

7. Model building and model choice

Sonja Greven

Summer Term 2015

General recommendations

- As $E(\mathbf{b}_i) = \mathbf{0}$, all covariates in \mathbf{Z}_i should be linear transformations of covariates in \mathbf{X}_i .
- If \mathbf{Z}_i contains x^p , it should also contain x^k , $k = 0, \dots, (p - 1)$.
- The more complex the structure for the fixed and random effects is, the simpler the covariance structure in Σ_i should be.
- A saturated model with one β coefficient per covariate combination can be useful if there are few discrete covariates and few equal timepoints t_j . This avoids parametric assumptions. On the downside, no trends are modeled and several parameters per covariate may have to be tested (loss of power).

Overview Chapter 7 - Model building and model choice

7.1 Model diagnostics

7.2 Model selection

Residual diagnostics 1

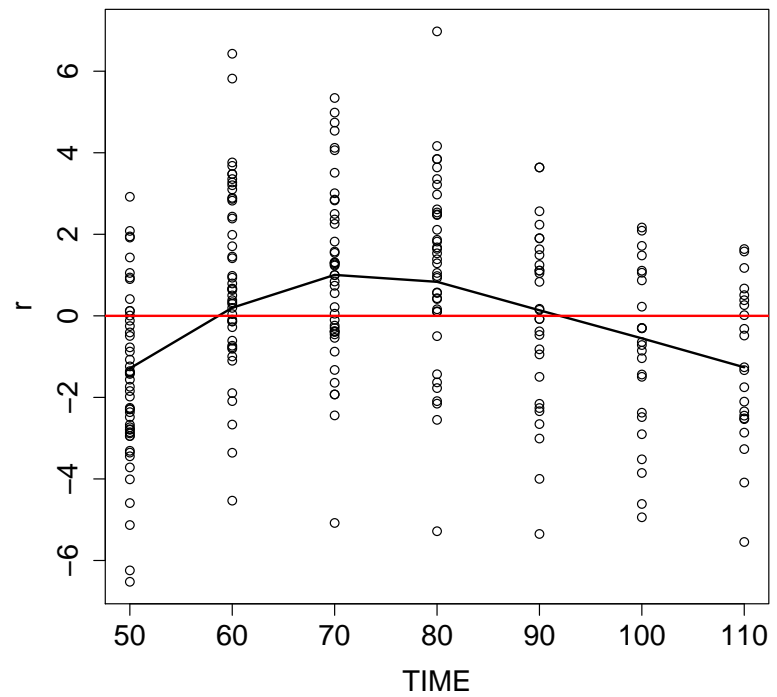
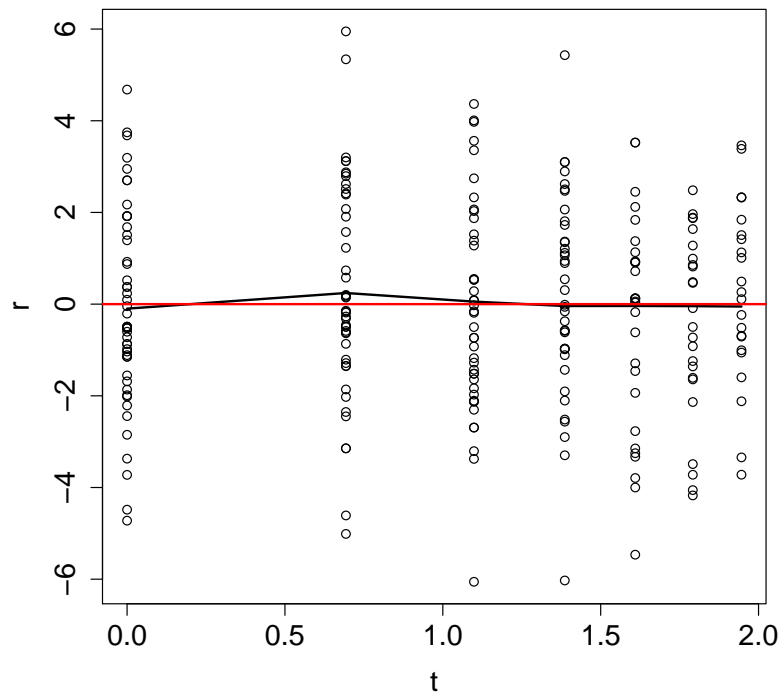
Plotting the residuals $r_{ij} = y_{ij} - \mathbf{x}_{ij}^T \hat{\boldsymbol{\beta}}$ against covariates can help in diagnosing a misspecified mean structure, e.g. an omitted variable or a missing quadratic term. There should be no systematic trend!

Example rat data, random intercept model with linear trend in transformed time variable $t = \log(1 + (TIME - 50)/10)$:

```
> lme1 <- lme(RESPONSE ~ group * t - group,
             random = ~ 1 | SUBJECT, data = rats)
> r <- resid(lme1, level = 0) # 0 - without random effects
> plot(rats$t, r, xlab = "t")
> lines(lowess(rats$t, r))
```

Analogously for the original untransformed time variable *TIME*.

Residual diagnostics 1



Residual diagnostics 2

When plotting the residuals against the estimated mean, there should be no systematic trend.

CD4 example, random intercept, linear time trend with breakpoint in 0:

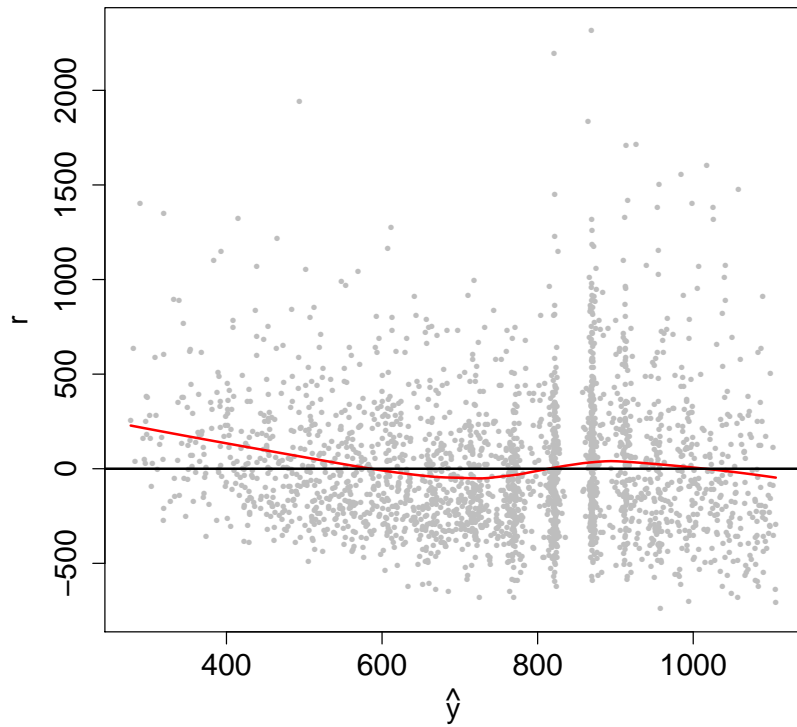
```
> cd4$Timesc <- cd4$Time * (cd4$Time > 0) # for breakpoint
> lme1 <- lme(CD4 ~ Time + Timesc, data = cd4, random = ~ 1|ID)
> yhat <- predict(lme1, level = 0) # 0 - predictions with-
> r <- resid(lme1, level = 0) # out random effects
> plot(yhat, r)
> lines(lowess(yhat, r, iter = 0))
> abline(h = 0)
```

