of a Z-score.
- We consider chronic malnutrition measured in terms of **insufficient height for age (stunting)**, i.e.

$$zscore_i = \frac{cheight_i - med}{s},$$

where  $med$  and  $s$  are the median and standard deviation of (age-stratified) height in a reference population.

- Additive mixed model for stunting:

$$\begin{aligned} zscore_i = & \mathbf{x}'_i \boldsymbol{\beta} + m_1(cage_i) + m_2(cfeed_i) + m_3(mage_i) + m_4(mbmi_i) \\ & + m_5(mheight_i) + b_{s_i} + \varepsilon_i, \end{aligned}$$

with covariates

<i>csex</i>	gender of the child (1 = male, 0 = female)
<i>cfeed</i>	duration of breastfeeding (in months)
<i>cage</i>	age of the child (in months)
<i>mage</i>	age of the mother (at birth, in years)
<i>mheight</i>	height of the mother (in cm)
<i>mbmi</i>	body mass index of the mother
<i>medu</i>	education of the mother (1 = no education, 2 = primary school, 3 = elementary school, 4 = higher)
<i>mwork</i>	employment status of the mother (1 = employed, 0 = unemployed)
<i>s</i>	residential district (54 districts in total)

- The random effect  $b_{s_i}$  captures **spatial variability** induced by unobserved spatially varying covariates.

- **Marginal perspective** on a mixed model:

$$\mathbf{y} \sim \text{N}(\mathbf{X}\boldsymbol{\beta}, \mathbf{V})$$

where

$$\mathbf{V} = \sigma^2 \mathbf{I} + \mathbf{ZDZ}'$$

- Interpretation: The random effects induce a **correlation structure** and therefore enable a proper statistical analysis of correlated data.
- **Conditional perspective** on a mixed model:

$$\mathbf{y}|\mathbf{b} \sim \text{N}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Zb}, \sigma^2 \mathbf{I}).$$

- Interpretation: Random effects are **additional regression coefficients** (for example subject-specific effects in longitudinal data) that are estimated subject to a regularisation penalty.

- Interest in the following is on the selection of random effects: Compare

$$M_1 : \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon}, \quad \mathbf{b} \sim \mathbf{N}(\mathbf{0}, \tau^2 \boldsymbol{\Sigma})$$

and

$$M_2 : \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}.$$

- Equivalent: Compare model with random effects ( $\tau^2 > 0$ ) and without random effects ( $\tau^2 = 0$ ).
- Random Intercept:  $\tau^2 > 0$  versus  $\tau^2 = 0$  corresponds to the **inclusion and exclusion of the random intercept** and therefore to the presence or absence of intra-individual correlations.
- Penalised splines:  $\tau^2 > 0$  versus  $\tau^2 = 0$  differentiates between a spline model and a simple polynomial model. In particular, we can compare **linear versus nonlinear models**.

## Akaike Information Criteria in Linear Mixed Models

- In linear mixed models, **two variants of AIC** are conceivable based on either the marginal or the conditional distribution.
- The **marginal AIC relies** on the marginal model

$$\mathbf{y} \sim \text{N}(\mathbf{X}\boldsymbol{\beta}, \mathbf{V})$$

and is defined as

$$\text{mAIC} = -2l(\mathbf{y}|\hat{\boldsymbol{\beta}}, \hat{\tau}^2, \hat{\sigma}^2) + 2(p + 2),$$

where the **marginal likelihood** is given by

$$l(\mathbf{y}|\hat{\boldsymbol{\beta}}, \hat{\tau}^2, \hat{\sigma}^2) = -\frac{1}{2} \log(|\hat{\mathbf{V}}|) - \frac{1}{2}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})' \hat{\mathbf{V}}^{-1} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

and  $p = \dim(\boldsymbol{\beta})$ .



- The **conditional AIC** relies on the conditional model

$$\mathbf{y}|\mathbf{b} \sim \text{N}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b}, \sigma^2\mathbf{I})$$

and is defined as

$$\text{cAIC} = -2l(\mathbf{y}|\hat{\boldsymbol{\beta}}, \hat{\mathbf{b}}, \hat{\tau}^2, \sigma^2) + 2(\rho + 1),$$

where

$$l(\mathbf{y}|\hat{\boldsymbol{\beta}}, \hat{\mathbf{b}}, \hat{\tau}^2, \sigma^2) = -\frac{n}{2} \log(\hat{\sigma}^2) - \frac{1}{2\hat{\sigma}^2} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{Z}\hat{\mathbf{b}})' (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{Z}\hat{\mathbf{b}})$$

is the **conditional likelihood** and

$$\rho = \text{tr} \left( \left( \begin{pmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{Z} \\ \mathbf{Z}'\mathbf{X} & \mathbf{Z}'\mathbf{Z} + \sigma^2/\tau^2\boldsymbol{\Sigma} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{Z} \\ \mathbf{Z}'\mathbf{X} & \mathbf{Z}'\mathbf{Z} \end{pmatrix} \right) \right)$$

are the **effective degrees of freedom** (trace of the hat matrix).

