

# The Linear Mixed-Effects Probability Model

So far:

- Micro Introduction
- Examples of different models and datasets

In this talk:

- Formalize notation
- Introduce new definitions
- Model fitting (as done in the *lme4*-package)

## Recall a basic example of a model

(workers and machines)

### Model:

$$y_{ijk} = \mu + \beta_j + b_i + \epsilon_{ijk}, \quad i = 1, \dots, 6 \quad j = 1, 2, 3 \quad k = 1, 2, 3$$

$\beta_j$  : effect of machine (fixed)

$b_i$ : effect of worker (random)

### Assumption:

$$b_i \sim \mathcal{N}(0, \sigma_b^2) \quad \epsilon_{ijk} \sim \mathcal{N}(0, \sigma^2)$$

## Generalization:

Two random variables:

$\mathcal{Y}$ : the  $n$ -dimensional response vector

$\mathcal{B}$ : the  $q$ -dimensional vector of random effects

$$\mathcal{B} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_\theta). \quad (1)$$

$$(\mathcal{Y}|\mathcal{B} = \mathbf{b}) \sim \mathcal{N}(\mathbf{Z}\mathbf{b} + \mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n) \quad (2)$$

So the *linear predictor* is

$$\mathbf{Z}\mathbf{b} + \mathbf{X}\boldsymbol{\beta} \quad (3)$$

With the *model matrices*  $\mathbf{Z}$  of dimension  $n \times q$

and  $\mathbf{X}$  of dimension  $n \times p$ ,

where  $p$  is the dimension of the *fixed-effects* parameter vector  $\boldsymbol{\beta}$

## More definitions:

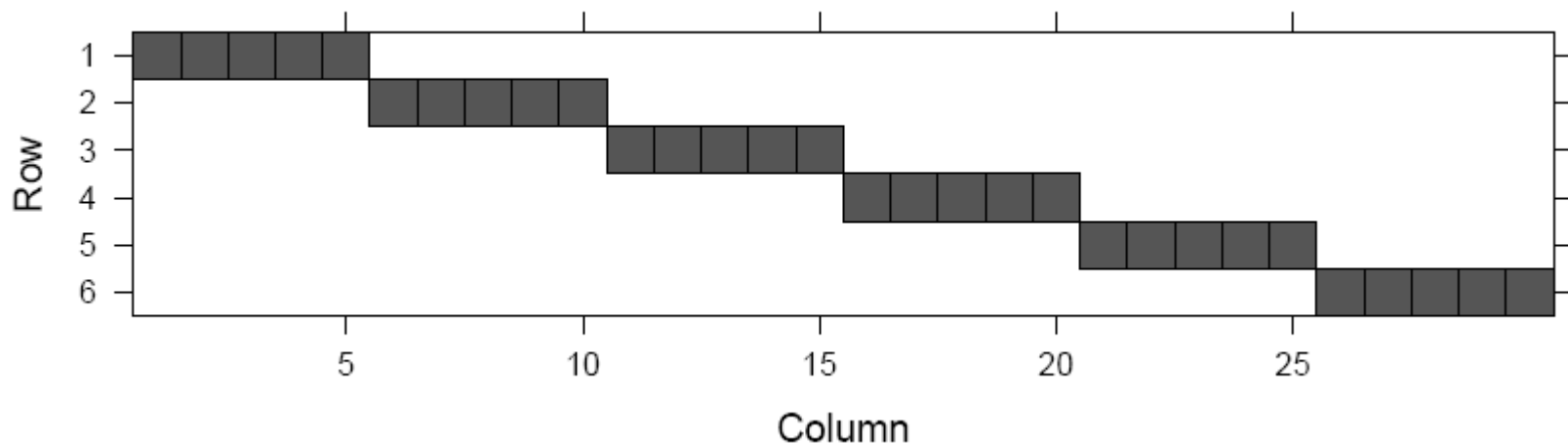
$\theta$  : *variance-component parameter vector*

$\Sigma_{\theta}$  : *variance-covariance matrix*

$\sigma$  : *common scale parameter*

The form of the random-effects model matrix,  $\mathbf{Z}$ , and the form of the variance-covariance matrix,  $\Sigma_{\theta}$ , and the method by which  $\Sigma_{\theta}$  is determined from the value of  $\theta$  are all based on the random-effects terms in the model formula.

$\mathbf{Z}$  can be large, but it is sparse (i.e. most elements in the matrix are zero).



## More definitions:

$\Lambda_\theta$ : relative covariance factor, defined so that

$$\Sigma_\theta = \sigma^2 \Lambda_\theta \Lambda_\theta^T$$

where  $\sigma^2$  is the same variance parameter as in  $(Y|B = \mathbf{b})$ .