For comparison, random intercept model with smooth time trend:

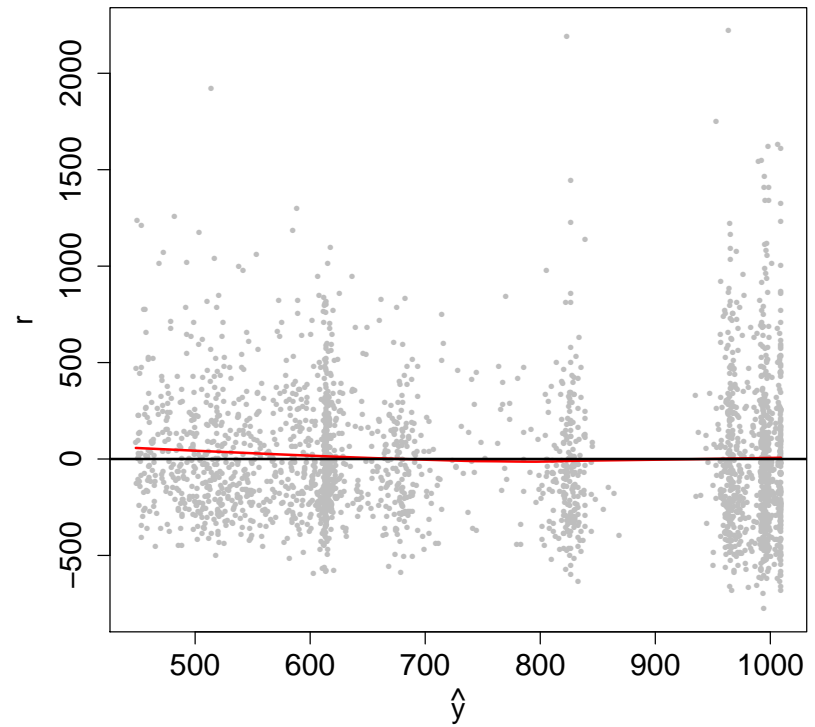
```
> mygamm <- gamm(CD4 ~ s(Time), random = list(ID = ~ 1),
  data = cd4, method = "REML")
> r <- resid(mygamm$lme, level = 1) # 1 - include random
  # effects for smooth, not for subjects
> yhat <- predict(mygamm$lme, level = 1)
> plot(yhat, r)
> lines(lowess(yhat, r, iter = 0))
> abline(h = 0)
```

Residual diagnostics 2

Linear trend with breakpoint



Smooth trend



Residual diagnostics 3

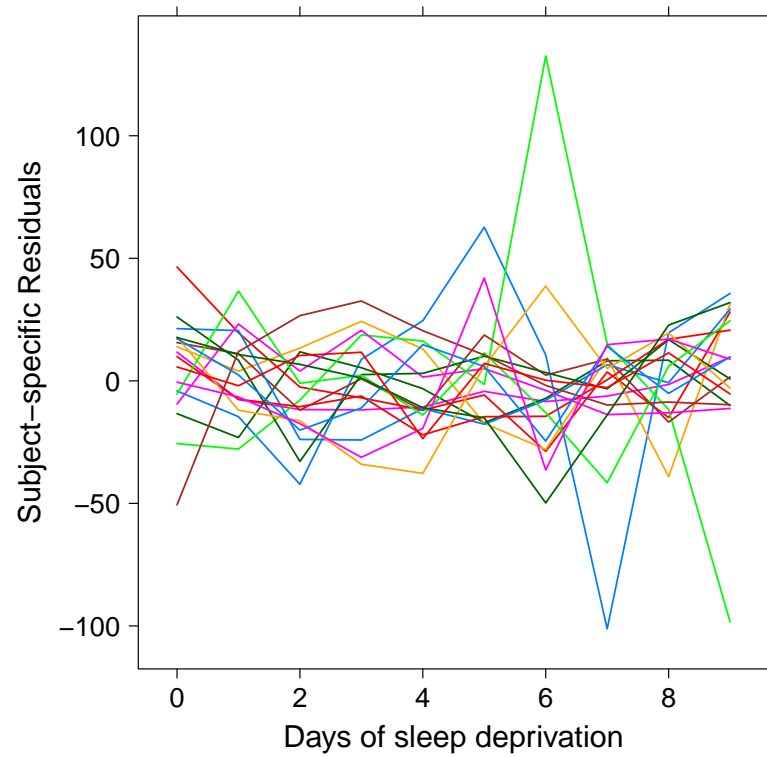
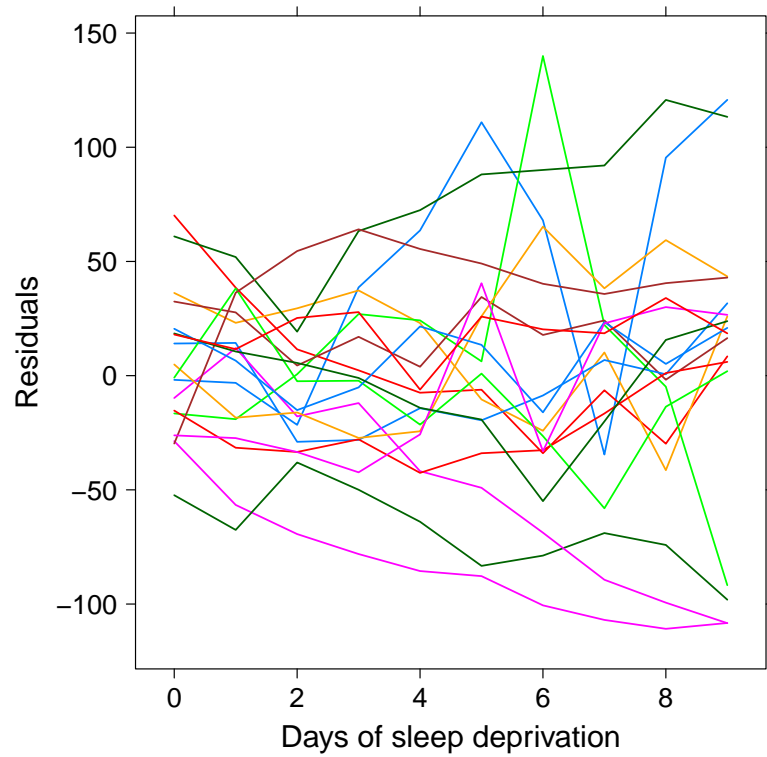
Plotting the residuals $r_{ij} = y_{ij} - \mathbf{x}_{ij}^T \hat{\boldsymbol{\beta}}$ against covariates, e.g. time, can also indicate a missing random slope.