- The conditional AIC seems to be recommended when the model shall be used for predictions with the **same set of random effects** (for example in penalised spline smoothing).
- The marginal AIC is more plausible when observations with **new random effects** shall be predicted (e.g. new individuals in longitudinal studies).

## Marginal AIC

- Model  $M_1$  ( $\tau^2 > 0$ ) is preferred over  $M_2$  ( $\tau^2 = 0$ ) when

$$\begin{aligned} \text{mAIC}_1 < \text{mAIC}_2 &\Leftrightarrow -2l(\mathbf{y}|\hat{\boldsymbol{\beta}}_1, \hat{\tau}^2, \hat{\sigma}_1^2) + 2(p+2) < -2l(\mathbf{y}|\hat{\boldsymbol{\beta}}_2, \mathbf{0}, \hat{\sigma}_2^2) + 2(p+1) \\ &\Leftrightarrow 2l(\mathbf{y}|\hat{\boldsymbol{\beta}}_1, \hat{\tau}^2, \hat{\sigma}_1^2) - 2l(\mathbf{y}|\hat{\boldsymbol{\beta}}_2, \mathbf{0}, \hat{\sigma}_2^2) > 2. \end{aligned}$$

- The left hand side is simply the test statistic for a **likelihood ratio test on  $\tau^2 = 0$  versus  $\tau^2 > 0$** .
- Under standard asymptotics, we would have

$$2l(\mathbf{y}|\hat{\boldsymbol{\beta}}_1, \hat{\tau}^2, \hat{\sigma}_1^2) - 2l(\mathbf{y}|\hat{\boldsymbol{\beta}}_2, \mathbf{0}, \hat{\sigma}_2^2) \stackrel{a, H_0}{\sim} \chi_1^2$$

and the marginal AIC would have a type 1 error of

$$P(\chi_1^2 > 2) \approx 0.1572992$$

- Common interpretation: AIC selects **rather too many than too few effects**.

- In contrast to the regularity conditions for likelihood ratio tests,  $\tau^2$  is on the **boundary of the parameter space** for model  $M_2$ .
- The likelihood ratio test statistic is no longer  $\chi^2$ -distributed but (approximately) follows a mixture of a **point mass in zero and a scaled  $\chi_1^2$  variable**.
- The point mass in zero corresponds to the probability

$$P(\hat{\tau}^2 = 0)$$

that is typically **larger than 50%**.

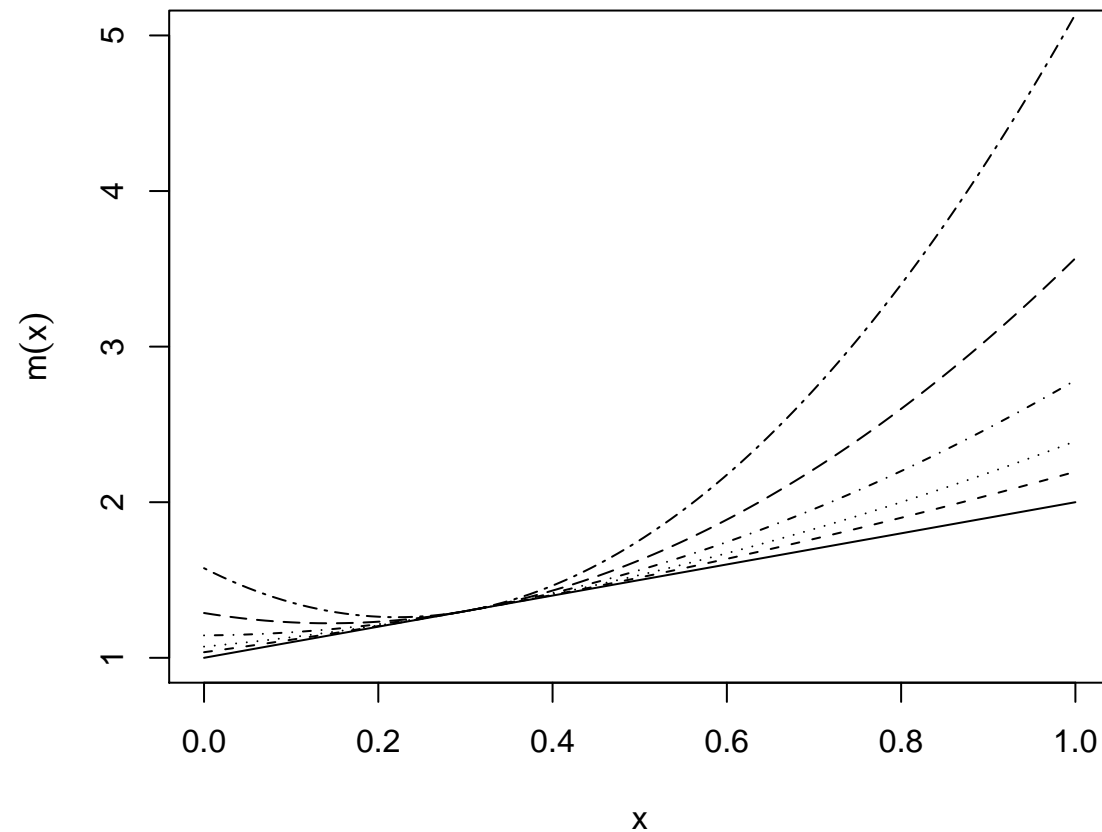
- Similar difficulties appear in more complex models with several variance components when deciding on zero variances.