With the spherical random effects:

$$\mathcal{U} \sim \mathcal{N}(\mathbf{0}, \sigma^2 I_q)$$

we get:

$$\mathcal{B} = \Lambda_\theta \mathcal{U}$$

So we really have:

$$E[\mathcal{B}] = \Lambda_\theta E[\mathcal{U}] = \Lambda_\theta \mathbf{0} = \mathbf{0}$$

and:

$$\begin{aligned} \text{Var}(\mathcal{B}) &= E[(\mathcal{B} - E[\mathcal{B}])(\mathcal{B} - E[\mathcal{B}])^T] = E[\mathcal{B}\mathcal{B}^T] \\ &= E[\mathbf{\Lambda}_\theta \mathcal{U}\mathcal{U}^T \mathbf{\Lambda}_\theta^T] = \mathbf{\Lambda}_\theta E[\mathcal{U}\mathcal{U}^T] \mathbf{\Lambda}_\theta^T = \mathbf{\Lambda}_\theta \text{Var}(\mathcal{U}) \mathbf{\Lambda}_\theta^T \\ &= \mathbf{\Lambda}_\theta \sigma^2 I_q \mathbf{\Lambda}_\theta^T = \sigma^2 \mathbf{\Lambda}_\theta \mathbf{\Lambda}_\theta^T = \mathbf{\Sigma}_\theta \end{aligned}$$

In our discussion we will concentrate on  $\mathbf{\Lambda}_\theta$  (not  $\mathbf{\Sigma}_\theta$ ) and  $\mathcal{U}$  (not  $\mathcal{B}$ ):

So we look at:

$$\mathcal{U} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_q). \quad (4)$$

$$(\mathcal{Y} | \mathcal{U} = \mathbf{u}) \sim \mathcal{N}(\mathbf{Z}\mathbf{\Lambda}_\theta \mathbf{u} + \mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n) \quad (5)$$



So the *linear predictor* becomes:

$$\boldsymbol{\gamma} = \mathbf{Z}\boldsymbol{\Lambda}_\theta\boldsymbol{u} + \mathbf{X}\boldsymbol{\beta} \quad (6)$$

And the *conditional mean* of  $\mathcal{Y}$ , given  $\mathcal{U} = \boldsymbol{u}$ :

$$\boldsymbol{\mu} = E[\mathcal{Y}|\mathcal{U} = \boldsymbol{u}] \quad (7)$$

Note: For a linear mixed model, we obviously have  $\boldsymbol{\mu} = \boldsymbol{\gamma}$ .

In other forms of mixed models this may not be the case anymore.

## Conditional Distribution:

### Notation:

$y_{obs}$  : an observed data vector (an actual realization of  $\mathcal{Y}$ )

$y$  : an arbitrary value of  $\mathcal{Y}$

Now we are interested in the conditional distribution of  $(\mathcal{U}|\mathcal{Y} = \mathbf{y})$

Our model parameters are  $\boldsymbol{\theta}$ ,  $\boldsymbol{\beta}$  and  $\boldsymbol{\sigma}$ .

The likelihood of those parameters, given the observed data,  $\mathbf{y}_{obs}$ , is the probability density of  $\mathcal{Y}$ , evaluated at  $\mathbf{y}_{obs}$ .

Parameters fixed,  $y$  varying

$y$  fixed at  $\mathbf{y}_{obs}$ , parameters varying

## Natural approach for evaluating the likelihood:

1. Determine marginal distribution of  $\mathcal{Y}$ 
  - Determine joint density of  $\mathcal{U}$  and  $\mathcal{Y}$ :  $f_{\mathcal{Y},\mathcal{U}}(\mathbf{y}, \mathbf{u})$
  - Integrate this density wrt.  $\mathbf{u}$  to get the marginal density:  $f_{\mathcal{Y}}(\mathbf{y})$
2. Evaluate that density at  $\mathbf{y}_{obs}$ .

But here we choose a different order of the steps, that is:

Evaluate the joint density at  $\mathbf{y}_{obs}$  to produce an intermediate function  $h(\mathbf{u})$ .

And then integrate this function  $h(\mathbf{u})$  along  $\mathbf{u}$ .

This does not work generally, it could even happen that the joint density does not exist (think of a joint distribution that is discrete wrt.  $\mathbf{y}$  and continuous wrt.  $\mathbf{u}$ )

We define:

$$h(\mathbf{u}) = f_{\mathcal{Y}, \mathcal{U}}(\mathbf{y}_{obs}, \mathbf{u}) \quad (9)$$

the *unnormalized conditional density*.