Example sleepstudy data, models without and with random slope:

```
> lme1 <- lme(Reaction ~ Days, random = ~ 1 | Subject)
> r <- resid(lme1, level = 0) # 0: residuals w/o random effects
> xyplot(r ~ Days, groups = Subject, type = "l")

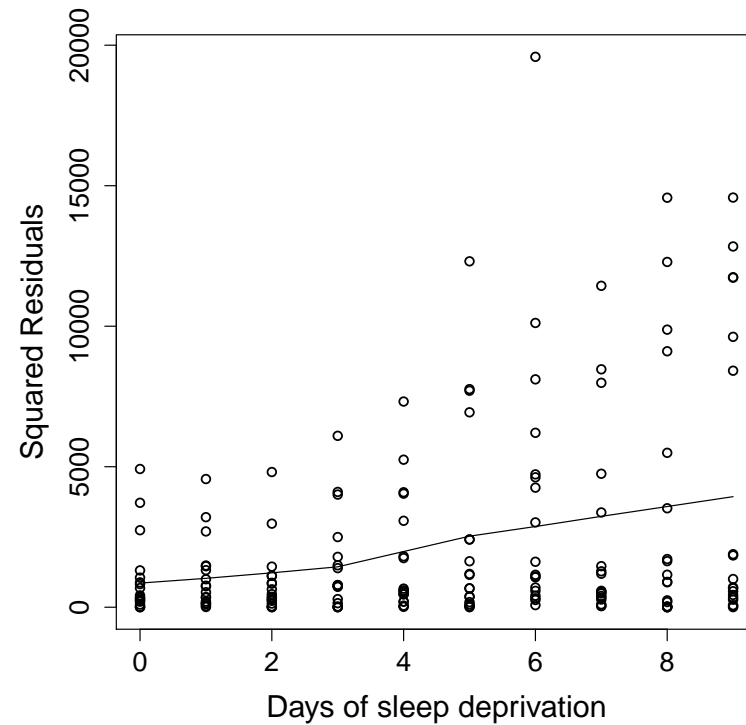
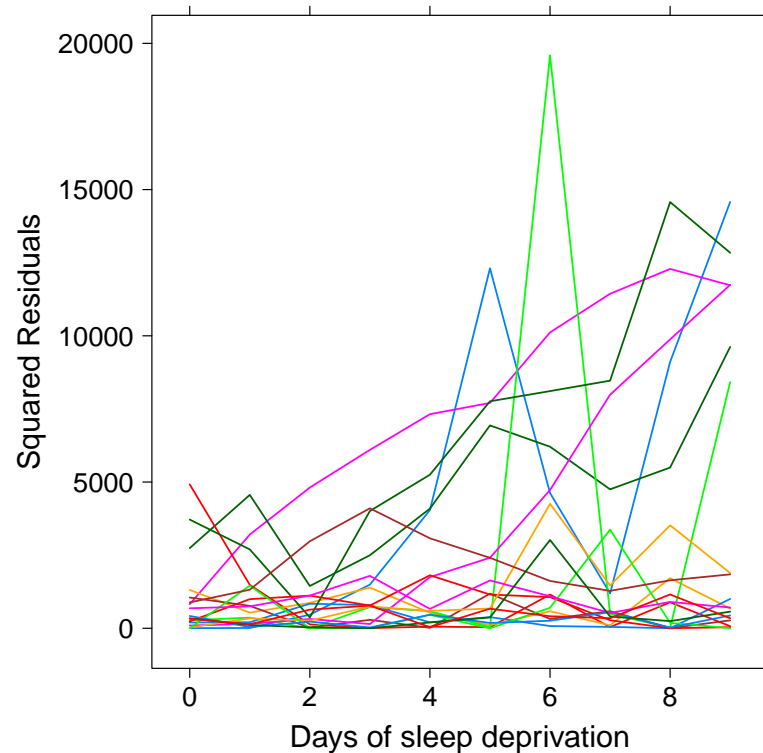
> lme2 <- lme(Reaction ~ Days, random = ~ Days | Subject)
> r <- resid(lme2, level = 1) # 1: residuals with random effects
      # to see difference when including random slope
> xyplot(r ~ Days, groups = Subject, type = "l")
```

Residual diagnostics 3



Residual diagnostics 3

We can alternatively plot the squared residuals against a potential random slope variable. For the sleep data and the random intercept model:



Transformed residuals

Remember that

$$\text{Cov}(\mathbf{Y}_i - \mathbf{X}_i\boldsymbol{\beta}) = \mathbf{V}_i.$$

Thus, the residual vector $\mathbf{r}_i = \mathbf{y}_i - \mathbf{X}_i\hat{\boldsymbol{\beta}}$ will have zero mean, but will be correlated and heteroscedastic. We need to keep this in mind for diagnostics.

One could consider the subject-specific residuals $\mathbf{y}_i - \mathbf{X}_i\hat{\boldsymbol{\beta}} - \mathbf{Z}_i\hat{\mathbf{b}}_i$. However, $\hat{\mathbf{b}}_i$ very much depends on the normality assumption for \mathbf{b}_i , and is also influenced by the assumed structure for \mathbf{V}_i .

Diagnostics are thus often based on transformed residuals $\mathbf{r}_i^* = \mathbf{L}_i^{-1}\mathbf{r}_i$, where $\hat{\mathbf{V}}_i = \mathbf{L}_i\mathbf{L}_i^T$ is the Cholesky decomposition with lower triangular matrix \mathbf{L}_i . \mathbf{r}_i^* are approximately uncorrelated with unit variance.

Transformed residuals

The transformed residuals \mathbf{r}_i^* can be shown to have the following interpretation:

- The first element is the standardized residual for y_{i1} .
- The j th element is an estimate of

$$\frac{Y_{ij} - \mathbf{E}(Y_{ij} | Y_{i1}, \dots, Y_{i(j-1)})}{\text{Var}(Y_{ij} | Y_{i1}, \dots, Y_{i(j-1)})},$$

i.e. the standardized deviation from the conditional mean given all previous observations.