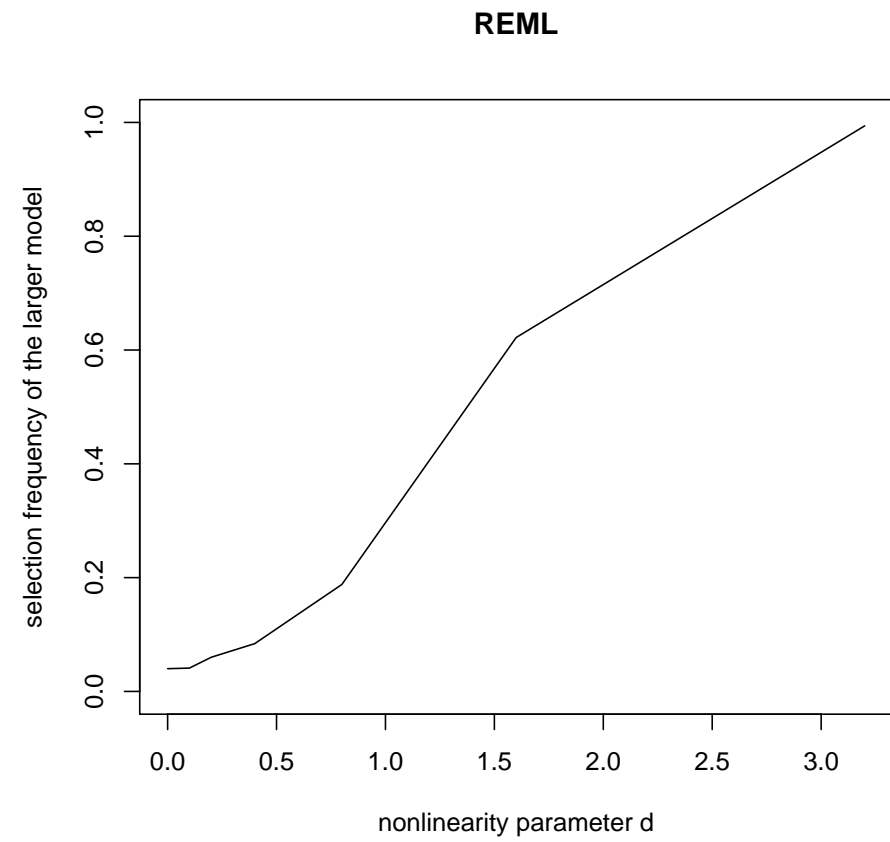
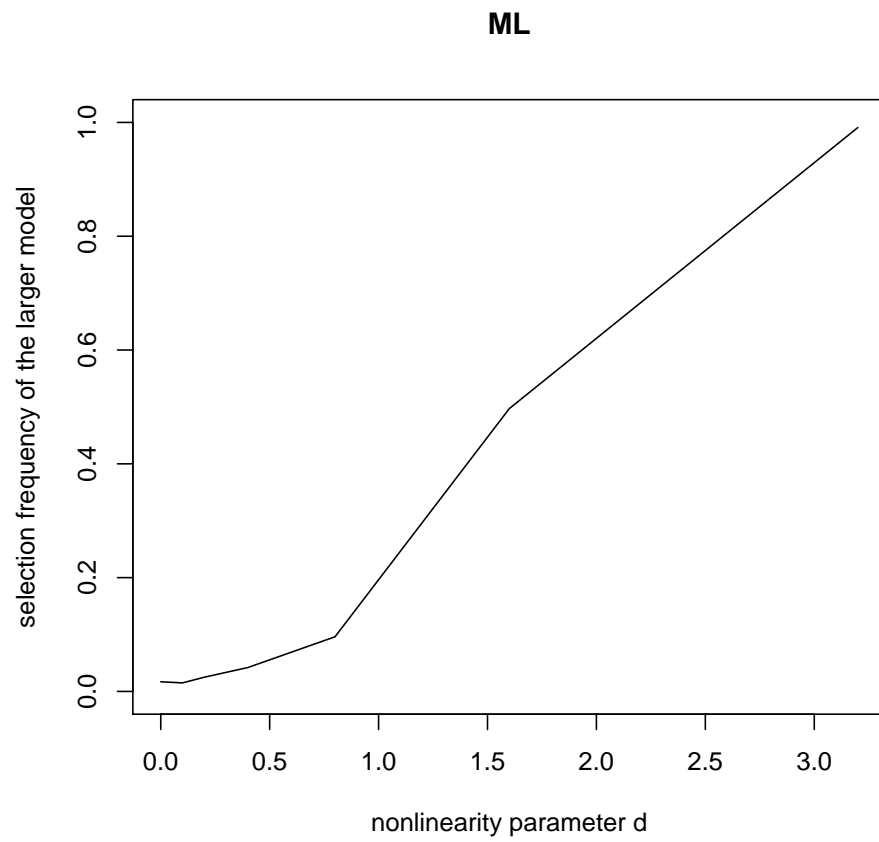
- The classical assumptions underlying the derivation of AIC are also not fulfilled.
- The high probability of estimating a zero variance yields a **bias towards simpler models**:
  - The marginal AIC is positively biased for twice the expected relative Kullback-Leibler-Distance.
  - The bias is dependent on the true unknown parameters in the random effects covariance matrix and this dependence does not vanish asymptotically.
  - Compared to an unbiased criterion, the marginal AIC favors smaller models excluding random effects.
- This contradicts the usual intuition that the AIC picks rather too many than too few effects.