We see that:

$$f_{\mathcal{U}|\mathcal{Y}}(\mathbf{u}|\mathbf{y}_{obs}) = \frac{h(\mathbf{u})}{\int_{R^q} h(\mathbf{u}) d\mathbf{u}} \quad (10)$$

and thus the likelihood is:

$$L(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma | \mathbf{y}_{obs}) = \int_{R^q} f_{\mathcal{Y}, \mathcal{U}}(\mathbf{y}_{obs}, \mathbf{u}) d\mathbf{u} = \int_{R^q} h(\mathbf{u}) d\mathbf{u}. \quad (11)$$

Tools to evaluate the likelihood in general situations:

$\tilde{\mathbf{u}}$  : the conditional mode of  $\mathbf{u}$ , given  $\mathcal{Y} = \mathbf{y}_{obs}$ :

$$\tilde{\mathbf{u}} = \arg \max_{\mathbf{u}} f_{\mathcal{U}|\mathcal{Y}}(\mathbf{u}|\mathbf{y}_{obs}) = \arg \max_{\mathbf{u}} h(\mathbf{u}) = \arg \max_{\mathbf{u}} f_{\mathcal{Y}|\mathcal{U}}(\mathbf{y}_{obs}|\mathbf{u}) f_{\mathcal{U}}(\mathbf{u}) \quad (12)$$

Recall (4) and (5):

$$\mathcal{U} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_q). \quad (4)$$

$$(\mathcal{Y}|\mathcal{U} = \mathbf{u}) \sim \mathcal{N}(\mathbf{Z}\Lambda_\theta \mathbf{u} + \mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n) \quad (5)$$

Thus we have:

$$f_{\mathcal{Y}|\mathcal{U}}(\mathbf{y}|\mathbf{u}) = \frac{\exp(-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\Lambda_\theta \mathbf{u}\|^2)}{(2\pi\sigma^2)^{n/2}} \quad (13)$$

$$f_{\mathcal{U}}(\mathbf{u}) = \frac{\exp(-\frac{1}{2\sigma^2} \|\mathbf{u}\|^2)}{(2\pi\sigma^2)^{q/2}} \quad (14)$$

And so the product:

$$h(\mathbf{u}) = \frac{\exp(- [\|\mathbf{y}_{obs} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\boldsymbol{\Lambda}_\theta \mathbf{u}\|^2 + \|\mathbf{u}\|^2] / (2\sigma^2))}{(2\pi\sigma^2)^{(n+q)/2}} \quad (15)$$

Looking at the negative log density, we get:

$$-2 \log(h(\mathbf{u})) = (n + q) \log(2\pi\sigma^2) + \frac{\|\mathbf{y}_{obs} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\boldsymbol{\Lambda}_\theta \mathbf{u}\|^2 + \|\mathbf{u}\|^2}{\sigma^2} \quad (16)$$

So we get:

$$\tilde{\mathbf{u}} = \arg \min_{\mathbf{u}} \|\mathbf{y}_{obs} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\boldsymbol{\Lambda}_\theta \mathbf{u}\|^2 + \|\mathbf{u}\|^2 \quad (17)$$

↑  
sum of squared residuals

↑  
penalty for high complexity

The expression to be minimized  $\|\mathbf{y}_{obs} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\boldsymbol{\Lambda}_\theta \mathbf{u}\|^2 + \|\mathbf{u}\|^2$  is called the *objective function*, here it is a *penalized residual sum of squares* (**PRSS**).

The minimizer  $\tilde{\mathbf{u}}$  is called the *penalized least squares* (**PLS**) solution



We think of the **PRSS** criterion as a function of the parameters, given the data, ie.:

$$\mathbf{r}_{\theta, \beta}^2 = \min_{\mathbf{u}} [\|\mathbf{y}_{obs} - \mathbf{X}\beta - \mathbf{Z}\Lambda_{\theta} \mathbf{u}\|^2 + \|\mathbf{u}\|^2] \quad (18)$$

We can also minimize this expression wrt  $\beta$ .  
And we will see that this can even be done simultaneously wrt  $\mathbf{u}$  and  $\beta$  without using iterations. The minimum value we get is:

$$\mathbf{r}_{\theta}^2 = \min_{\mathbf{u}, \beta} [\|\mathbf{y}_{obs} - \mathbf{X}\beta - \mathbf{Z}\Lambda_{\theta} \mathbf{u}\|^2 + \|\mathbf{u}\|^2] \quad (19)$$

$\tilde{\beta}$ : *conditional estimate* of  $\beta$

the value of  $\beta$  for which the minimum in (19) is attained.