Transformed residuals

After the transformation, the residuals can be used for the same kind of diagnostics as in the linear model, e.g.

- to identify **outlying observations**
- to identify skewness
- to plot the transformed residuals r_{ij}^* against the transformed predicted values $\hat{\mu}_{ij}^*$ with

$$\hat{\mu}_i^* = \mathbf{L}_i^{-1} \hat{\mu}_i = \mathbf{L}_i^{-1} \mathbf{X}_i \hat{\boldsymbol{\beta}},$$

or against a selected transformed covariate (such as e.g. time).

Outlier diagnostics

Define the **Mahalanobis distance**

$$d_i = \mathbf{r}_i^{*T} \mathbf{r}_i^*.$$

as a summary measure of multivariate distance between observed and fitted values for individual i . If the model is correctly specified, we have the approximate distribution

$$d_i \sim \chi_{n_i}^2, \quad \text{for } i = 1, \dots, N.$$

This can be used to identify **outlying individuals**: p-values can be computed for each subject and used to compare subjects, keeping in mind that p-values smaller α are expected to occur αN times.

Transformed residuals in R

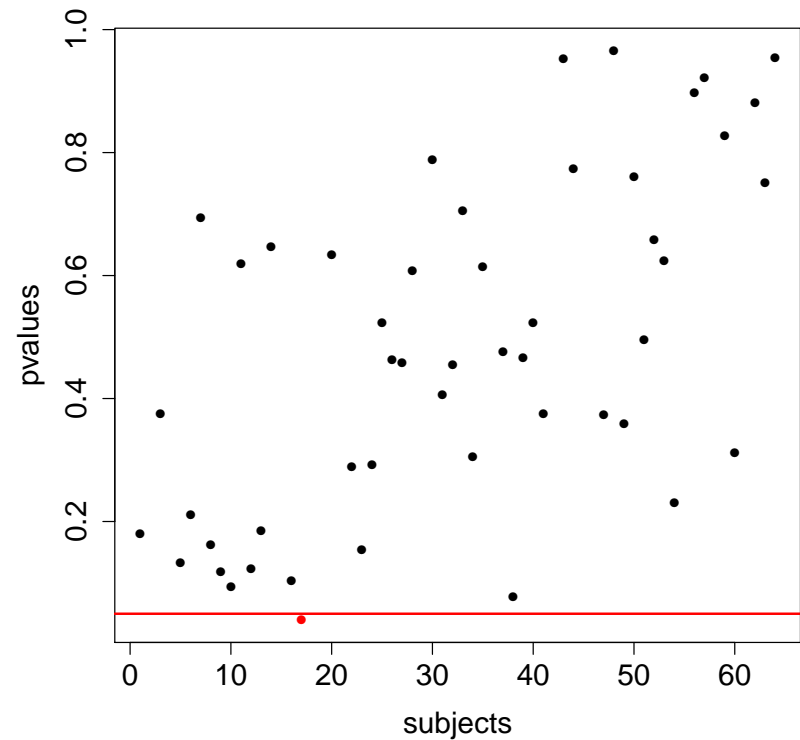
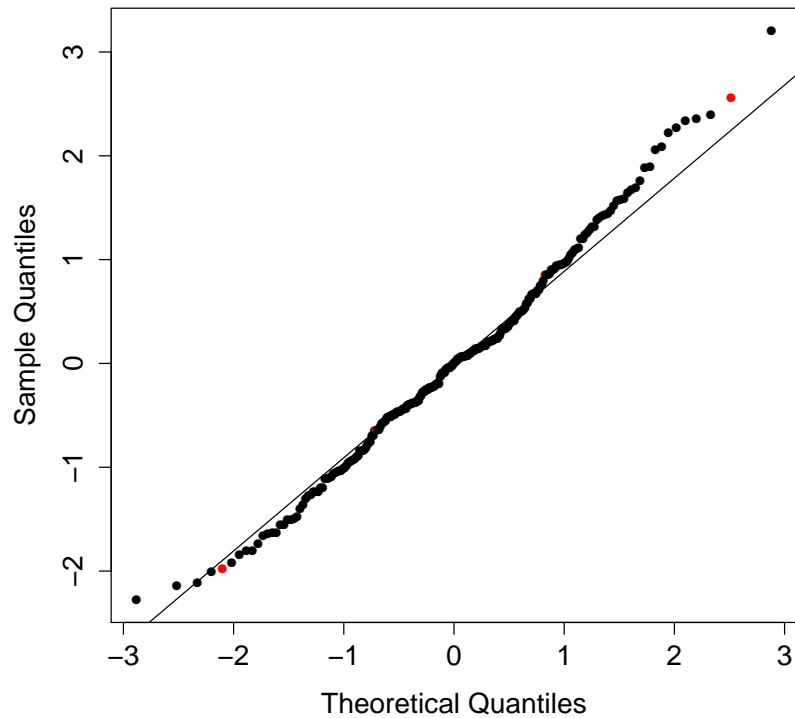
```
> library(RLRsim) # useful to extract lme model components
> r.star <- function(m){ # takes an lme object
+   design <- extract.lmeDesign(m)
+   Z <- design$Z
+   D <- design$Vr * design$sigma^2
+   R <- design$sigma^2 * diag(nrow(Z))
+   V <- Z %*% D %*% t(Z) + R
+   L <- t(chol(V))
+   r.star <- solve(L, resid(m, level = 0))
+   return(r.star) # returns the transformed residuals
+ }
```