- Simulated example:  $y_i = m(x) + \varepsilon$  where

$$m(x) = 1 + x + 2d(0.3 - x)^2.$$

- The parameter  $d$  determines the **amount of nonlinearity**.

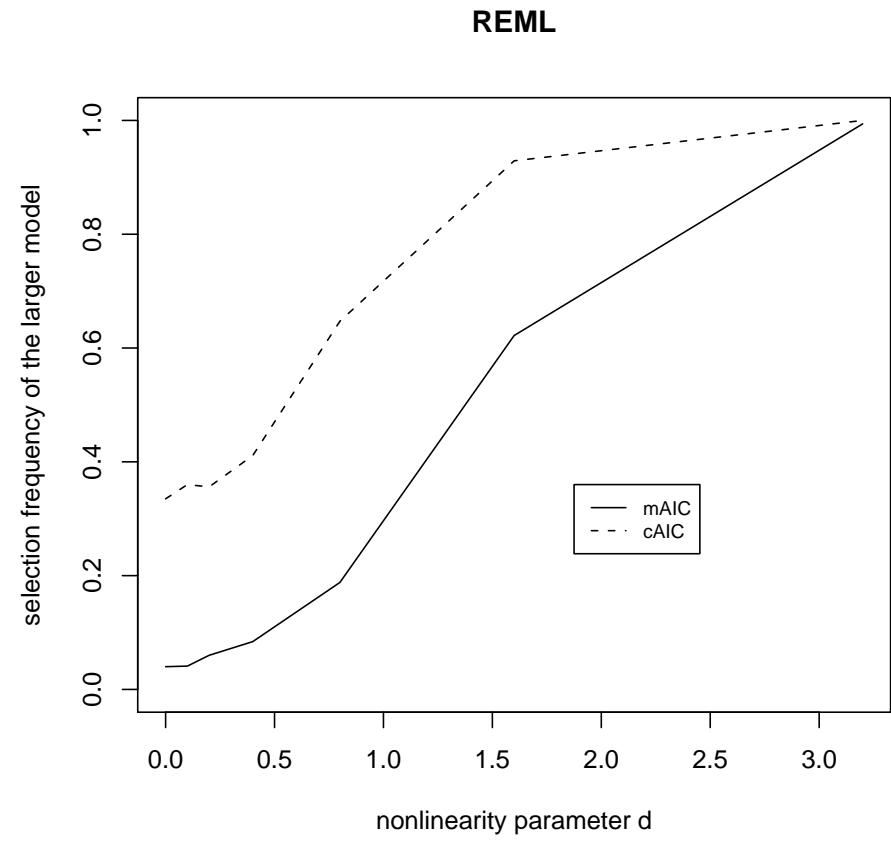
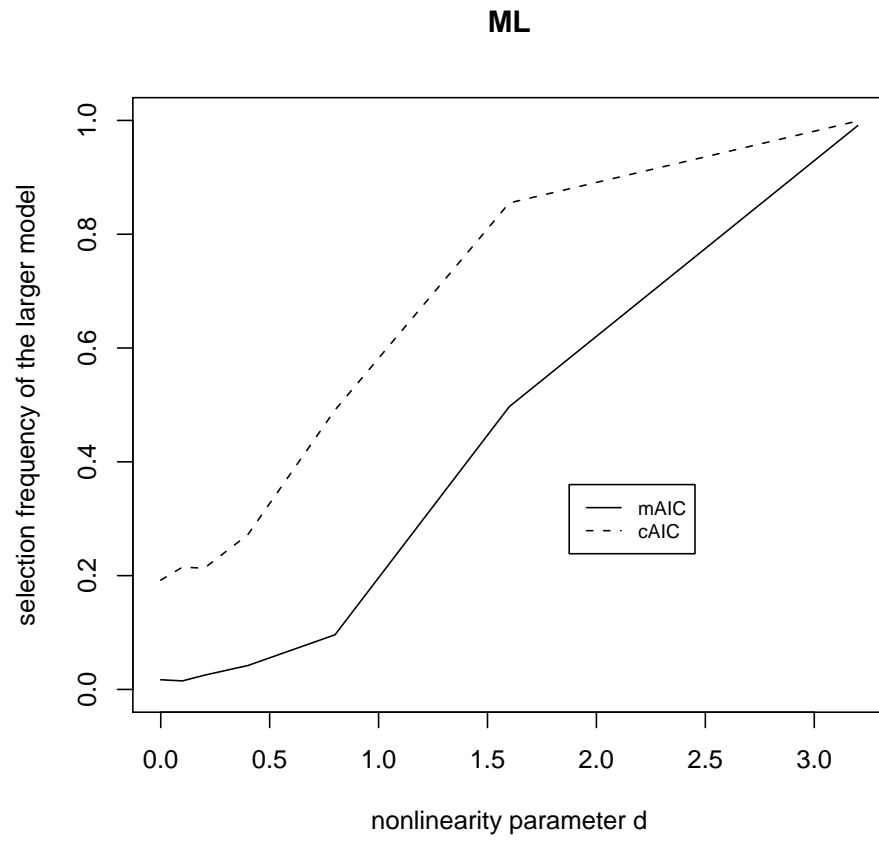