One way to determine the solution is to rephrase it as a linear least squares problem for an extended residual vector

$$\tilde{\mathbf{u}} = \arg \min_{\mathbf{u}} \left\| \begin{bmatrix} \mathbf{y}_{obs} - \mathbf{X}\boldsymbol{\beta} \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{Z}\Lambda_{\theta} \\ \mathbf{I}_q \end{bmatrix} \mathbf{u} \right\|^2 \quad (20)$$

This is called a pseudo-data approach because we create the effect of the penalty term,  $\|\mathbf{u}\|^2$ , by adding “pseudo-observations” to the response vector and to the predictor (adding zeros and  $\mathbf{I}_q\mathbf{u}$ ).

For this linear least squares problem, we can give the solution by solving the normal equations. So we get that the solution satisfies:

$$(\Lambda_{\theta}^T \mathbf{Z}^T \mathbf{Z} \Lambda_{\theta} + \mathbf{I}_q) \tilde{\mathbf{u}} = \Lambda_{\theta}^T \mathbf{Z}^T (\mathbf{y}_{obs} - \mathbf{X}\boldsymbol{\beta}) \quad (21)$$

We want fast evaluation of  $\tilde{\mathbf{u}}$  for different inputs, so we form the *sparse Cholesky factor*,  $\mathbf{L}_\theta$ .

It is a lower  $q \times q$  matrix with:

$$\mathbf{L}_\theta \mathbf{L}_\theta^T = (\Lambda_\theta^T \mathbf{Z}^T \mathbf{Z} \Lambda_\theta + \mathbf{I}_q) \quad (22)$$

In order to get a sparse Cholesky factor  $\mathbf{L}_\theta$  we might want to permute the columns of our data.

This is done through a so-called *Permutation matrix*  $\mathbf{P}$ .

We also call them *fill-reducing permutations* as we want to avoid positions in the factor getting filled, where the matrix being decomposed is zero.

(22) thus becomes:

$$\mathbf{L}_\theta \mathbf{L}_\theta^T = \mathbf{P}(\Lambda_\theta^T \mathbf{Z}^T \mathbf{Z} \Lambda_\theta + \mathbf{I}_q) \mathbf{P}^T \quad (23)$$

The pseudo-data representation in (20) becomes:

$$\tilde{\mathbf{u}} = \arg \min_{\mathbf{u}} \left\| \begin{bmatrix} \mathbf{y}_{obs} - \mathbf{X}\boldsymbol{\beta} \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{Z}\Lambda_{\theta}\mathbf{P}^T \\ \mathbf{P}^T \end{bmatrix} \mathbf{P}\mathbf{u} \right\|^2 \quad (24)$$

And the system of linear equations for  $\tilde{\mathbf{u}}$  accordingly:

$$\mathbf{L}_{\theta}\mathbf{L}_{\theta}^T\mathbf{P}\tilde{\mathbf{u}} = \mathbf{P}(\Lambda_{\theta}^T\mathbf{Z}^T(\mathbf{y}_{obs} - \mathbf{X}\boldsymbol{\beta}))\mathbf{P}^T\mathbf{P}\tilde{\mathbf{u}} = \mathbf{P}\Lambda_{\theta}^T\mathbf{Z}^T(\mathbf{y}_{obs} - \mathbf{X}\boldsymbol{\beta}) \quad (25)$$

Note: Once we evaluate  $L_{\theta}$  it is straight forward to solve (25) for  $\tilde{\mathbf{u}}$ . Thus this step is very crucial, and the ability to evaluate  $L_{\theta}$  rapidly for many different values of  $\boldsymbol{\theta}$  is what makes the methods in `lme4` feasible.

## Back to the evaluation of the likelihood:

We've seen in (11) and (15) that:

$$L(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma | \mathbf{y}_{obs}) = \int_{R^q} f_{\mathcal{Y}, \mathcal{U}}(\mathbf{y}_{obs}, \mathbf{u}) d\mathbf{u} = \int_{R^q} h(\mathbf{u}) d\mathbf{u}.$$

$$h(\mathbf{u}) = \frac{\exp(- [\|\mathbf{y}_{obs} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\boldsymbol{\Lambda}_\theta \mathbf{u}\|^2 + \|\mathbf{u}\|^2] / (2\sigma^2))}{(2\pi\sigma^2)^{(n+q)/2}}$$