Example rat data

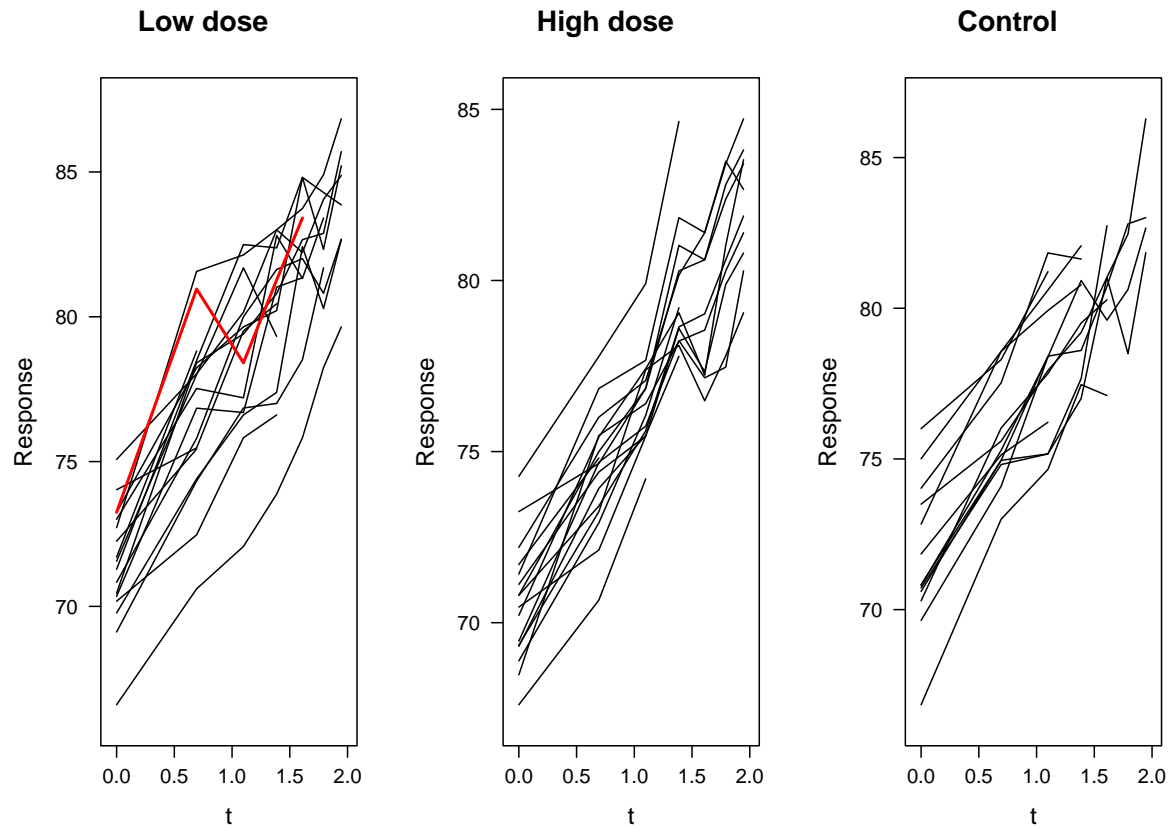
Random intercept model:

```
### Transformed residuals ###
> lme1 <- lme(RESPONSE ~ group * t - group,
              random = ~ 1 | SUBJECT, data = rats)
> r.star1 <- r.star(lme1) # transformed model residuals
### QQ-Plot ###
> qqnorm(r.star1)
> qqline(r.star1)
### Outlier Diagnostics ###
> subjects <- unique(sort(rats$SUBJECT)) # for each subject
> di <- sapply(subjects, FUN = function(subj)
               crossprod(r.star2[(rats$SUBJECT == subj)])) # compute d_i
> ni <- sapply(subjects, FUN = function(subj)
               sum(rats$SUBJECT == subj)) # and n_i
```

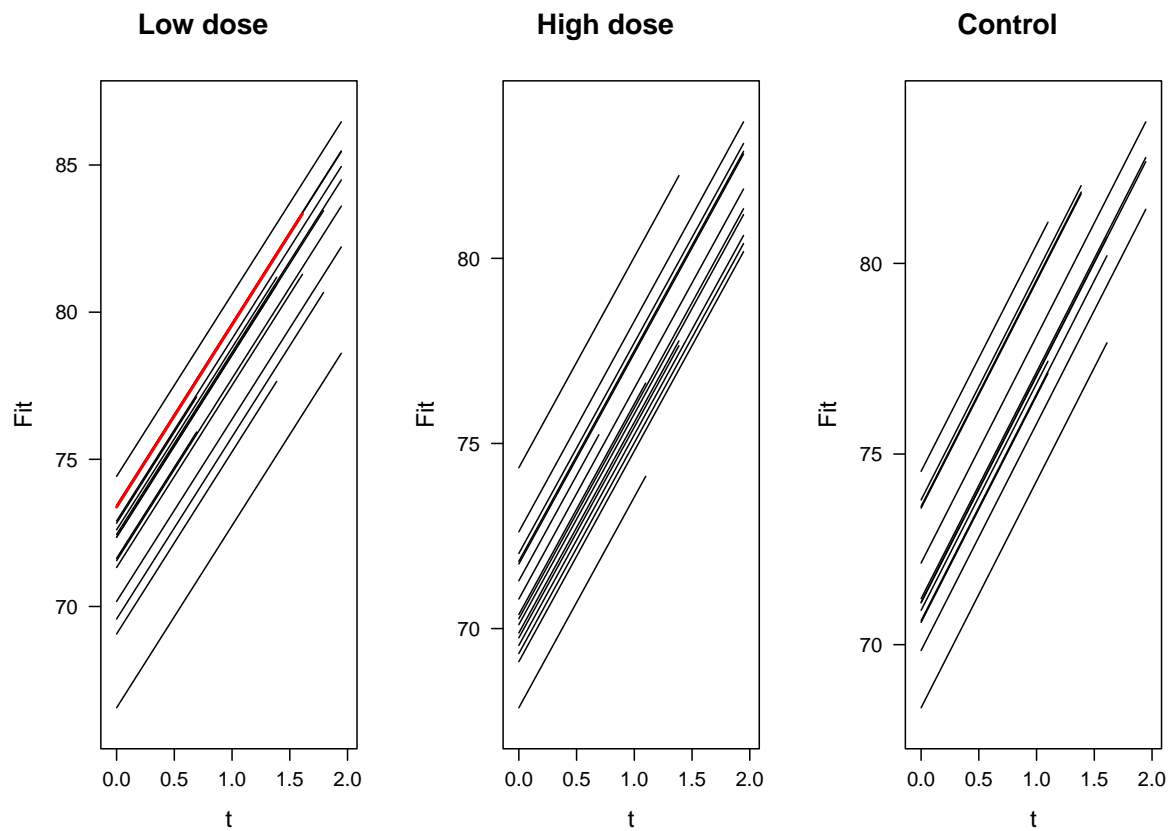
```
> pvalues <- pchisq(di, ni, lower = FALSE) # chi^2_{n_i} p-values  
> plot(subjects, pvalues); abline(h = 0.05)
```