## Conditional AIC

- Vaida & Blanchard (2005) have shown that the conditional AIC from above is asymptotically unbiased for the expected relative Kullback Leibler distance for **given random effects covariance matrix**.
- Intuition: Result should carry over when using a **consistent estimate**.
- Simulation results indicate that this is not the case.





- Surprising result of the simulation study: The complex model including the random effect is chosen **whenever  $\hat{\tau}^2 > 0$** .
- If  $\hat{\tau}^2 = 0$ , the conditional AICs of the simple and the complex model coincide (despite the additional parameters included in the complex model).
- The observed phenomenon can be shown to be a general property of the conditional AIC:

$$\hat{\tau}^2 > 0 \quad \Leftrightarrow \quad \text{cAIC}(\hat{\tau}^2) < \text{cAIC}(0)$$

$$\hat{\tau}^2 = 0 \quad \Leftrightarrow \quad \text{cAIC}(\hat{\tau}^2) = \text{cAIC}(0).$$

- Principal difficulty: The degrees of freedom in the cAIC are **estimated from the same data as the model parameters**.

- Liang et al. (2008) propose a **corrected conditional AIC**, where the degrees of freedom  $\rho$  are replaced by

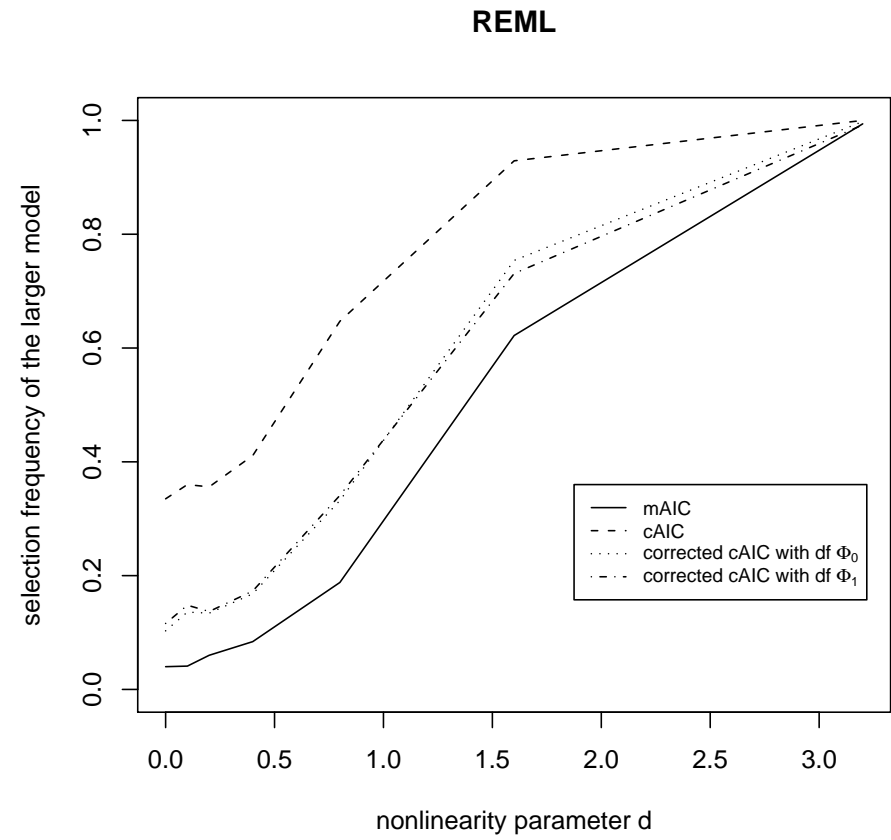
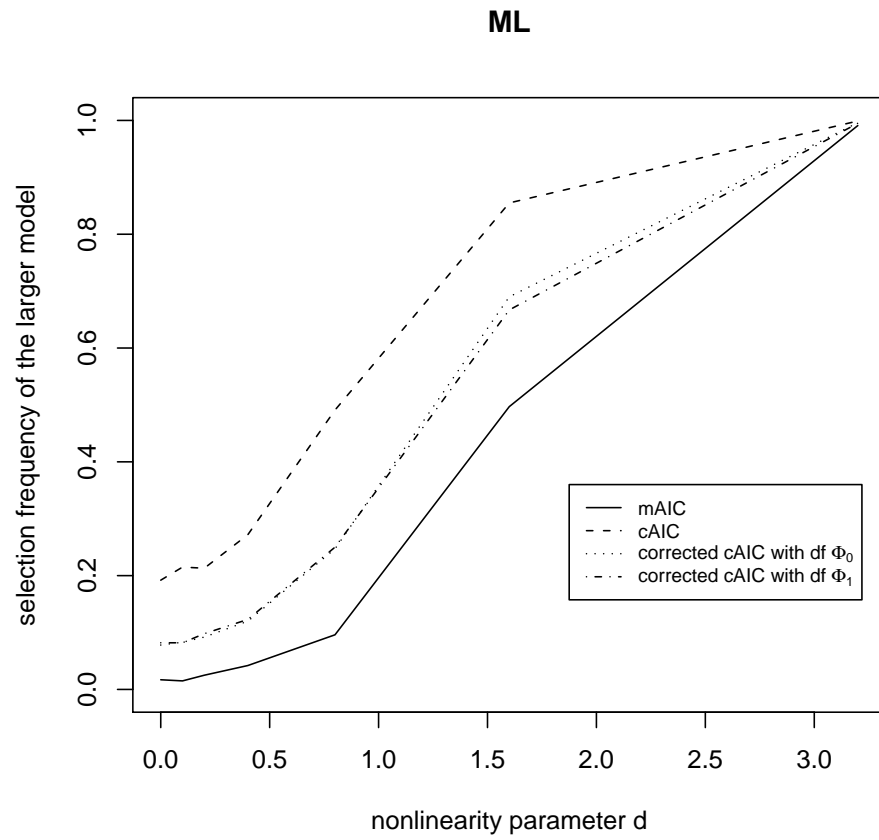