We can now write the **PRSS** for general  $\mathbf{u}$  as:

$$\|\mathbf{y}_{obs} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\boldsymbol{\Lambda}_\theta \mathbf{u}\|^2 + \|\mathbf{u}\|^2 = r_{\theta, \beta}^2 + \|\mathbf{L}_\theta^T(\mathbf{u} - \tilde{\mathbf{u}})\|^2 \quad (26)$$

Plugging this into the definition of  $h(\mathbf{u})$  and using the change-of-variable:

$$z = \frac{\mathbf{L}_\theta^T(\mathbf{u} - \tilde{\mathbf{u}})}{\sigma} \quad (27)$$

We get after a calculation (Bates[10], ch. 5.4.2 - available on <http://lme4.r-forge.r-project.org/book/>):

$$L(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma | \mathbf{y}_{obs}) = \frac{\exp(-\frac{r_{\boldsymbol{\theta}, \boldsymbol{\beta}}^2}{2\sigma^2})}{(2\pi\sigma^2)^{n/2} |\mathbf{L}_{\boldsymbol{\theta}}|}$$

So the deviance (negative twice the log-likelihood) becomes:

$$d(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma | \mathbf{y}_{obs}) = -2 \log(L(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma | \mathbf{y}_{obs})) = n \log(2\pi\sigma^2) + \frac{r_{\boldsymbol{\theta}, \boldsymbol{\beta}}^2}{\sigma^2} + 2 \log(|\mathbf{L}_{\boldsymbol{\theta}}|^2)$$

So the maximum-likelihood estimates for the parameters are those that minimize this deviance.

We can even further simplify this expression by using the facts, that  $\beta$  only occurs in  $r_{\theta, \beta}^2$  and minimizing this expression wrt  $\beta$  for any value of  $\theta$  goes back to the penalized least square problems.

So let  $\hat{\beta}_\theta$  be the value of  $\beta$  that minimizes PRSS wrt to  $\beta$  and  $u$ .

And  $r_\theta^2$  the PRSS at these minimizing values.

Furthermore let  $\hat{\sigma}_\theta^2 = r_\theta^2/n$ , the value of  $\sigma^2$  that minimizes the above deviance or a given  $r_\theta^2$ .

Then the *profiled deviance*, which is now only a function of  $\theta$ , becomes:

$$\tilde{d}(\theta|\mathbf{y}_{obs}) = 2 \log |\mathbf{L}_\theta| + n \left[ 1 + \log \left( \frac{2\pi r_\theta^2}{n} \right) \right]$$

Now minimization of  $\tilde{d}(\theta|\mathbf{y}_{obs})$  wrt  $\theta$  determines the MLE,  $\tilde{\theta}$ .  
The MLEs for  $\hat{\beta}$  and  $\hat{\sigma}$  then are the corresponding conditional estimates evaluated at  $\hat{\theta}$ .



Simultaneously evaluating  $\tilde{\mathbf{u}}$  and  $\beta_\theta$  uses the same approach we've already seen in (20), that is to rephrase the PLS problem into a linear least square problem.

Thus we rewrite:

$$\tilde{\mathbf{u}} = \arg \min_{\mathbf{u}} \left\| \begin{bmatrix} \mathbf{y}_{obs} - \mathbf{X}\beta \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{Z}\Lambda_\theta \\ \mathbf{I}_q \end{bmatrix} \mathbf{u} \right\|^2 \quad (20)$$

as

$$\begin{bmatrix} \tilde{\mathbf{u}} \\ \hat{\beta}_\theta \end{bmatrix} = \arg \min_{\mathbf{u}, \beta} \left\| \begin{bmatrix} \mathbf{y}_{obs} \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{Z}\Lambda_\theta \mathbf{P}^T & \mathbf{X} \\ \mathbf{P}^T & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{P}\mathbf{u} \\ \beta \end{bmatrix} \right\|^2 \quad (28)$$

Which now yields the equation:

$$\begin{bmatrix} P(\Lambda_\theta^T Z^T Z \Lambda_\theta + I_q) P^T & P \Lambda_\theta^T Z^T X \\ X^T Z \Lambda_\theta P^T & X^T X \end{bmatrix} \begin{bmatrix} P \tilde{u} \\ \hat{\beta}_\theta \end{bmatrix} = \begin{bmatrix} P \Lambda_\theta^T Z^T y_{obs} \\ X^T y_{obs} \end{bmatrix}$$

The Matrix on the LHS can be decomposed into a ‘*Cholesky-like*’ decomposition.

This way we also found a fast way to get  $\tilde{\mathbf{u}}$  and  $\hat{\beta}_\theta$ , thus we have found all MLE estimators of the parameters.

That is we fitted the model to the actual data.

Outlook: Complete Analysis of several data examples