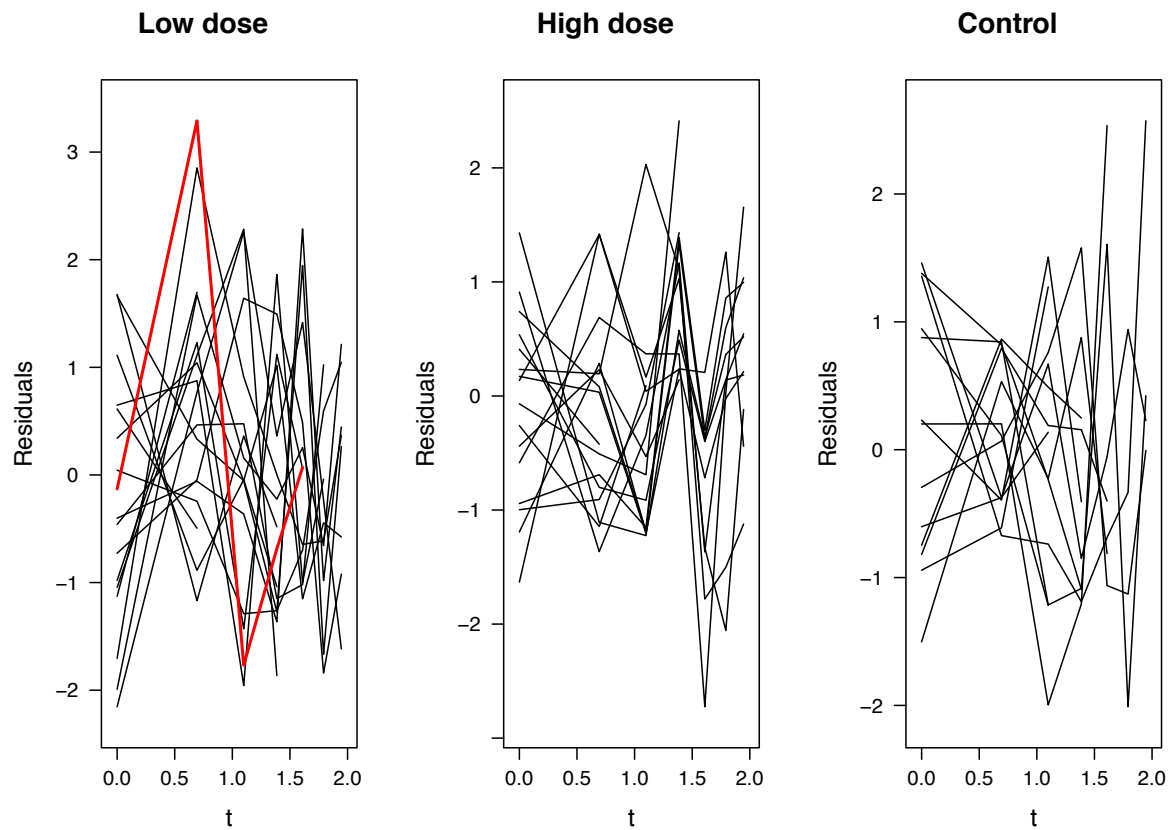
Example rat data - Data



Example rat data - Fit



Example rat data - Residuals



The choice of the covariance structure

A good model for the covariance structure is important for inference on the fixed effects, interpretation and prediction.

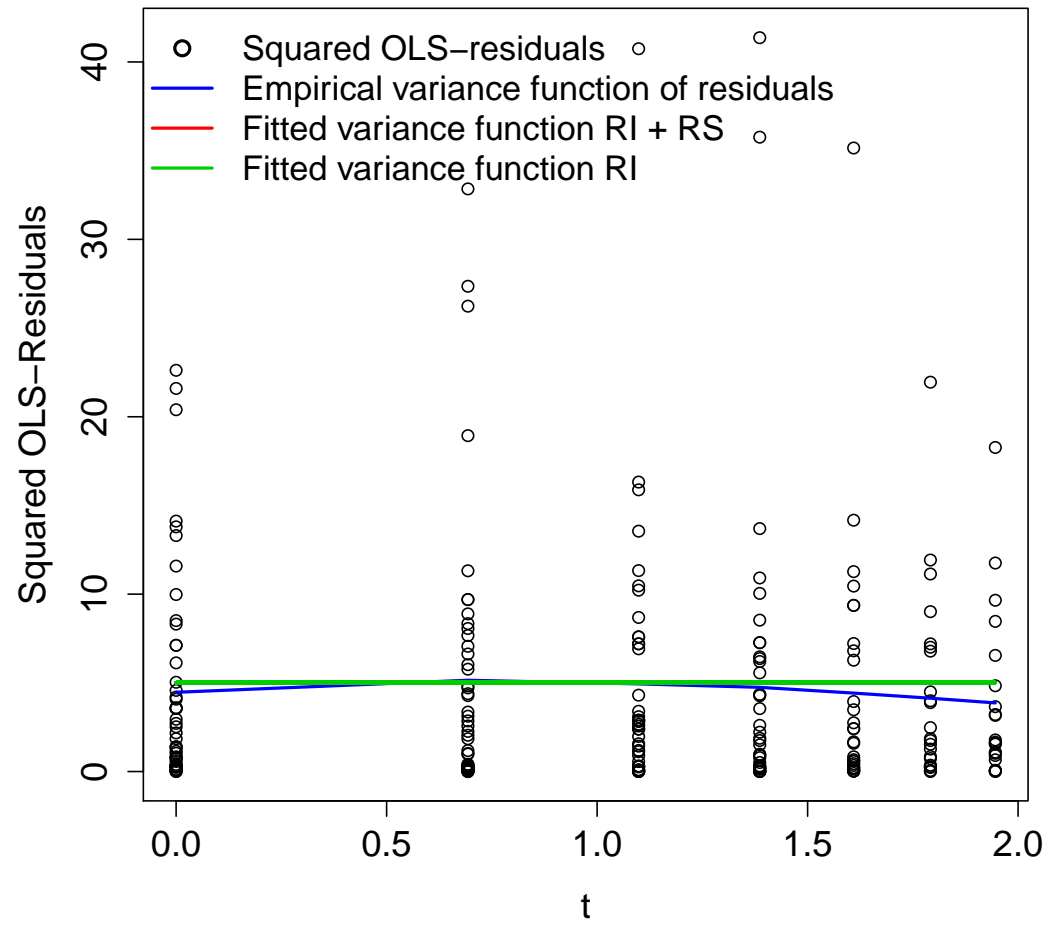
An informal check is to plot the squared OLS residuals

$$\mathbf{r}_{OLS,i} = \mathbf{y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}}_{OLS}$$

and the fitted variance function against t . The fitted variance function corresponds to the diagonal entries of $\mathbf{Z} \hat{\mathbf{D}} \mathbf{Z}^T + \hat{\mathbf{R}}$.

Example rat data with random intercept and slope: The fitted variance function is

$$(1 \ t) \hat{\mathbf{D}} \begin{pmatrix} 1 \\ t \end{pmatrix} + \hat{\sigma}^2 = \hat{d}_{11} + 2\hat{d}_{12}t + \hat{d}_{22}t^2 + \hat{\sigma}^2.$$



The semi-variogram revisited

A more comprehensive check for the covariance structure is the following.