$$\Phi_0 = \sum_{i=1}^n \frac{\partial \hat{y}_i}{\partial y_i} = \text{tr} \left( \frac{\partial \hat{\mathbf{y}}}{\mathbf{y}} \right)$$

if  $\sigma^2$  is known.

- For unknown  $\sigma^2$ , they propose to replace  $\rho + 1$  by

$$\Phi_1 = \frac{\tilde{\sigma}^2}{\hat{\sigma}^2} \text{tr} \left( \frac{\partial \hat{\mathbf{y}}}{\mathbf{y}} \right) + \tilde{\sigma}^2 (\hat{\mathbf{y}} - \mathbf{y})' \frac{\partial \hat{\sigma}^{-2}}{\partial \mathbf{y}} + \frac{1}{2} \tilde{\sigma}^4 \text{tr} \left( \frac{\partial^2 \hat{\sigma}^{-2}}{\partial \mathbf{y} \partial \mathbf{y}'} \right),$$

where  $\tilde{\sigma}^2$  is an estimate for the true error variance.



- The corrected conditional AIC shows **satisfactory theoretical properties**.
- However, it is **computationally cumbersome**:
  - Liang et al. suggested to approximate the derivatives numerically (by adding small perturbations to the data).
  - Numerical approximations require  $n$  and  $2n$  model fits. In our application, computing the corrected conditional AICs would take about 110 days.
  - In addition, the numerical derivatives were found to be instable in some situations (for example the random intercept model with small cluster sizes).
- We have developed a **closed form representation of  $\Phi_0$**  that is available almost instantaneously.

## Application: Childhood Malnutrition in Zambia

- Model equation:

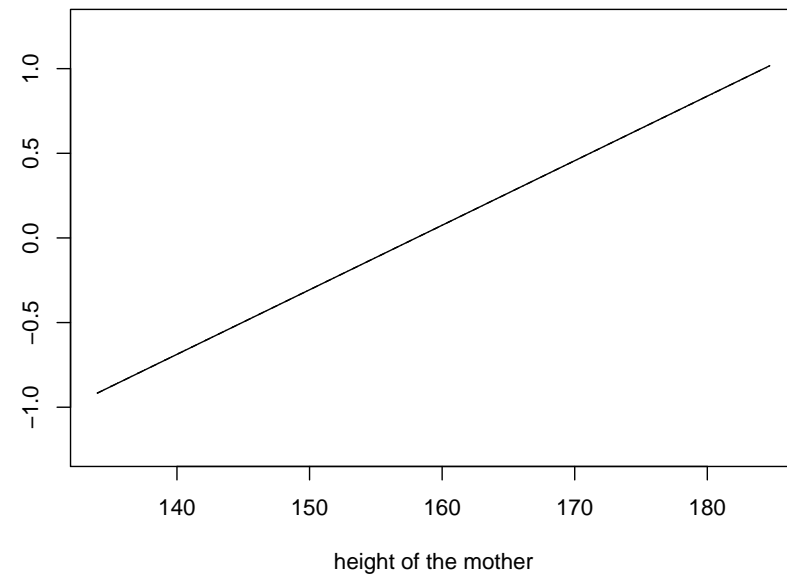
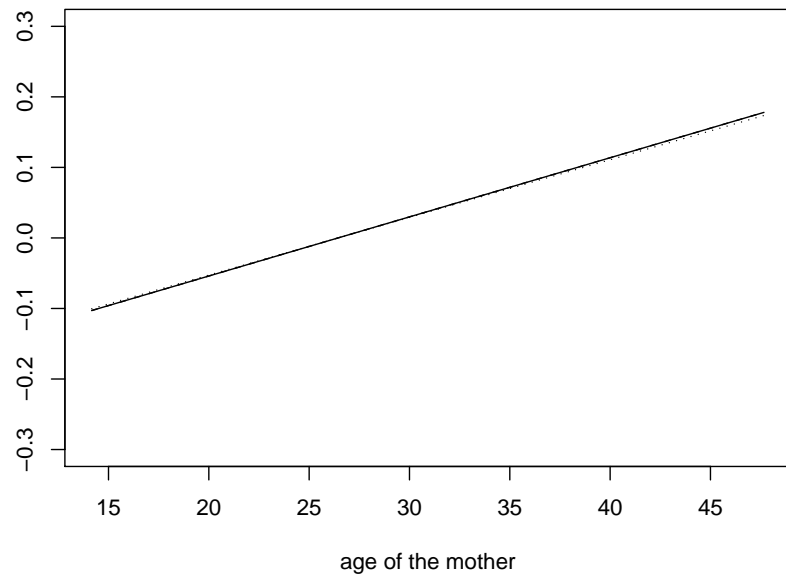
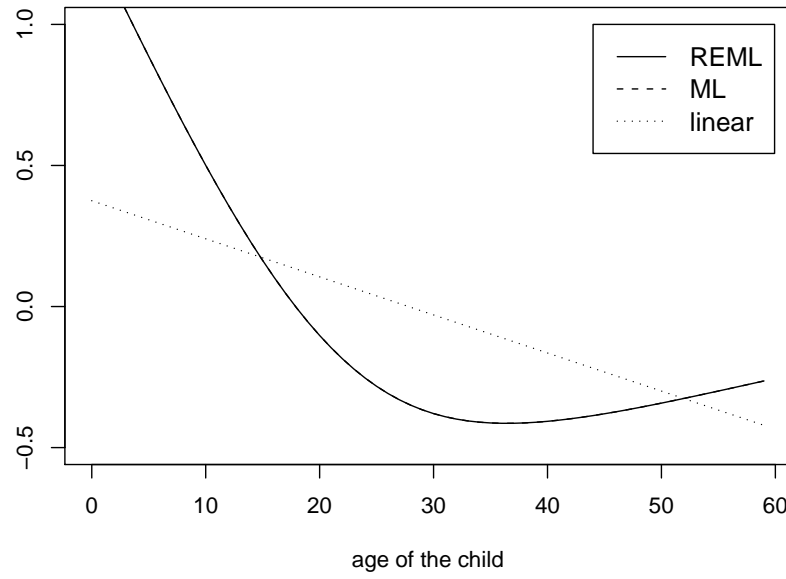
$$\begin{aligned} zscore_i = & \mathbf{x}'_i \boldsymbol{\beta} + m_1(cage_i) + m_2(cfeed_i) + m_3(mage_i) + m_4(mbmi_i) \\ & + m_5(mheight_i) + b_{s_i} + \varepsilon_i. \end{aligned}$$

- Parametric effects in  $\mathbf{x}'\boldsymbol{\beta}$  are not subject to model selection.  
 $\Rightarrow 2^6 = 64$  models to consider in the model comparison.