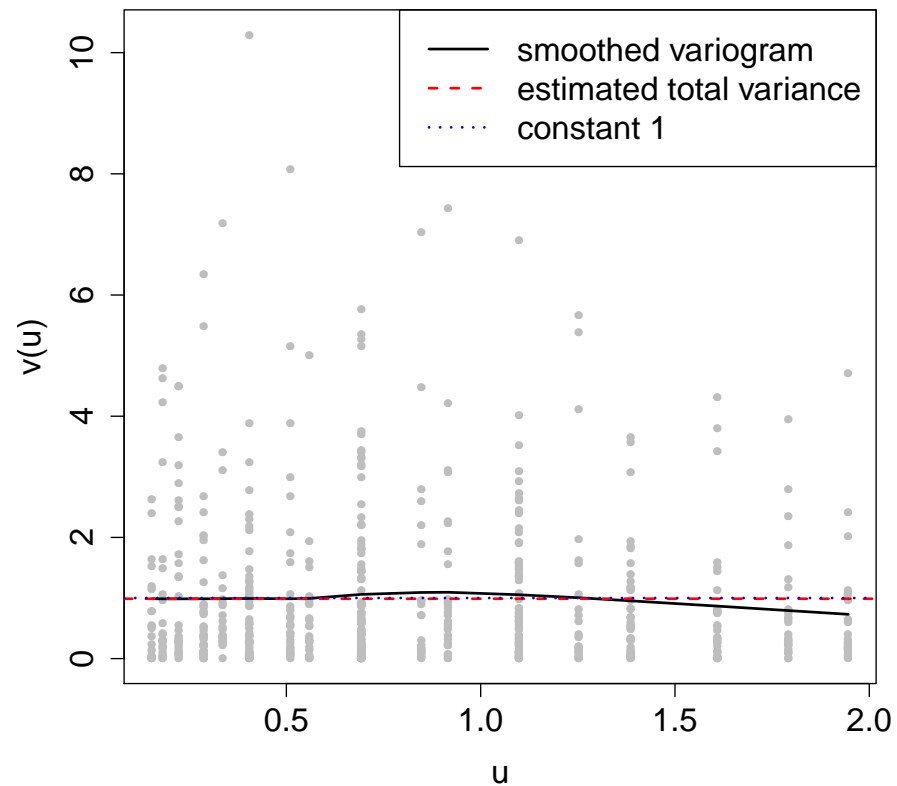
As the transformed residuals are uncorrelated with mean zero and variance one, we have

$$\begin{aligned}\frac{1}{2}\mathbf{E}[(r_{ij}^* - r_{ik}^*)^2] &= \frac{1}{2} [\mathbf{Var}(r_{ij}^*) + \mathbf{Var}(r_{ik}^*) - 2\mathbf{Cov}(r_{ij}^*, r_{ik}^*)] \\ &= \frac{1}{2} \cdot 1 + \frac{1}{2} \cdot 1 - 0 = 1.\end{aligned}$$

Thus, if the model for the covariance structure is correct, the empirical semi-variogram for the transformed residuals should randomly fluctuate around the constant 1.

Example rat data

Semi-variogram for the transformed residuals, random intercept model:



The normality assumption for the random effects

It would be of interest to look at the distribution of the \mathbf{b}_i a) to check the normality assumption and b) to find outlying individuals. However, the $\hat{\mathbf{b}}_i$

- all have different distributions unless all \mathbf{X}_i and \mathbf{Z}_i are equal.
- can look normal even if the true distribution of \mathbf{b}_i is not normal (e.g. bimodal). This is due to the shrinkage effect.

Fitting a model with a mixture distribution for the random effects (see Section 6.3) allows to check for normality of the random effects. To find outlying individuals, we can use the Mahalanobis distance, see slide 14.

Overview Chapter 7 - Model building and model choice

7.1 Model diagnostics

7.2 Model choice

Model choice

Often, there are several possible model specifications. To compare two models M_1 and M_2 , one can

- directly compare the likelihood if the numbers of parameters in M_1 and M_2 are the same.

Examples:

- Gaussian vs. exponential serial correlation
- different transformations of a covariate in the fixed effects
- conduct a test if M_1 and M_2 are nested, see Chapter 5.
- use information criteria for model selection.

Information criteria

- **Goal:** Comparison of models M_1 and M_2 with potentially different numbers of parameters (potentially non-nested).
- Denote by l_1 and l_2 the maximized log-likelihood for models M_1 and M_2 .
- If M_1 is nested in M_2 , we can conduct a likelihood ratio test ($H_0 : \boldsymbol{\theta} \in \Theta_1$ vs. $H_A : \boldsymbol{\theta} \in \Theta_2$, with $\Theta_1 \subset \Theta_2$). H_0 is rejected if

$$-2(l_1 - l_2) > \chi_{df_2 - df_1; 1 - \alpha}^2,$$

where df_2 and df_1 are the number of parameters for models M_2 and M_1 , respectively, and $\chi_{d; 1 - \alpha}^2$ is the $(1 - \alpha)$ -Quantile of the χ_d^2 distribution.

Information criteria

$$-2(l_1 - l_2) > \chi_{df_2 - df_1; 1 - \alpha}^2$$

can be alternatively written as

$$-2l_1 + \mathcal{F}(df_1) > -2l_2 + \mathcal{F}(df_2)$$

for a function $\mathcal{F}(df)$ with $\mathcal{F}(df_2) - \mathcal{F}(df_1) = \chi_{df_2 - df_1; 1 - \alpha}^2$.

An LR-test can only be conducted if the two models are nested, i.e. $\Theta_1 \subset \Theta_2$. If $\Theta_1 \not\subset \Theta_2$, information criteria of the form $-2l + \mathcal{F}(df)$ can be used to compare the models.

Information criteria

The two most commonly used criteria $-2l + \mathcal{F}(df)$ are

Criterion	$\mathcal{F}(df)$
Akaike (AIC)	$2df$
Schwarz (BIC)	$\ln(n)df$

where n denotes the sample size and $df = \dim(\Theta)$ the number of parameters. The AIC chooses more complex models than the BIC.

For the linear mixed model, the question is: which are the correct log-likelihood and number of parameters to use?