- The eight best fitting models:

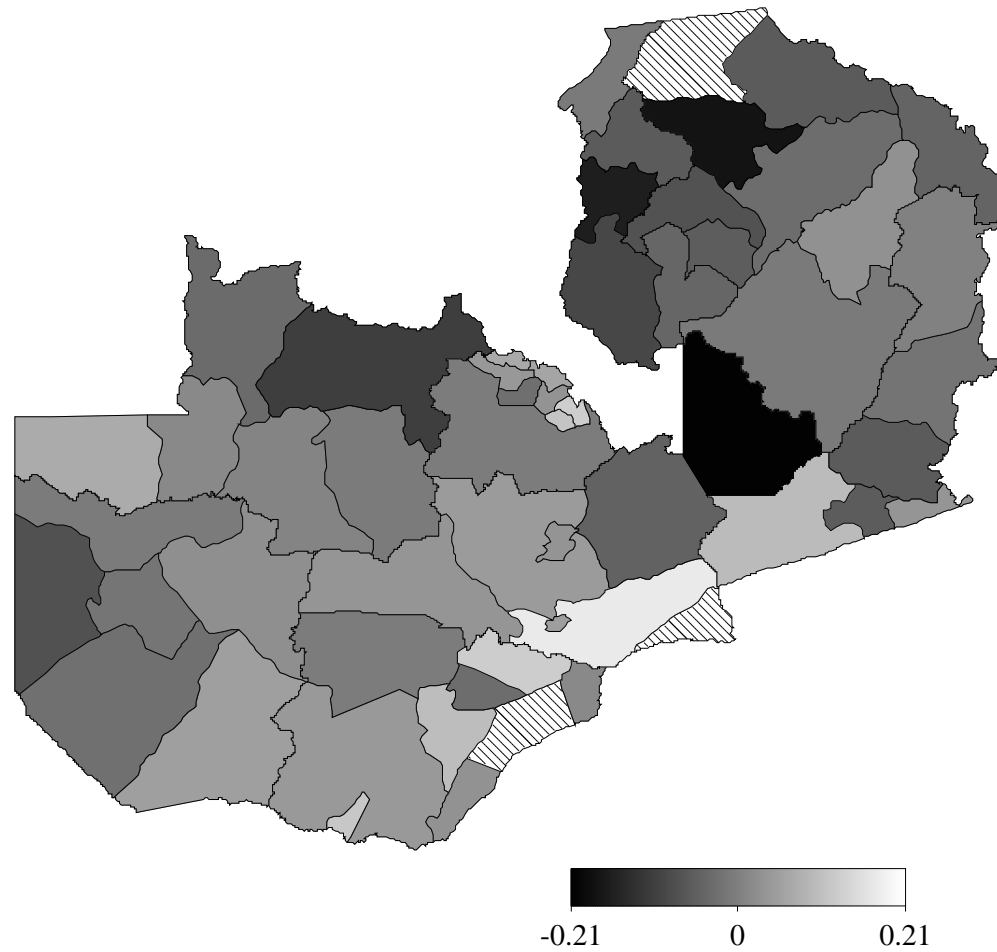
cfeed	cage	mage	mheight	mbmi	s	<i>ML</i>		<i>REML</i>	
						<i>cAIC</i>	<i>mAIC</i>	<i>cAIC</i>	<i>mAIC</i>
+	+	-	-	-	+	<b>4125.78</b>	<b>4151.10</b>	<b>4125.78</b>	<b>4173.72</b>
+	+	+	-	-	+	<b>4125.78</b>	4153.10	<b>4125.78</b>	4175.72
+	+	-	+	-	+	<b>4125.78</b>	4153.10	<b>4125.78</b>	4175.72
+	+	-	-	+	+	<b>4125.78</b>	4153.10	<b>4125.78</b>	4175.72
+	+	+	+	-	+	<b>4125.78</b>	4155.10	<b>4125.78</b>	4177.72
+	+	+	-	+	+	<b>4125.78</b>	4155.10	<b>4125.78</b>	4177.72
+	+	-	+	+	+	<b>4125.78</b>	4155.10	<b>4125.78</b>	4177.72
+	+	+	+	+	+	<b>4125.78</b>	4157.10	<b>4125.78</b>	4179.72

- Linear effects are selected for age, height and body mass index of the mother.
- Some nonlinearity is detected for age of the child and duration of breastfeeding.





- Inclusion of the region-specific random effect is required to capture spatial variation.



## Summary

- The marginal AIC suffers from the same theoretical difficulties as likelihood ratio tests on the boundary of the parameter space.
- The marginal AIC is biased towards simpler models excluding random effects.
- The conventional conditional AIC tends to select too many variables.
- Whenever a random effects variance is estimated to be positive, the corresponding effect will be included.
- The corrected conditional AIC rectifies this difficulty and is now available in closed form.

- References:
  - Greven, S. & Kneib, T. (2009): On the Behavior of Marginal and Conditional Akaike Information Criteria in Linear Mixed Models. Technical Report.
  - Liang, H., Wu, H. & Zou, G. (2008): A note on conditional AIC for linear mixed-effects models. *Biometrika* 95, 773–778.
  - Vaida, F. & Blanchard, S. (2005): Conditional Akaike information for mixed-effects models. *Biometrika* 92, 351–370.
- A place called home:

`http://www.staff.uni-oldenburg.de/thomas.kneib`