The Akaike information criterion (AIC) - Background

- Suppose data \mathbf{y} is generated from a true underlying model with density $g(\cdot)$. We approximate $g(\cdot)$ by a parametric class of models $f_{\boldsymbol{\theta}}(\cdot) = f(\cdot|\boldsymbol{\theta})$.
- Measure the discrepancy by the Kullback-Leibler (KL) distance

$$\begin{aligned} D(g, f_{\boldsymbol{\theta}}) &= \mathbb{E}_{g(\mathbf{z})} \log \left[\frac{g(\mathbf{z})}{f_{\boldsymbol{\theta}}(\mathbf{z})} \right] \\ &= \int \log \left[\frac{g(\mathbf{z})}{f_{\boldsymbol{\theta}}(\mathbf{z})} \right] g(\mathbf{z}) d\mathbf{z} \\ &= \mathbb{E}_{g(\mathbf{z})} \{ \log[g(\mathbf{z})] - \log[f_{\boldsymbol{\theta}}(\mathbf{z})] \} \geq 0. \end{aligned}$$

- Out of a sequence of models, choose the one that minimizes $D(g, f_{\theta})$.
In practice: estimate θ by $\hat{\theta}(\mathbf{y})$, minimize the **expected** KL distance

$$E_{g(\mathbf{y})}(D(g, f_{\hat{\theta}(\mathbf{y})})) = E_{g(\mathbf{y})}(E_{g(\mathbf{z})}\{\log[g(\mathbf{z})] - \log[f_{\hat{\theta}(\mathbf{y})}(\mathbf{z})]\})$$

Equivalently, minimize the so-called **Akaike information** (AI)

$$\Leftrightarrow -2E_{g(\mathbf{y})}(E_{g(\mathbf{z})}\{\log[f_{\hat{\theta}(\mathbf{y})}(\mathbf{z})]\}) \quad (7.1)$$

-

$$AIC = -2\log(f_{\hat{\theta}(\mathbf{y})}(\mathbf{y})) + 2df$$

is an asymptotically unbiased estimator of the AI (7.1) under certain regularity conditions.

- Minimizing the AIC over a set of possible models can thus be viewed as minimizing the average distance of an approximating model to the underlying truth.
- Note: (7.1) is a **predictive** quantity, with independent replications \mathbf{y} and \mathbf{z} . The maximized log-likelihood $\log(f_{\hat{\boldsymbol{\theta}}(\mathbf{y})}(\mathbf{y}))$ uses the **same data** \mathbf{y} for both model fitting and evaluation.
 - \Rightarrow The blue bias correction term $+ 2df$ is necessary to avoid overoptimism in the model fit and resulting overfitting.
- The assumed regularity conditions include that the parameter space Θ is open in \mathbb{R}^{df} and thus do not hold if $\boldsymbol{\theta}$ includes variance parameters.

The marginal AIC

The first option is to base the AIC in the linear mixed model on the marginal log-likelihood for the marginal model (3.5),

$$\log f(\mathbf{y}|\boldsymbol{\beta}, \boldsymbol{\alpha}) = \ell_{ML}(\boldsymbol{\theta}) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^N \log |\mathbf{V}_i(\boldsymbol{\alpha})| - \frac{1}{2} \left\{ \sum_{i=1}^N (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})^T \mathbf{V}_i(\boldsymbol{\alpha})^{-1} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}) \right\}.$$

Statistical software (e.g. lme) often returns a marginal AIC using $\ell_{ML}(\hat{\boldsymbol{\theta}}_{ML})$ and with df set to the total number of parameters in $\boldsymbol{\theta}$.

The marginal AIC

- The marginal AIC as predictive quantity assumes that two independent replications z and y come from the same marginal distribution, but **do not share the same random effects**. It is thus appropriate when the focus is on the population-level **fixed effects**.
- The parameter space Θ for θ is not open (e.g. $d_{kk} \geq 0$). As for the LRT, the usual χ^2 asymptotic distributions become χ^2 mixtures and the bias is no longer $2df = 2E(\chi_{df}^2)$, as assumed in the marginal AIC. This induces a preference for models without random effects ([Greven & Kneib, 2010](#)). The selection of fixed effects is likely not or not much affected.

The marginal AIC

For REML estimation, an AIC based on $\ell_{REML}(\hat{\alpha}_{REML})$ is often returned by statistical software (e.g. lme). The **marginal AIC should not be used with REML estimation** as

- a) the REML-likelihoods for different fixed effects are not comparable
- b) the fixed effects do not even occur in the REML-likelihood
- c) additionally, the used degrees of freedom often incorrectly still include the number of fixed effects.

The conditional AIC

An alternative is to base the AIC on the conditional log-likelihood

$$\begin{aligned} \log f(\mathbf{y}|\mathbf{b}, \boldsymbol{\beta}, \boldsymbol{\alpha}) &= -\frac{n}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^N \log |\boldsymbol{\Sigma}_i(\boldsymbol{\alpha})| \\ &\quad - \frac{1}{2} \left\{ \sum_{i=1}^N (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta} - \mathbf{Z}_i\mathbf{b}_i)^T \boldsymbol{\Sigma}_i(\boldsymbol{\alpha})^{-1} (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta} - \mathbf{Z}_i\mathbf{b}_i) \right\}. \end{aligned}$$

The conditional AIC uses $\log f(\mathbf{y}|\hat{\mathbf{b}}, \hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\alpha}})$ where the predicted or estimated quantities can be based on ML or REML estimation. The conditional log-likelihood is always based on \mathbf{Y} and valid with ML or REML estimation.

The conditional AIC

- The conditional AIC as a predictive quantity assumes that two independent replications z and y come from the same conditional distribution and **share the same random effects**. [Vaida & Blanchard \(2005\)](#) argue that it is appropriate when the focus is on the subjects or **random effects**.
- [Greven & Kneib \(2010\)](#) propose an unbiased estimator for the degrees of freedom to use in the conditional AIC (for $\mathbf{R} = \sigma^2 \mathbf{I}_n$), which is implemented in the R-package `cAIC4` for models fitted with `lme4` or `gamm4`. The random effects, due to shrinkage, contribute between 0 and Nq degrees of freedom.

Example rat data

Consider again the random intercept model for the rat data

$$Y_{ij} = \beta_0 + b_{1i} + \beta_{g_i} t_j + \epsilon_{ij}$$

with transformed time t_j and compare with the untransformed time $TIME_j$.

```
> lmet <- lme(RESPONSE ~ group * t - group,
             random = ~ 1 | SUBJECT, data = rats, method = "ML")
> lmeTIME <- lme(RESPONSE ~ group * TIME - group,
                random = ~ 1 | SUBJECT, data = rats, method = "ML")
> anova(lmet, lmeTIME)
```

	Model	df	AIC	BIC	logLik
lmet	1	6	931.9924	953.169	-459.9962
lmeTIME	2	6	1074.0125	1095.189	-531.0063

Interpretation?

Example sleep deprivation study

For the sleep deprivation data, compare a model with a random intercept with a model with random intercept and slope.

```
> library(lme4)
> library(cAIC4)
> M1 <- lmer(Reaction ~ Days + (1 | Subject), sleepstudy)
> M2 <- lmer(Reaction ~ Days + (1 + Days | Subject), sleepstudy)
> cAIC(M1)$caic
[1] 1767.118
> cAIC(M2)$caic
[1] 1711.618
```

Interpretation?